

# Aprendizado de Máquina

## Aprendizado não Supervisionado - Agrupamento

Inteligência Artificial

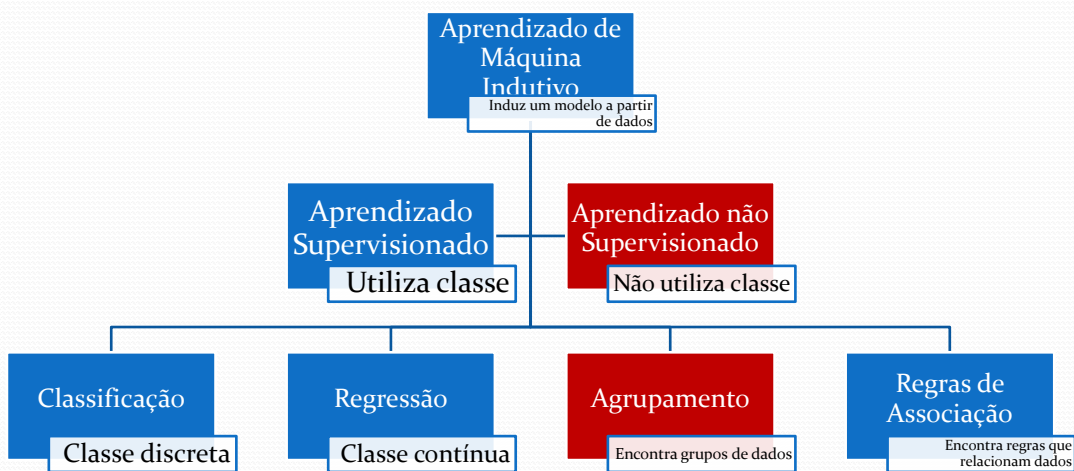
HAC

IA2022

1

1

## Hierarquia do Aprendizado de Máquina Indutivo



HAC

IA2022

2

2

## Agrupamento (clustering) - Discussão Geral

Tarefa de aprendizado não-supervisionado:

- Exemplos não estão rotulados – não existe um atributo especial conhecido como “atributo meta”

Cliente 1	renda	dívida	<del>classe</del>
xxx	50	10	<del>bom</del>

Não conhecida

exemplo	tempo	temperatura	umidade	vento	<del>classe</del>
E1	sol	2	72	forte	<del>sim</del>

HAC

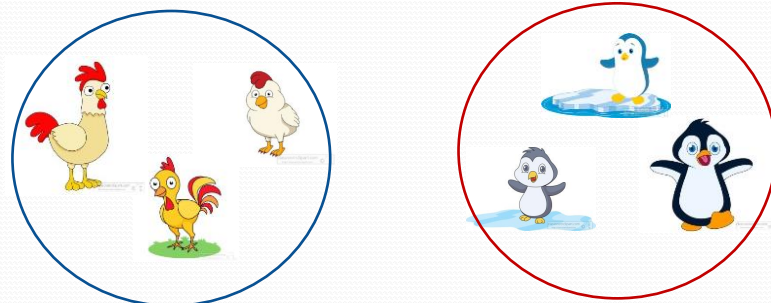
IA2022

3

3

## Agrupamento (clustering) - Discussão Geral

- Objetivo: agrupar objetos em **clusters** (grupos) de modo que objetos pertencentes a um mesmo **cluster** são **mais similares** entre si de acordo com alguma **medida de similaridade** pré-definida, enquanto que objetos pertencentes a clusters diferentes têm uma **similaridade menor**.



HAC

IA2022

4

4

## Agrupamento (clustering) - Aplicações



Encontrar grupos de documentos sobre um mesmo assunto

Descobrir funções de genes encontrando grupos de genes com características semelhantes



HAC

IA2022

5

5

## Agrupamento (clustering) - Aplicações



Encontrar grupos de usuários (visitantes) de um site e identificar suas características



Encontrar sub-populações de consumidores ou padrões de consumo por região

HAC

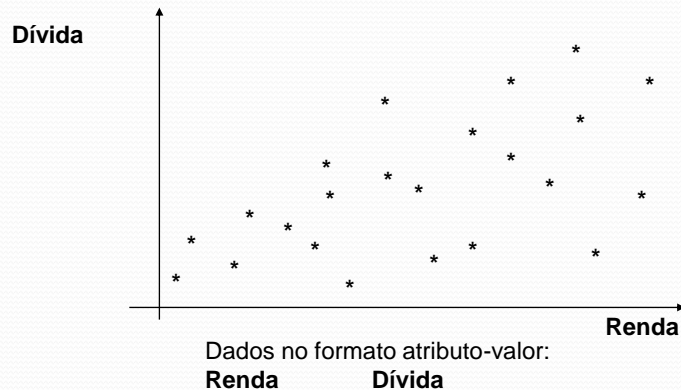
IA2022

6

6

## Agrupamento (clustering) - Discussão Geral

- Representação Gráfica – gráfico de dispersão



HAC

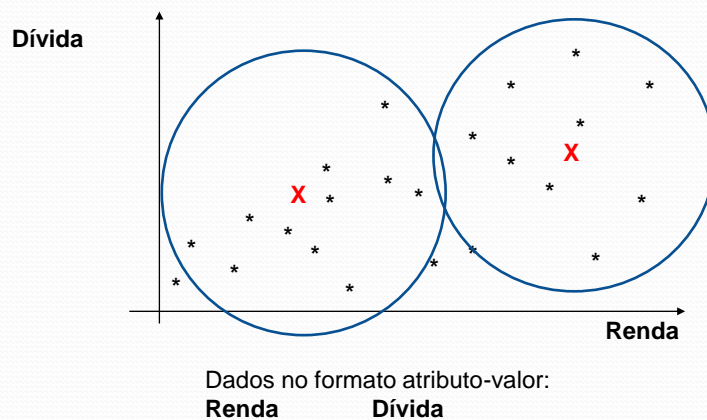
IA2022

7

7

## Agrupamento (clustering) - Discussão Geral

- Representação Gráfica – gráfico de dispersão



HAC

IA2022

8

8

## Agrupamento (clustering) - Discussão Geral

- Questões fundamentais para agrupamento:
  - O que é similaridade?
  - Como escolher uma medida de similaridade?
  - Qual o número ideal de grupos?
  - Como escolher um algoritmo?
  - Como validar e interpretar o resultado de um agrupamento?

HAC

IA2022

9

9

## Questões fundamentais para o agrupamento

- O que é similaridade?



HAC

IA2022

10

10

## Questões fundamentais para o agrupamento

- Como escolher uma medida de similaridade?
- **Medidas de proximidade:**
- São medidas de **similaridade** ou de **dissimilaridade** entre objetos
- Os algoritmos de agrupamento podem utilizar medidas de similaridade ou medidas de dissimilaridade
- Devem ser escolhidas de acordo com:
  - o tipo dos atributos envolvidos (contínuo, categórico)
  - a esparsidade dos dados

HAC

IA2022

11

11

## Questões fundamentais para o agrupamento

- Como escolher uma medida de proximidade?
  - São medidas de **dissimilaridade** ou **similaridade**

Entre objetos com atributos contínuos:

Medidas de distância  
(dissimilaridade)

Medidas de correlação  
(similaridade)  
Correlação de Pearson  
Medida de cosseno

Entre objetos com atributos discretos:

Coeficiente de casamento  
simples (CCS) (similaridade)

Coeficiente de  
Jaccard(similaridade)

(São muitas as medidas definidas na literatura)

HAC

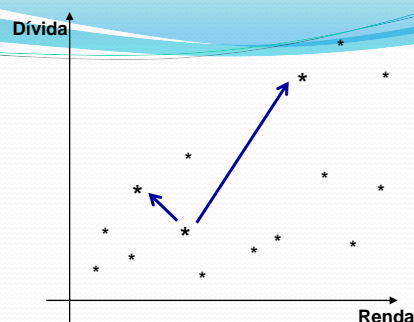
IA2022

12

12

## Medidas de dissimilaridade

- Medidas de distância



- atributos dos exemplos são considerados como dimensões de um espaço multidimensional
- cada exemplo corresponde a um ponto no espaço
- similaridade entre dois pontos é inversamente proporcional a distância entre eles

HAC

IA2022

13

13

## Medidas de dissimilaridade

- Medidas de distância

- Distância Euclidiana

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

	potência	peso	aceleração	consumo
E1	130	3504	12	18
E2	165	3693	11,5	15
E3	150	3436	11	18

$$d(E1, E2) = \sqrt{(130 - 165)^2 + (3504 - 3693)^2 + (12 - 11,5)^2 + (18 - 15)^2}$$

$$\text{Sim}(E1, E2) = 1 - d(E1, E2)$$

HAC

IA2022

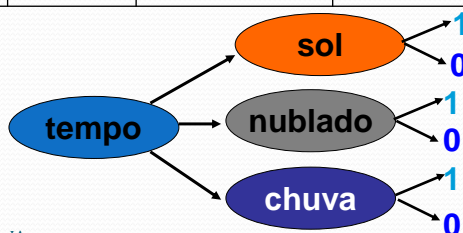
14

14

## Medidas de dissimilaridade

- Similaridade para variáveis nominais
  - Transformar cada valor do atributo nominal em uma variável binária fictícia.

	tempo	temperatura	umidade	vento
E1	sol	amena	alta	forte
E2	nublado	frio	média	forte
E3	sol	frio	alta	fraco



HAC

IA2022

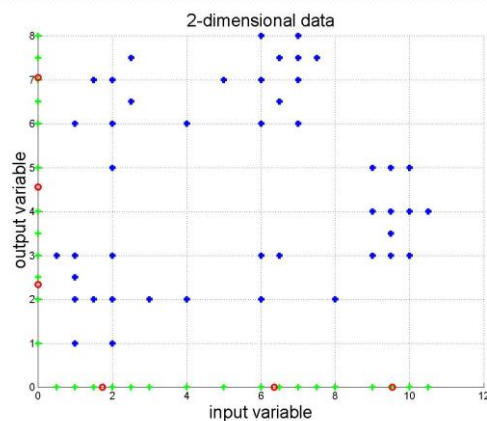
15

15

## Questões fundamentais para o agrupamento

- Qual o número ideal de grupos?
  - A noção de grupo não é determinística
  - Para um mesmo conjunto de dados, podemos encontrar números diferentes de grupos

**Em quantos grupos  
esses dados podem ser  
separados?**



HAC

IA2022

16

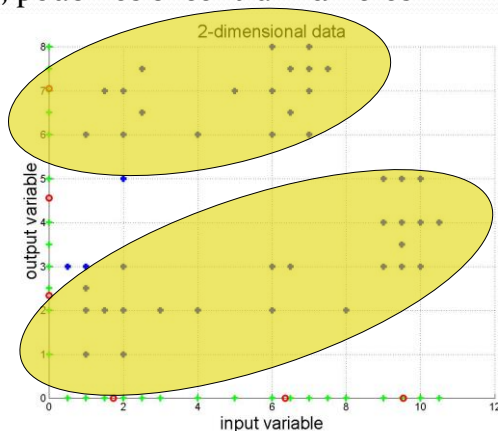
16



## Questões fundamentais para o agrupamento

- Qual o número ideal de grupos?
  - A noção de grupo não é determinística
  - Para um mesmo conjunto de dados, podemos encontrar números diferentes de grupos

**Em quantos grupos  
esses dados podem ser  
separados?**



HAC

IA2022

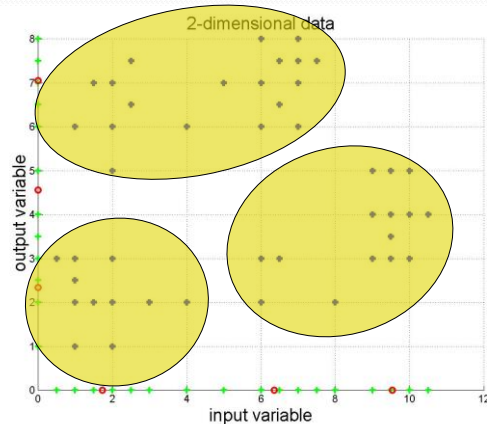
17

17

## Questões fundamentais para o agrupamento

- Qual o número ideal de grupos?
  - A noção de grupo não é determinística
  - Para um mesmo conjunto de dados, podemos encontrar números diferentes de grupos

**Em quantos grupos  
esses dados podem ser  
separados?**



HAC

IA2022

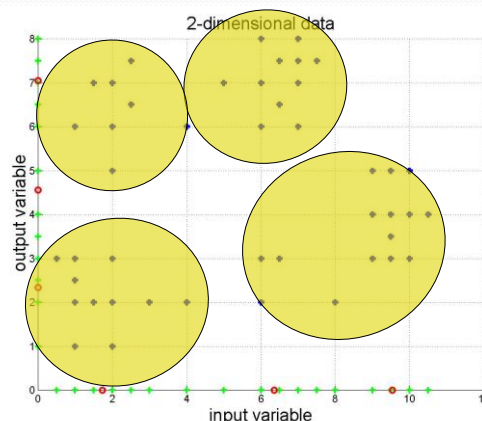
18

18

## Questões fundamentais para o agrupamento

- Qual o número ideal de grupos?
  - A noção de grupo não é determinística
  - Para um mesmo conjunto de dados, podemos encontrar números diferentes de grupos

**Em quantos grupos  
esses dados podem ser  
separados?**



HAC

IA2022

19

19

## Questões fundamentais para o agrupamento

- Qual o número ideal de grupos?
- O número de grupos encontrados, muitas vezes deve ser definido antes da aplicação do algoritmo de agrupamento (é um parâmetro do algoritmo).
- A forma do grupo está relacionada com a medida de similaridade ou dissimilaridade escolhida.
- O melhor número de grupos de um agrupamento em geral é encontrado por meio de experimentos

HAC

IA2022

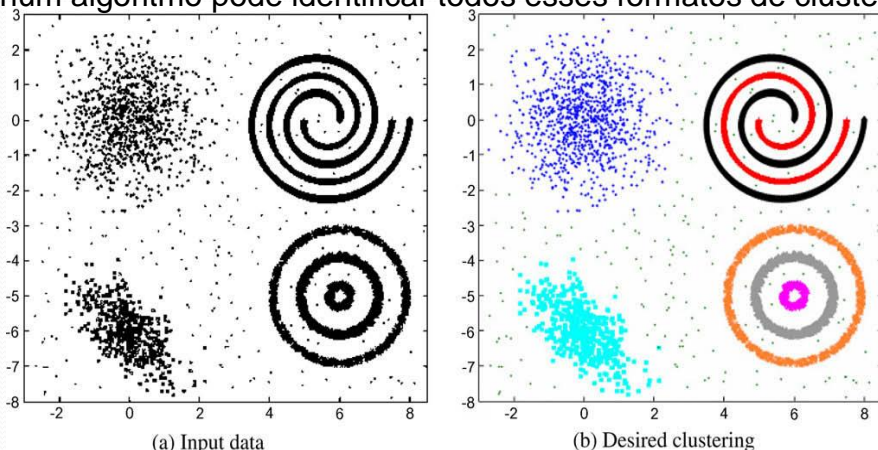
20

20

## Questões fundamentais para o agrupamento

- Como escolher um algoritmo?

Nenhum algoritmo pode identificar todos esses formatos de clusters



HAC

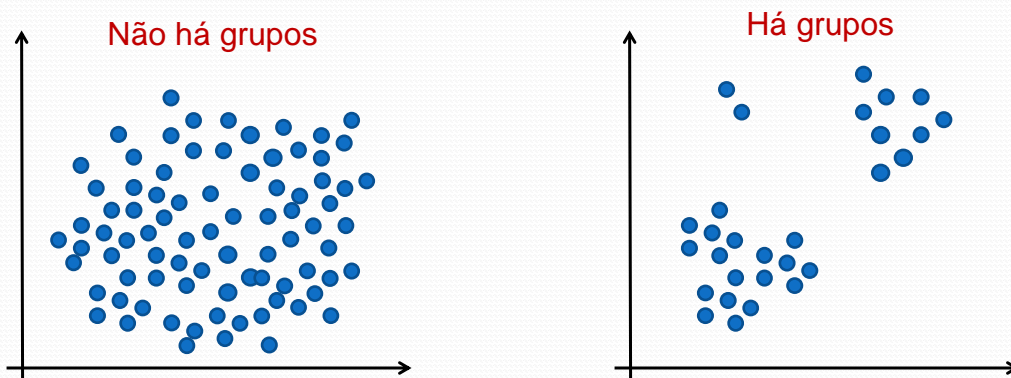
IA2022

21

21

## Questões fundamentais para o agrupamento

- Como validar e interpretar o resultado de um agrupamento?
- Todos os algoritmos de agrupamento, quando aplicados a dados, vão produzir grupos, independente dos dados possuírem grupos ou não.



HAC

IA2022

22

22

## Questões fundamentais para o agrupamento

- Como validar e interpretar o resultado de um agrupamento?
  - Uma estrutura é considerada **válida** se pode-se afirmar que não foi obtida por acaso ou por meio de artifícios do algoritmo;
  - Na **interpretação** dos grupos é feita a rotulação dos clusters, definindo sua natureza por meio da análise de seus objetos típicos;
  - Na interpretação, o papel do especialista é fundamental.

HAC

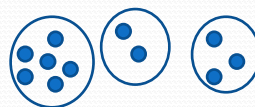
IA2022

23

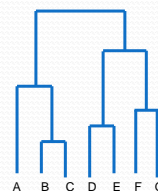
23

## Tipos de agrupamento

- **Agrupamento particional** – divide os dados em grupos sem sobreposição, de tal forma que cada dado pertence a apenas um grupo.



- **Agrupamento hierárquico** – divide os dados em grupos aninhados, que podem ser representados em uma estrutura de árvore chamada dendrograma



HAC

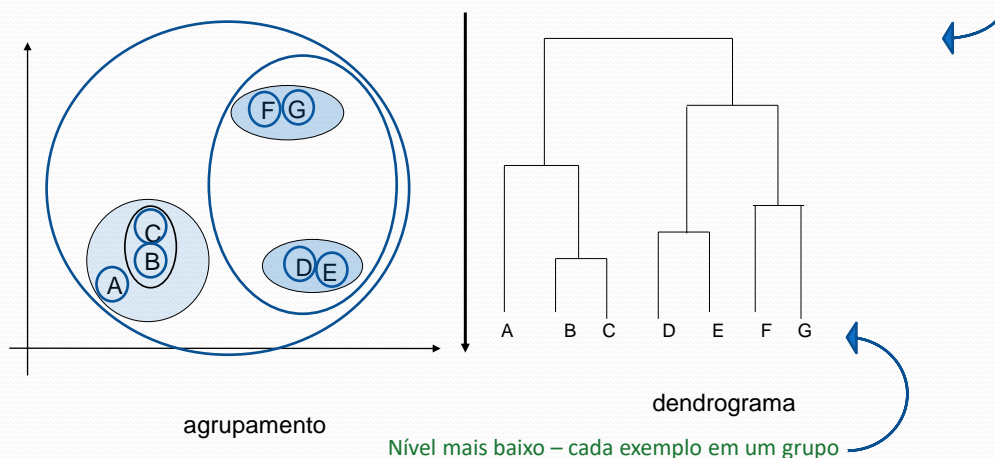
IA2022

24

24

## Agrupamento hierárquico

Nível mais alto – todos os exemplos em um só grupo



HAC

IA2022

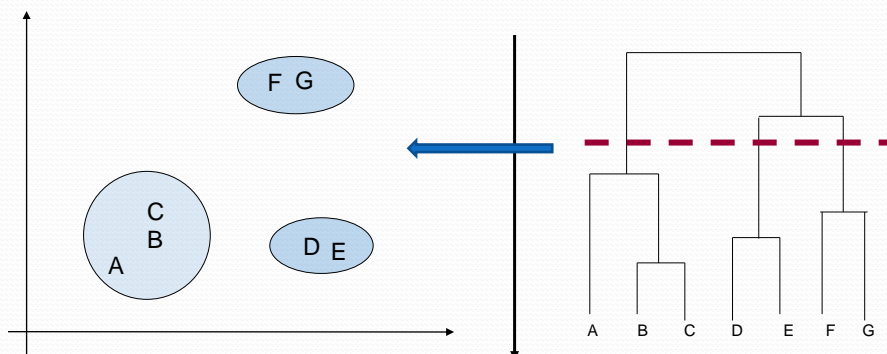
25

25

## Agrupamento hierárquico

- Seleção de um agrupamento

O corte do dendrograma em um determinado nível define um agrupamento em particular



HAC

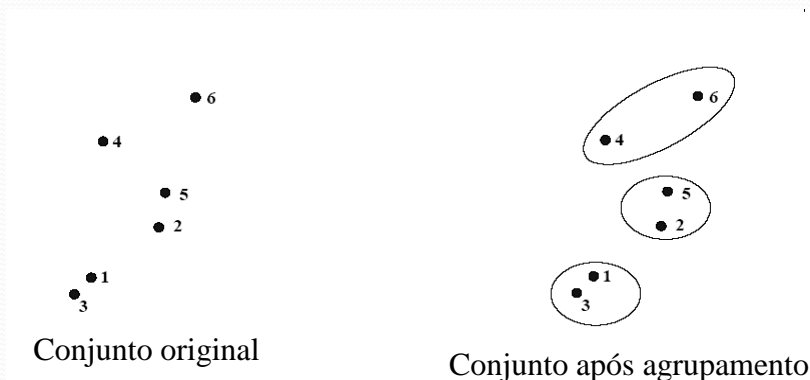
IA2022

26

26

## Agrupamento particional

- Dado um conjunto de dados finito  $X$  o problema de **agrupamento** em  $X$  consiste em encontrar vários **centros** de grupos (clusters) que possam caracterizar adequadamente categorias relevantes de  $X$ .



HAC

IA2022

27

27

## Agrupamento particional

- Gera uma única partição nos dados.
- Vantagem:
  - mais eficiente para conjunto de dados grande.
- Desvantagem:
  - é necessário definir previamente o número de grupos desejável.
- Gera grupos pela otimização de uma função critério (objetivo).

HAC

IA2022

28

28

## Agrupamento particional de erro quadrático

- A função critério mais utilizada em algoritmos de agrupamento particionais é o erro quadrático
- O erro quadrático de um agrupamento  $C$  com  $k$  clusters de um conjunto de padrões  $E$  é:

$$err2(E, C) = \sum_{j=1}^k \sum_{i=1}^{n_j} \|x_i^{(j)} - c_j\|^2$$

- onde  $x_i^{(j)}$  é o  $i$ -ésimo exemplo pertencente ao  $j$ -ésimo cluster e  $c_j$  é o centroide do  $j$ -ésimo cluster

HAC

IA2022

29

29

## Agrupamento k-means (k-médias)

- Mais simples e mais conhecido algoritmo de erro quadrático.
- É fácil de implementar e sua complexidade é  $O(n)$  com  $n$  sendo o número de exemplos;
- O usuário define previamente o número de grupos  $k$ ;
- Na versão original e na maioria das aplicações usa distância euclidiana quadrática para calcular similaridade entre os elementos;
- Problema: é sensível à partição inicial e pode convergir para um mínimo local do valor da função critério se a partição inicial não for escolhida apropriadamente.

HAC

IA2022

30

30

## Agrupamento k-means (k-médias)

Parâmetros de entrada:

- Conjunto de  $N$  exemplos não rotulados  $x_i, i=1, \dots, N$
- $k$  - número de grupos
- $dist$  - medida de distância

Parâmetros de saída:

- $k$  vetores que representam centroides de grupos

HAC

IA2022

31

31

## Agrupamento k-means (k-médias)

Escolha aleatoriamente um conjunto de vetores distintos para representar os centroides  $c_j, j = 1, \dots, k$

Repeat

For  $i=1$  to  $N$

Calcule a distância  $dist(x_i, c_j)$  de  $x_i$  a cada centroide  $c_j, j = 1, \dots, k$

Associe  $x_i$  ao centroide  $c_j$  que minimiza  $dist(x_i, c_j)$

End {For}

For  $j=1$  to  $k$

Atualize os centroides  $c_j$ , calculando a média dos exemplos que pertencem ao grupo com centro  $c_j$

End {For}

Until nenhuma mudança nos  $c_j$  ocorra entre duas iterações sucessivas

HAC

IA2022

32

32



## Agrupamento k-means - Exemplo

Considere os exemplos:

$$x_1 = [2, 5]$$

$$x_2 = [6, 4]$$

$$x_3 = [5, 3]$$

$$x_4 = [2, 2]$$

$$x_5 = [1, 4]$$

$$x_6 = [5, 2]$$

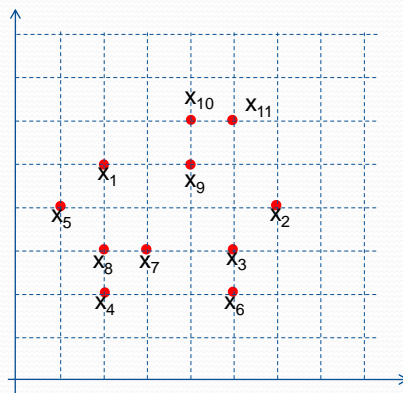
$$x_7 = [3, 3]$$

$$x_8 = [2, 3]$$

$$x_9 = [4, 5]$$

$$x_{10} = [4, 6]$$

$$x_{11} = [5, 6]$$



A distância entre exemplos é calculada com a distância euclidiana

33

HAC

IA2022

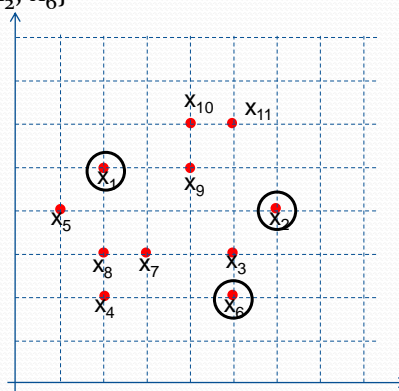
33

## Agrupamento k-means - Exemplo

- Aplicar o algoritmo k-means para  $k = 3$ 
  - Gerar aleatoriamente o vetor de centróides inicial:

$$\Theta = \{x_1, x_2, x_6\}$$

<b>c1</b>	2	5
<b>c2</b>	6	4
<b>c3</b>	5	2



HAC

IA2022

34

34

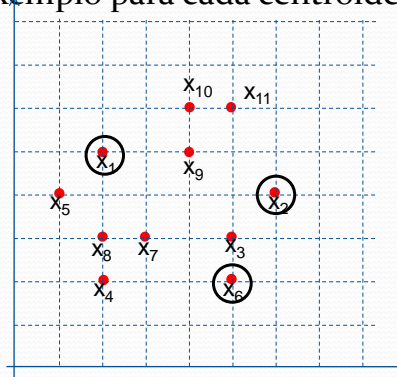
## Agrupamento k-means - Exemplo

- Gerar aleatoriamente o vetor de centróides inicial:

$$\Theta = \{x_1, x_2, x_6\}$$

- Calcular a distância de cada exemplo para cada centróide:

	c1	c2	c3
x1	0	4,1	4,2
x2	4,1	0,0	2,2
x3	3,6	1,4	1,0
x4	3,0	4,5	3,0
x5	1,4	5,0	4,5
x6	4,2	2,2	0,0
x7	2,2	3,2	2,2
x8	2,0	4,1	3,2
x9	2,0	2,2	3,2
x10	2,2	2,8	4,1
x11	3,2	2,2	4,0



HAC

IA2022

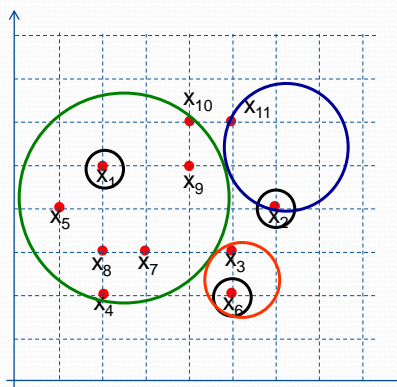
35

35

## Agrupamento k-means - Exemplo

- Atribuir cada exemplo a um cluster, pela menor distância ao centróide:

	c1	c2	c3	grupo
x1	0	4,1	4,2	c1
x2	4,1	0,0	2,2	c2
x3	3,6	1,4	1,0	c3
x4	3,0	4,5	3,0	c1
x5	1,4	5,0	4,5	c1
x6	4,2	2,2	0,0	c3
x7	2,2	3,2	2,2	c1
x8	2,0	4,1	3,2	c1
x9	2,0	2,2	3,2	c1
x10	2,2	2,8	4,1	c1
x11	3,2	2,2	4,0	c2



HAC

IA2022

36

36

## Agrupamento k-means - Exemplo

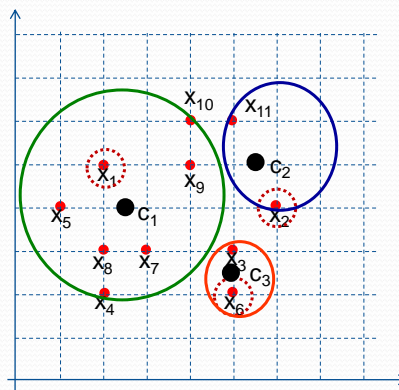
- Recalcular os centróides:

Atual:

<b>c1</b>	2	5
<b>c2</b>	6	4
<b>c3</b>	5	2

Novo:

<b>c1</b>	2,6	4
<b>c2</b>	5,5	5
<b>c3</b>	5	2,5



HAC

IA2022

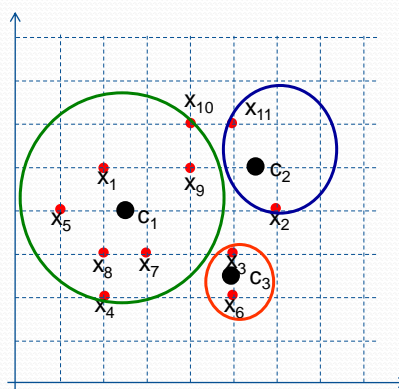
37

37

## Agrupamento k-means - Exemplo

- Calcular a distância de cada exemplo para cada centróide:
- Atribuir cada exemplo a um cluster, pela menor distância ao centróide:

	c1	c2	c3	
x1	1,2	3,5	3,9	c1
x2	3,4	1,1	1,8	c2
x3	2,6	2,1	0,5	c3
x4	2,1	4,6	3,0	c1
x5	1,6	4,6	4,3	c1
x6	3,1	3,0	0,5	c3
x7	1,1	3,2	2,1	c1
x8	1,2	4,0	3,0	c1
x9	1,7	1,5	2,7	c2
x10	2,4	1,8	3,6	c2
x11	3,1	1,1	3,5	c2



HAC

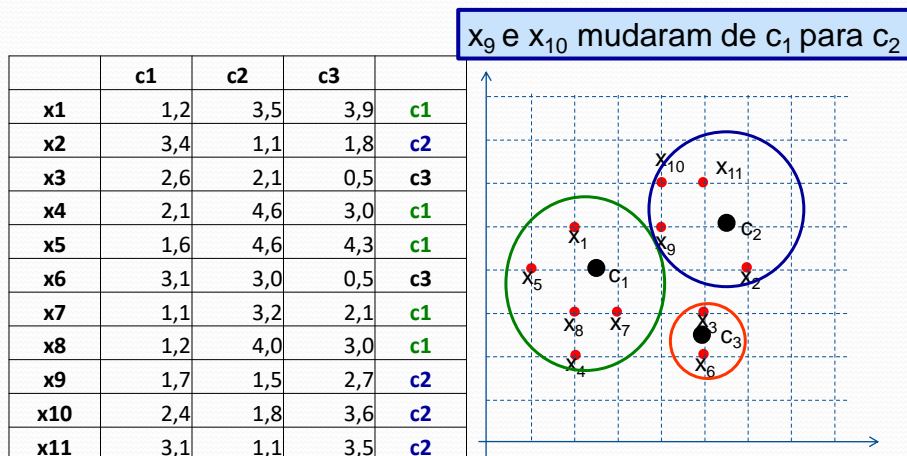
IA2022

38

38

## Agrupamento k-means - Exemplo

- Atribuir cada exemplo a um cluster, pela menor distância ao centróide:



HAC

IA2022

39

39

## Agrupamento k-means - Exemplo

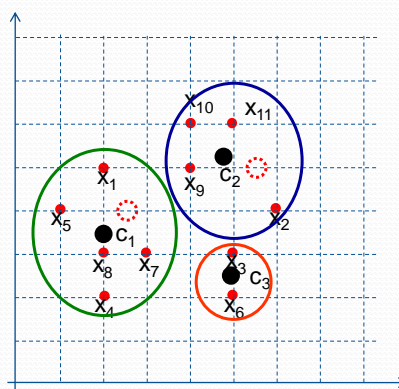
- Recalcular os centróides:

Atual:

c1	2,6	4
c2	5,5	5
c3	5	2,5

Novo:

c1	2,0	3,4
c2	4,8	5,3
c3	5	2,5



HAC

IA2022

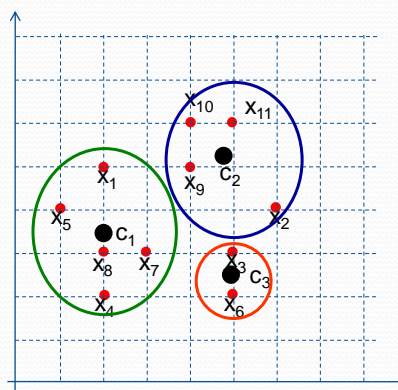
40

40

## Agrupamento k-means - Exemplo

- Calcular a distância de cada exemplo para cada centróide:
- Atribuir cada exemplo a um cluster, pela menor distância ao centróide:

	c1	c2	c3	
x1	1,6	2,8	3,9	c1
x2	4,0	1,77	1,80	c2
x3	3,0	2,3	0,5	c3
x4	1,4	4,3	3,0	c1
x5	1,2	4,0	4,3	c1
x6	3,3	3,3	0,5	c3
x7	1,1	2,9	2,1	c1
x8	0,4	3,6	3,0	c1
x9	2,6	0,9	2,7	c2
x10	3,3	1,1	3,6	c2
x11	4,0	0,7	3,5	c2



HAC

IA2022

41

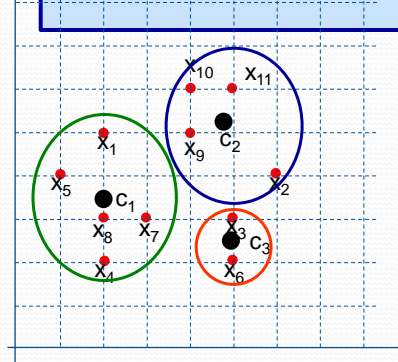
41

## Agrupamento k-means - Exemplo

- Atribuir cada exemplo a um cluster, pela menor distância ao centróide:

	c1	c2	c3	
x1	1,6	2,8	3,9	c1
x2	4,0	1,77	1,80	c2
x3	3,0	2,3	0,5	c3
x4	1,4	4,3	3,0	c1
x5	1,2	4,0	4,3	c1
x6	3,3	3,3	0,5	c3
x7	1,1	2,9	2,1	c1
x8	0,4	3,6	3,0	c1
x9	2,6	0,9	2,7	c2
x10	3,3	1,1	3,6	c2
x11	4,0	0,7	3,5	c2

Não houve mudanças  
Processo termina



HAC

IA2022

42

42

**Bons estudos!**

Até a próxima aula

