



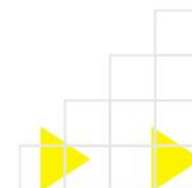
# Curso 2 – CD, AM e DM

## Mineração de Dados

### Parte 4

Extração de Padrões  
Agrupamento Particional

**Prof. Ricardo M. Marcacini**  
[ricardo.marcacini@icmc.usp.br](mailto:ricardo.marcacini@icmc.usp.br)

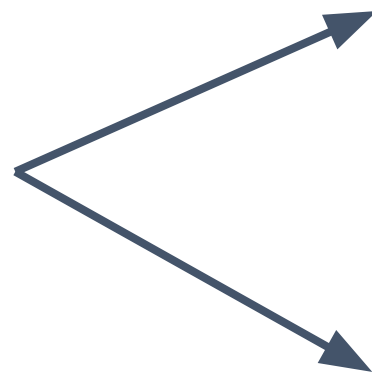
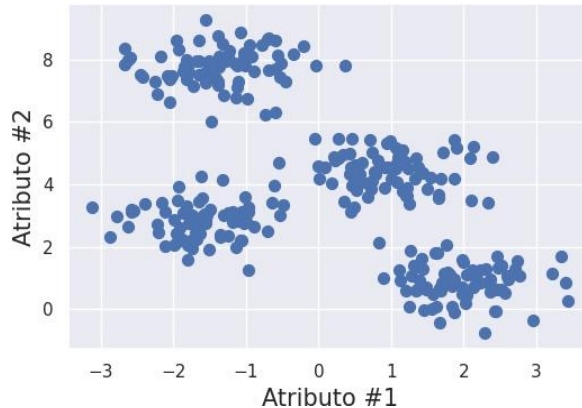


# Métodos para Agrupamento de Dados

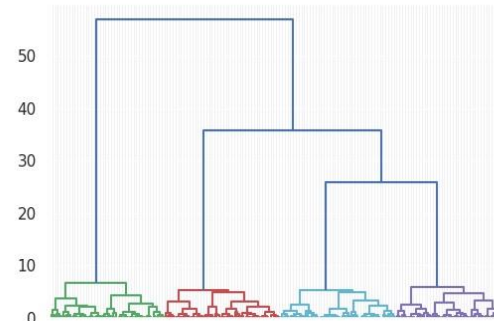


- **Particionais:** organizar dados em uma partição de  $k$  *clusters*
- **Hierárquicos:** organizar dados em uma decomposição hierárquica de *clusters* e *subclusters*

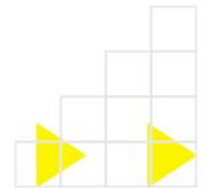
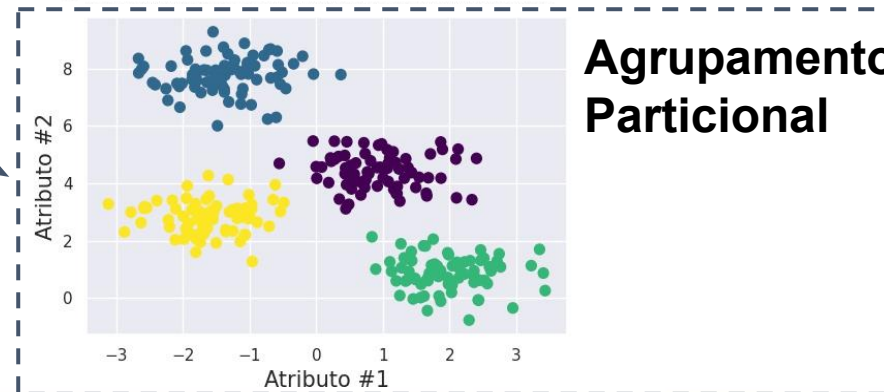
Conjunto de Dados



Agrupamento Hierárquico



Agrupamento Particional

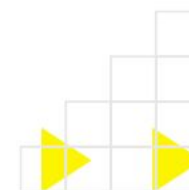


# Agrupamento Particional



- Falaremos sobre métodos de agrupamento para obter partições rígidas dos dados
- **Partição rígida:** clusters não possuem sobreposição
  - Dado um conjunto de  $n$  objetos

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$$



# Agrupamento Particional



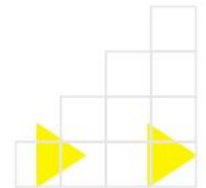
- Falaremos sobre métodos de agrupamento para obter partições rígidas dos dados
- **Partição rígida:** clusters não possuem sobreposição
  - Dado um conjunto de  $n$  objetos

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$$

- Obter um agrupamento  $C$  em  $k$  clusters

$$C = \{C_1, C_2, \dots, C_k\}$$

$$C_1 \cup C_2 \cup \dots \cup C_k = \mathbf{X}$$



# Agrupamento Particional



- Falaremos sobre métodos de agrupamento para obter partições rígidas dos dados
- **Partição rígida:** clusters não possuem sobreposição
  - Dado um conjunto de  $n$  objetos

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$$

- Obter um agrupamento  $C$  em  $k$  clusters

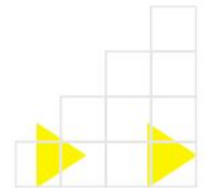
$$C = \{C_1, C_2, \dots, C_k\}$$

$$C_1 \cup C_2 \cup \dots \cup C_k = \mathbf{X}$$

- Sem clusters vazios e sem sobreposição

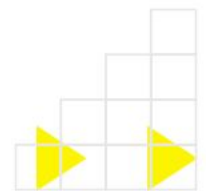
$$C_i \neq \emptyset$$

$$C_i \cap C_j = \emptyset \text{ para } i \neq j$$



# Algoritmo *k*-Médias ou *k-Means*

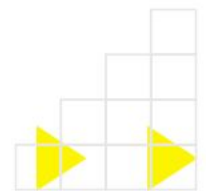
- Amplamente usado na indústria e academia
- Características desejáveis para Mineração de Dados
  - Simplicidade
  - Interpretabilidade
  - Eficiência Computacional



# Algoritmo *k*-Médias ou *k-Means*

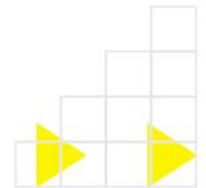
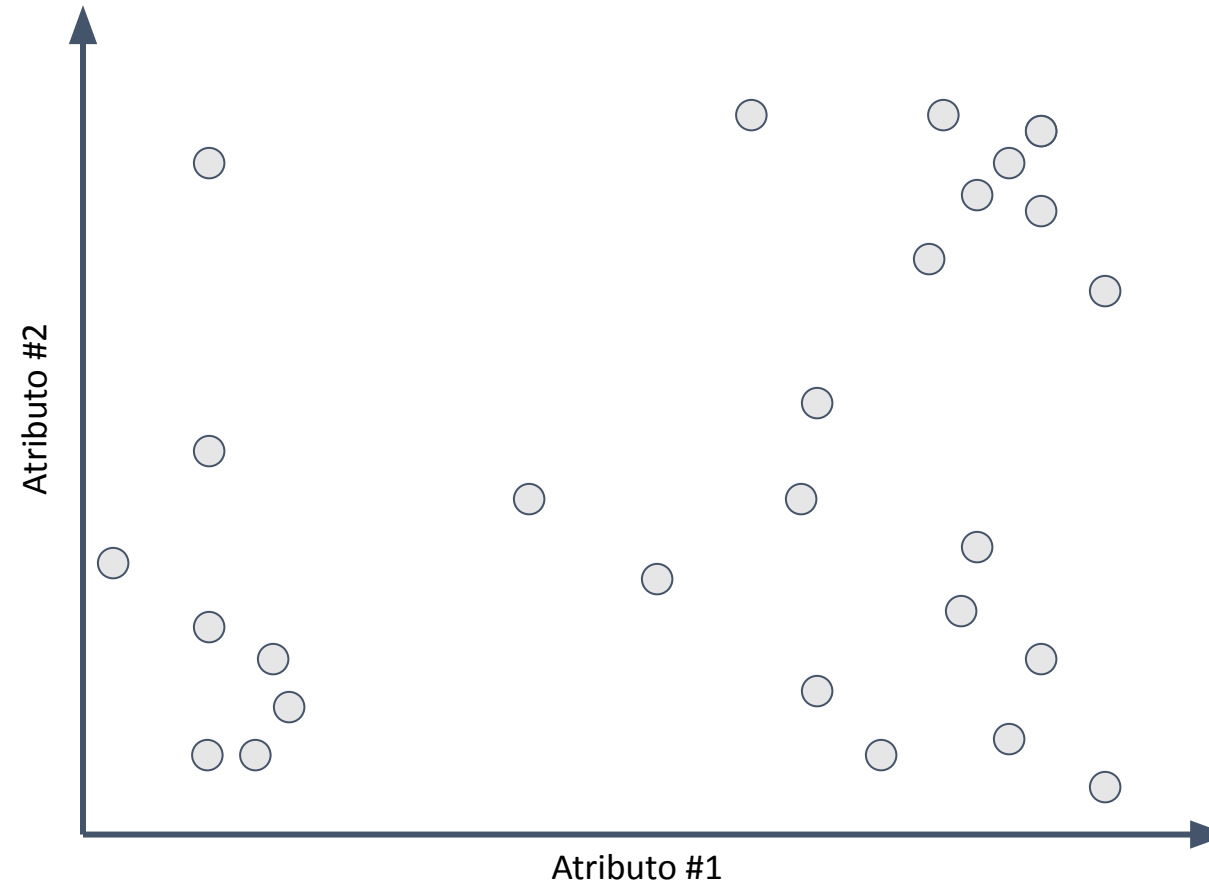
- Amplamente usado na indústria e academia
- Características desejáveis para Mineração de Dados
  - Simplicidade
  - Interpretabilidade
  - Eficiência Computacional

Vamos começar a estudar o *k-Means* a partir de um exemplo didático...





# Algoritmo *k-Means*

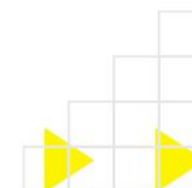
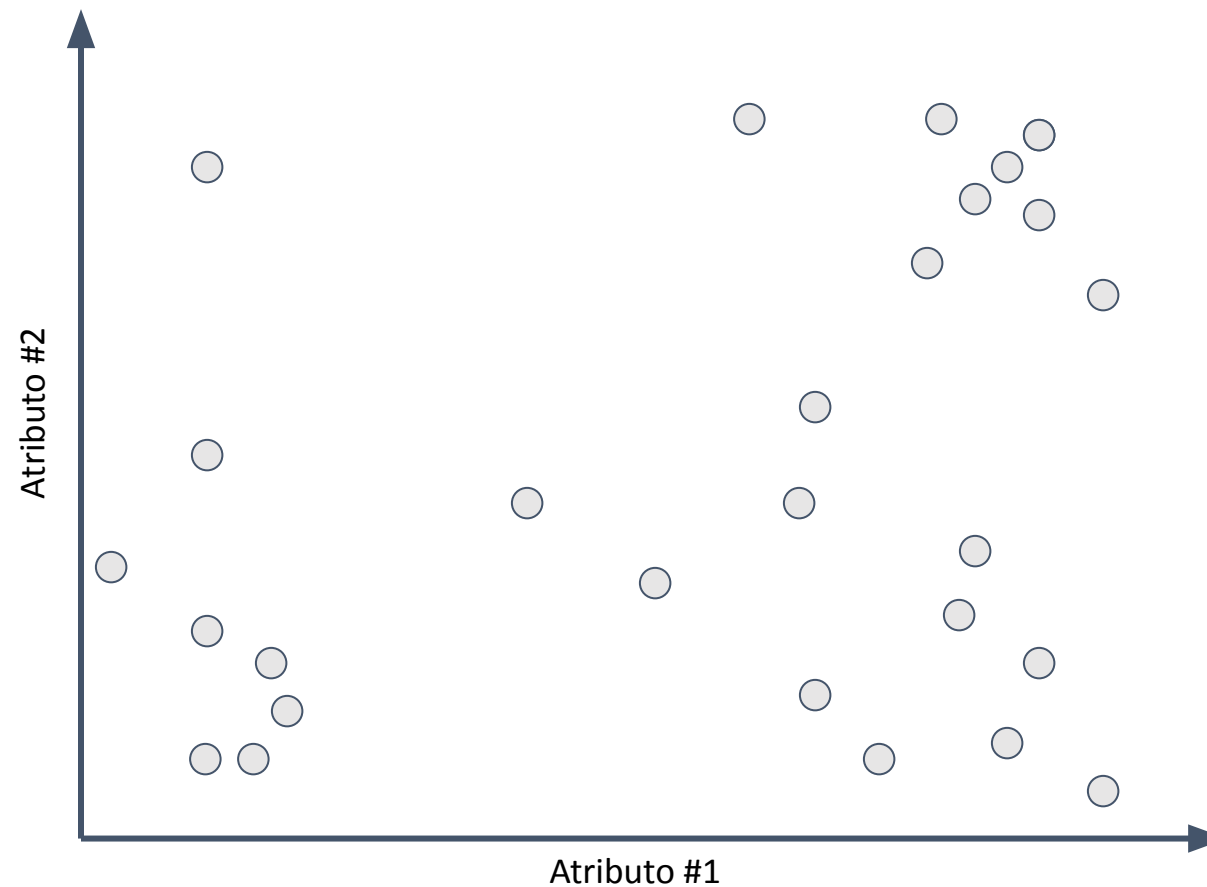


Exemplo adaptado de Gregory Piatetsky-Shapiro & Gary Parker ([www.kdnuggets.com](http://www.kdnuggets.com))



# Algoritmo *k-Means*

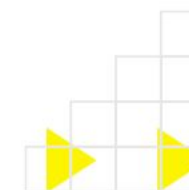
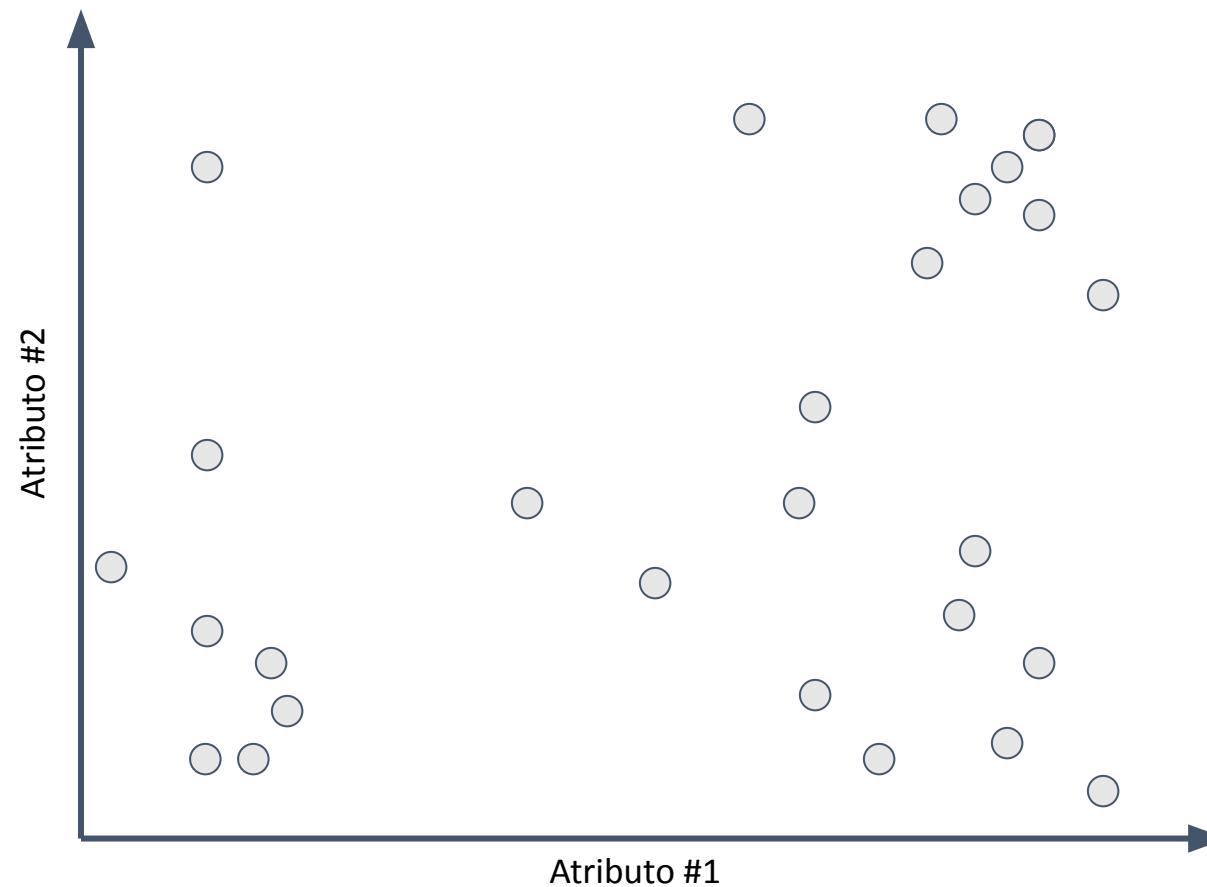
Primeiro passo é definir o número  $k$  de *clusters* que se deseja encontrar



# Algoritmo *k-Means*

Primeiro passo é definir o número  $k$  de *clusters* que se deseja encontrar

Vamos definir  **$k=3$**

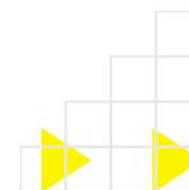
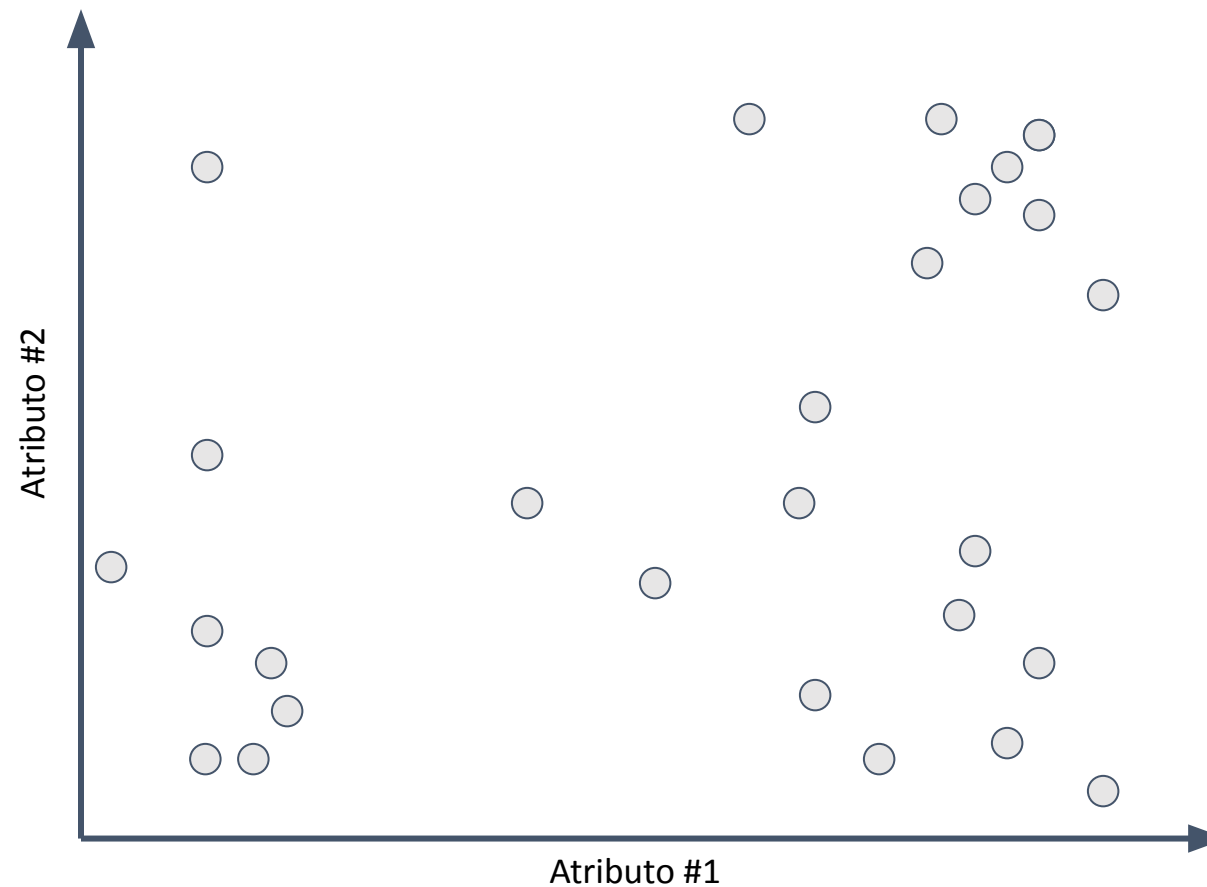


# Algoritmo *k-Means*

Com  $k=3$ , nós vamos inicializar  $k$  centroides.

Cada centroide é um ponto e representa um cluster.

Inicialização aleatória de centroides.

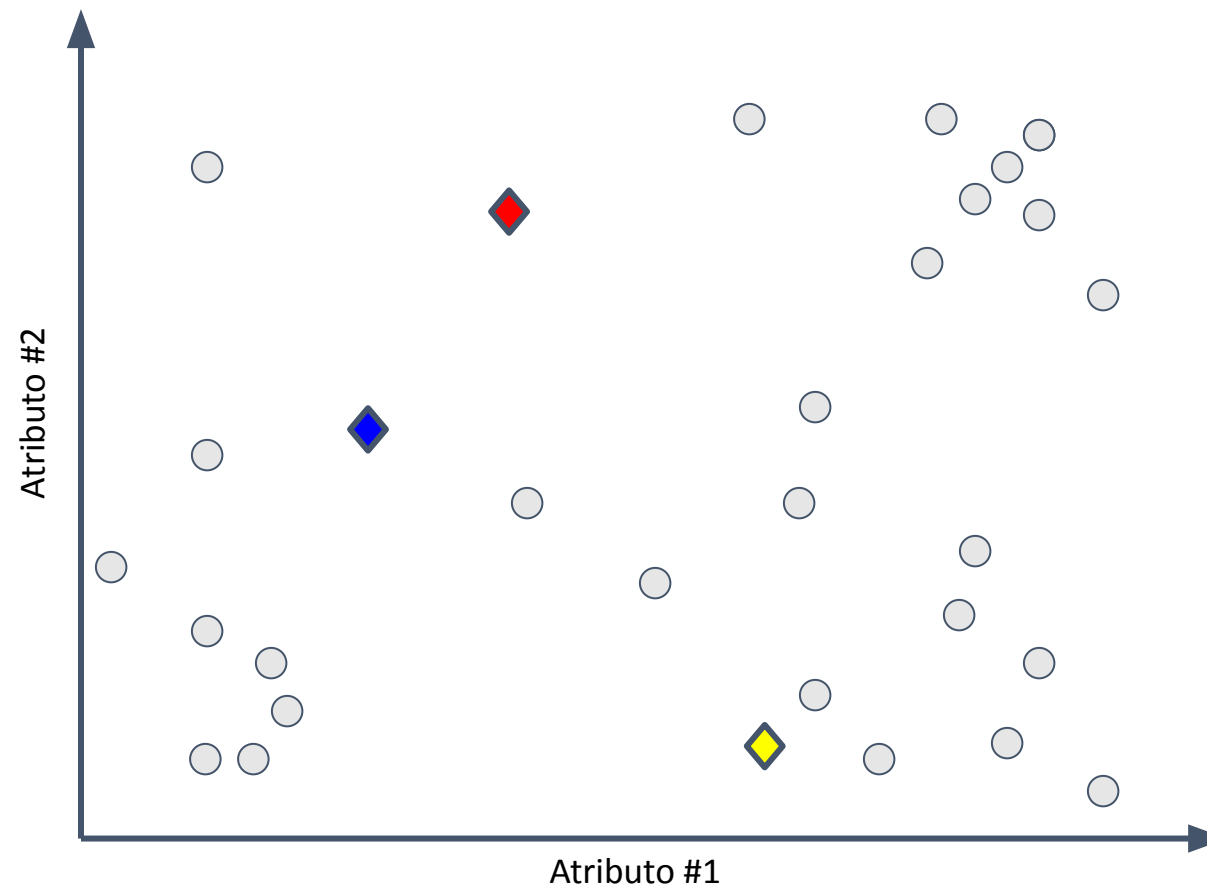


# Algoritmo *k-Means*

Com  $k=3$ , nós vamos inicializar  $k$  centroides.

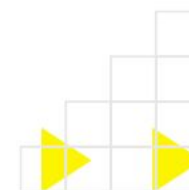
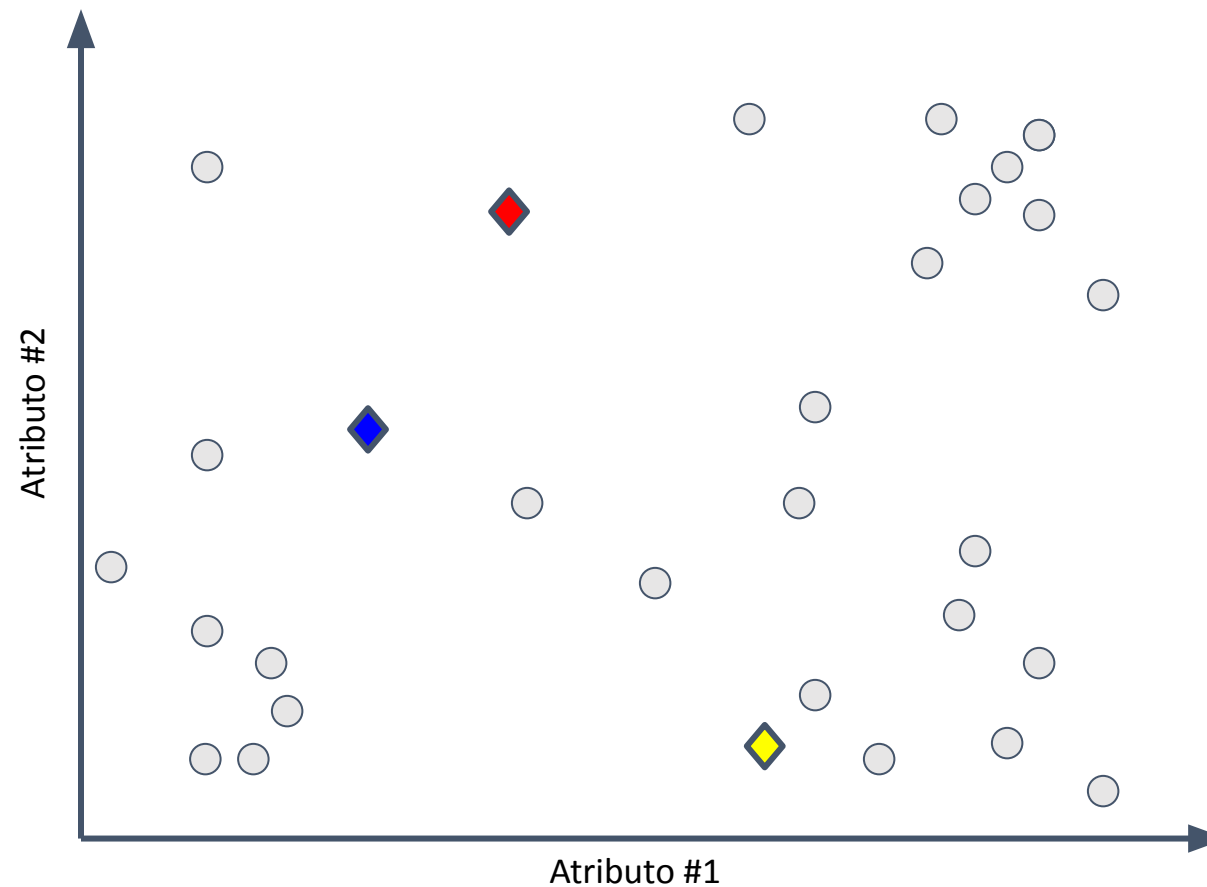
Cada centroide é um ponto e representa um cluster.

Inicialização aleatória de centroides.



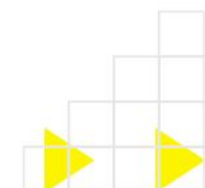
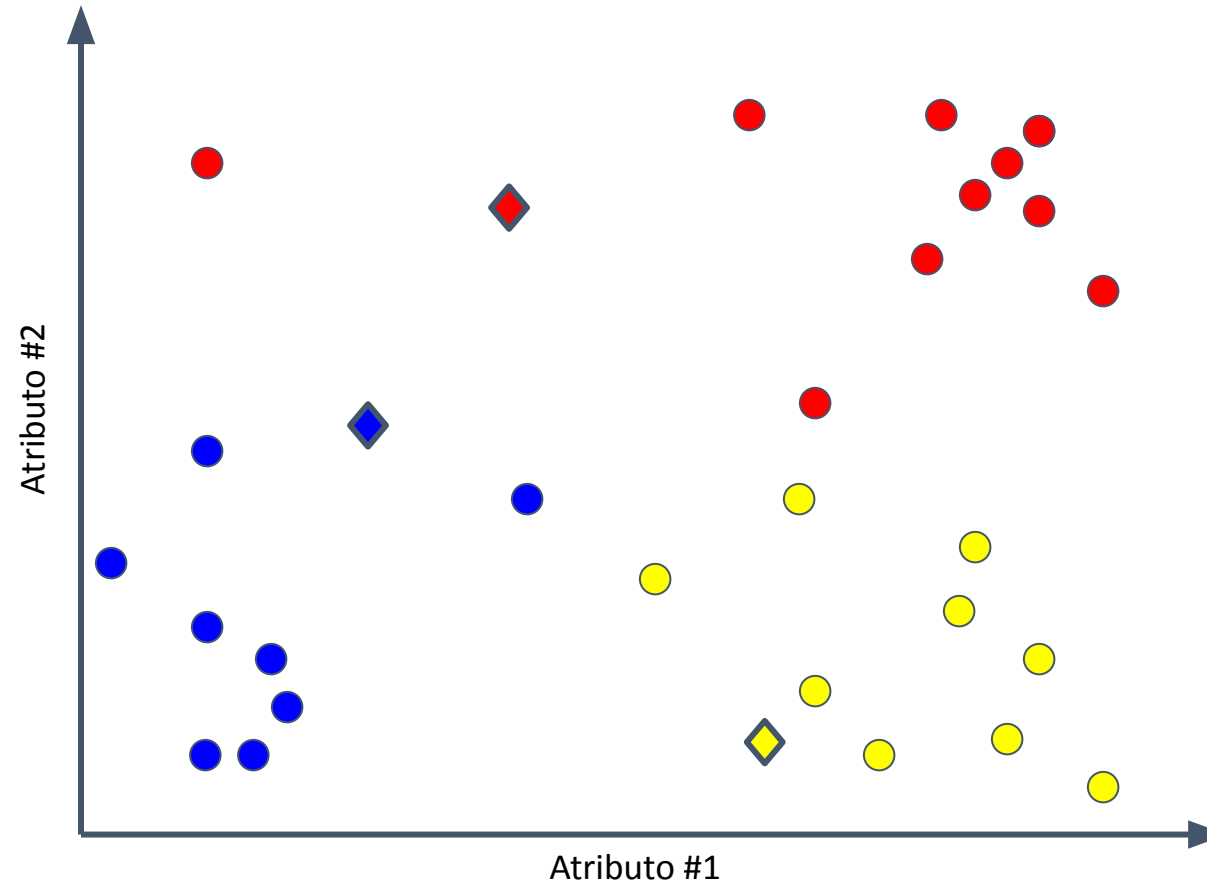
# Algoritmo *k-Means*

Agora, nós associamos cada objeto ao centróide mais próximo.



# Algoritmo *k-Means*

Agora, nós associamos cada objeto ao centróide mais próximo.

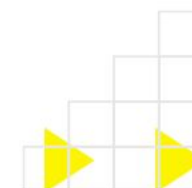
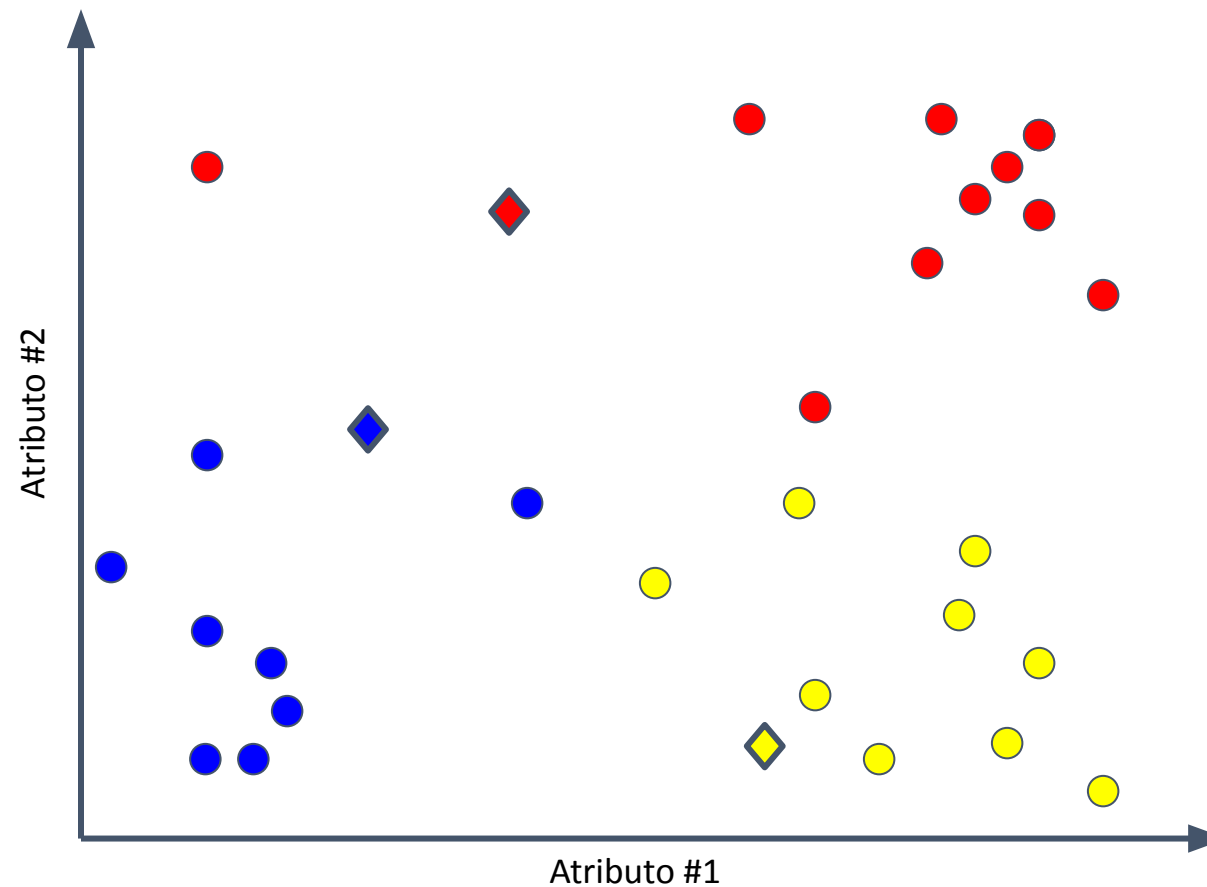


# Algoritmo *k-Means*



Agora, nós atualizamos o centroide de cada cluster.

O centroide é calculado como o vetor médio do cluster.



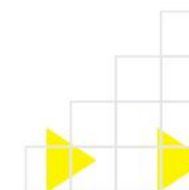
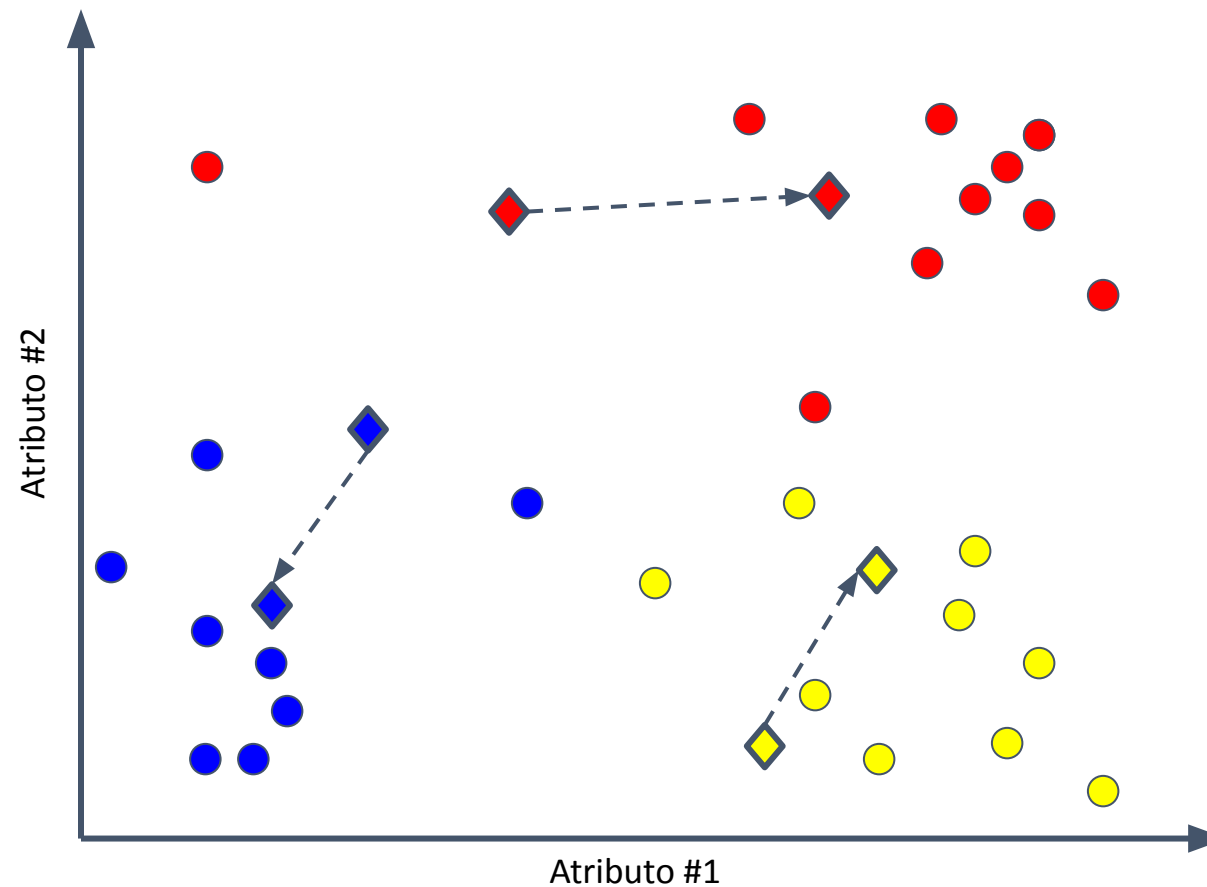


# Algoritmo *k-Means*



Agora, nós atualizamos o centroide de cada cluster.

O centroide é calculado como o vetor médio do cluster.



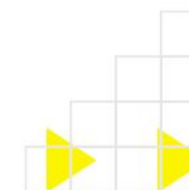
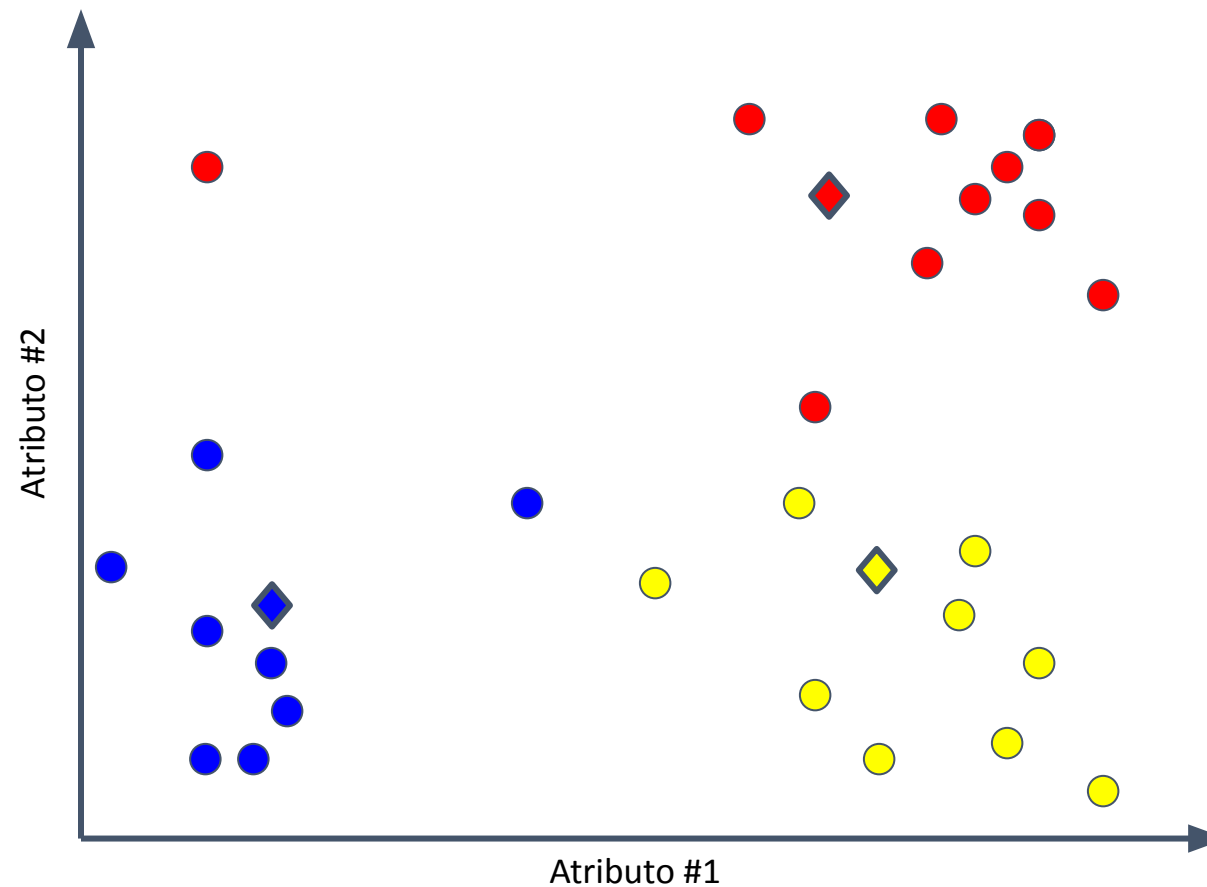
# Algoritmo *k-Means*



Agora, nós atualizamos o centroide de cada cluster.

O centroide é calculado como o vetor médio do cluster.

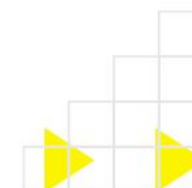
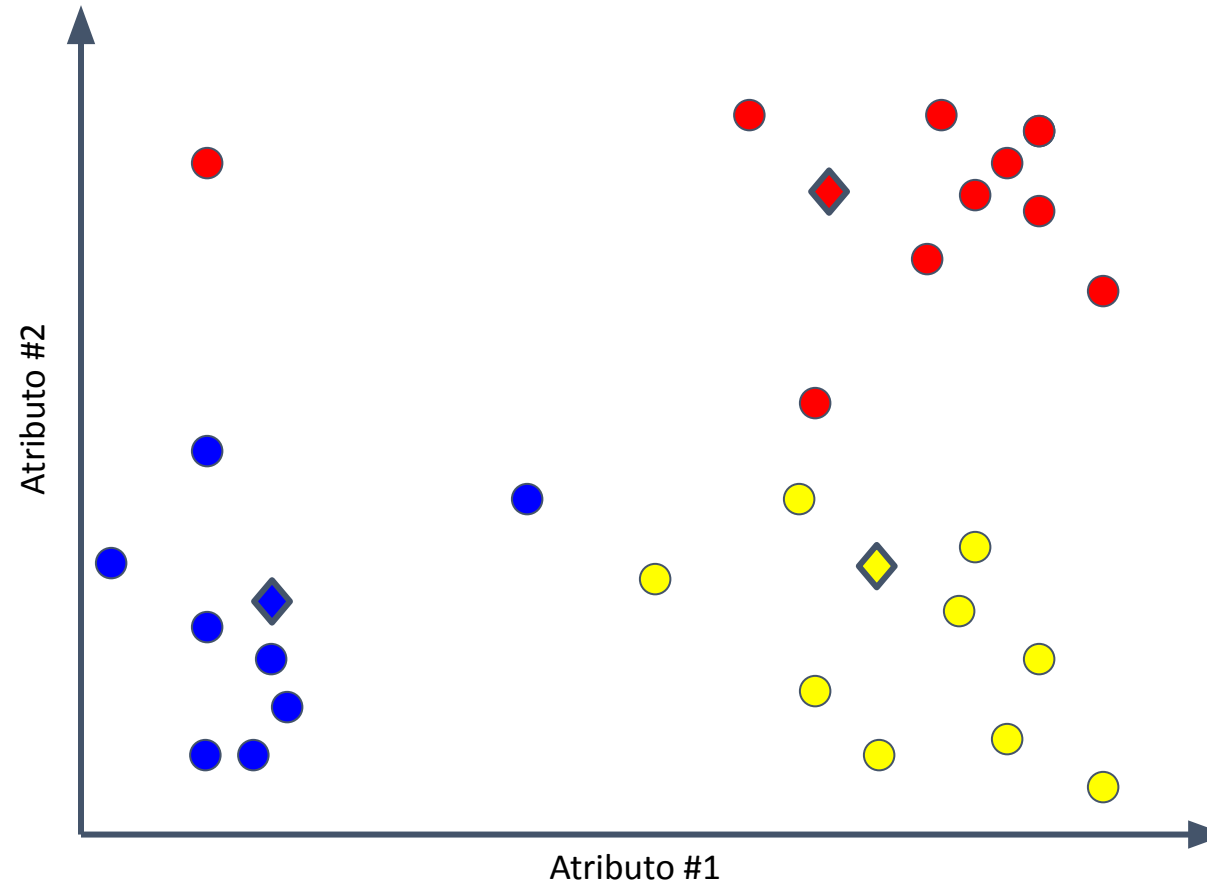
**Confira os novos centroides obtidos.**



# Algoritmo *k-Means*

Repetir até convergir:

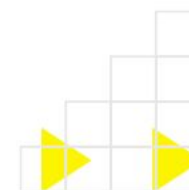
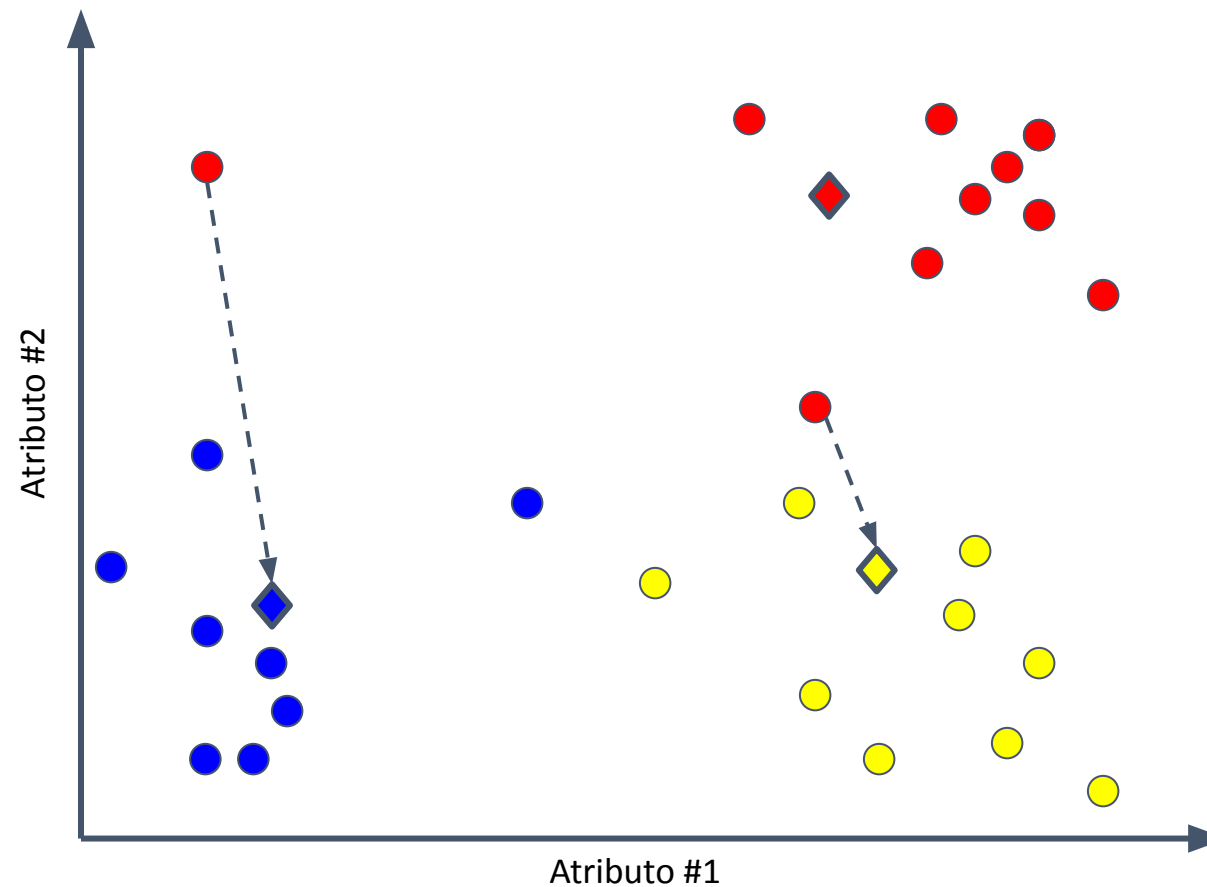
- Alocar objetos ao centróide mais próximo
- Atualizar centroides



# Algoritmo *k-Means*

Repetir até convergir:

- Alocar objetos ao centróide mais próximo
- Atualizar centroides

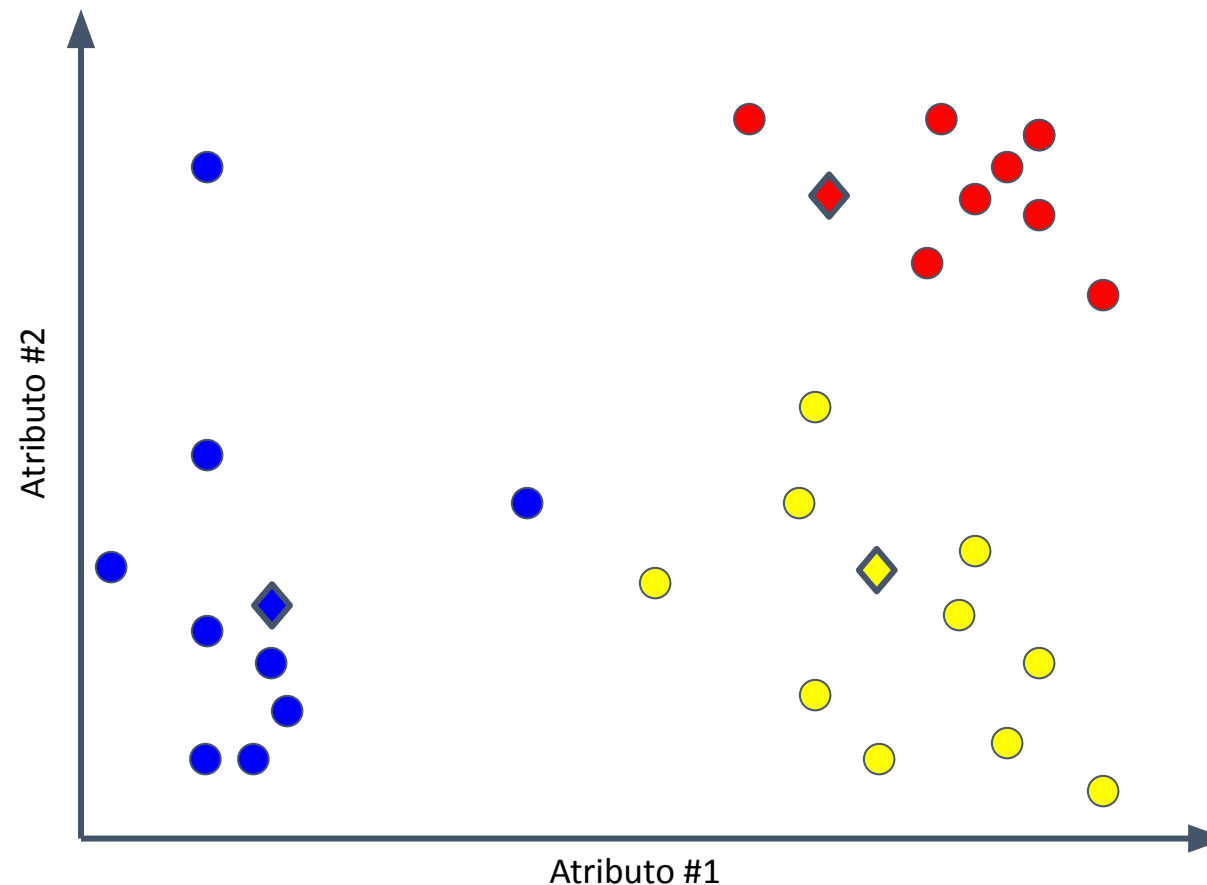


# Algoritmo *k-Means*

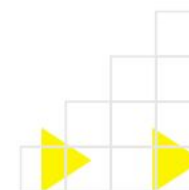


## Algoritmo:

1. Selecionar  $k$  centroides iniciais
2. Repetir até convergir:
  - 2.1. Formar  $k$  clusters atribuindo cada objeto ao centroide mais próximo
  - 2.2. Atualizar o centroide de cada cluster



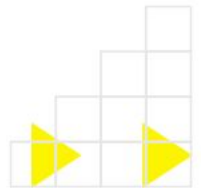
**Crítérios de convergência:** (1) poucas mudanças nos clusters/centroides;  
(2) número máximo de iterações.



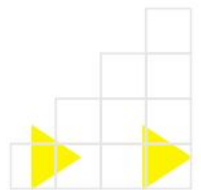
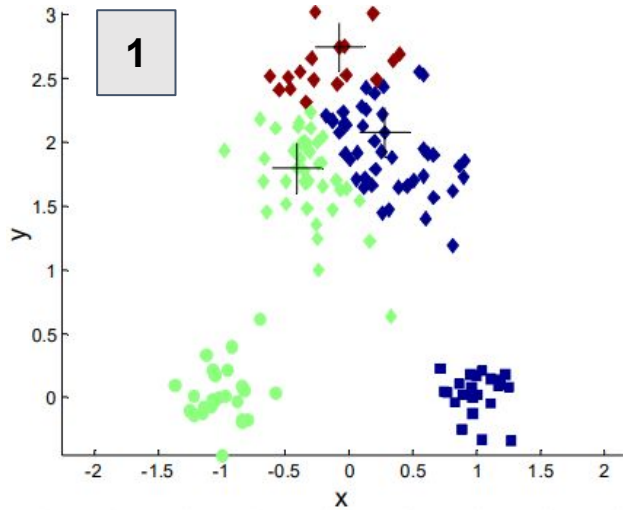
# Algoritmo *k-Means*

- Adequado para dados contínuos
- O cálculo da média faz sentido para seus dados?
- Medidas de proximidade para dados contínuos
  - Exemplo: distância euclidiana e dissimilaridade de cosseno
- Converge em poucas iterações

Em geral, os centroides iniciais são escolhidos aleatoriamente.  
*Clusters* podem ser diferentes em cada execução do k-means.

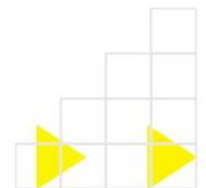
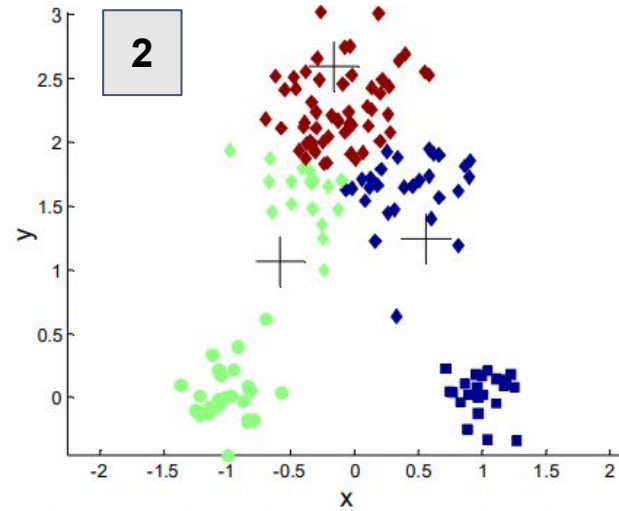
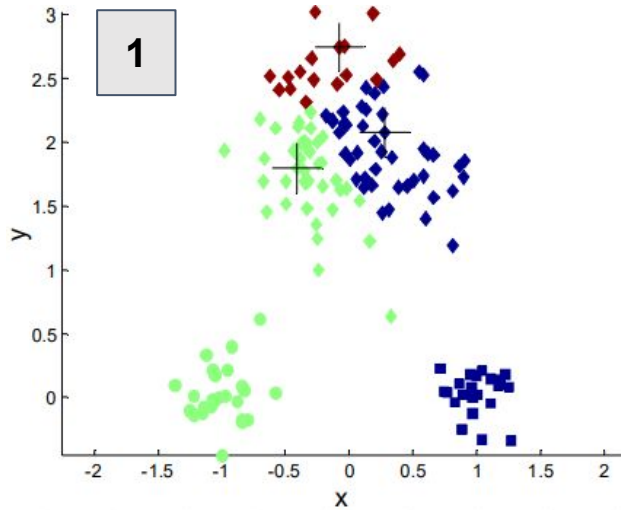


# Algoritmo *k-Means*

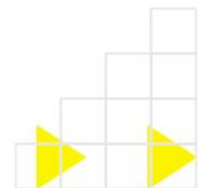
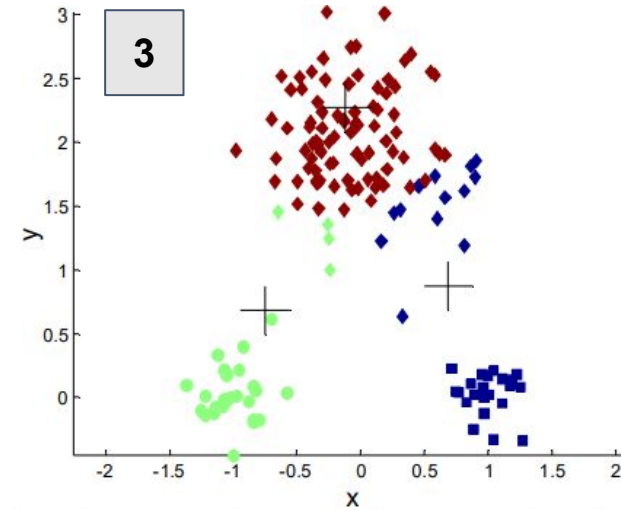
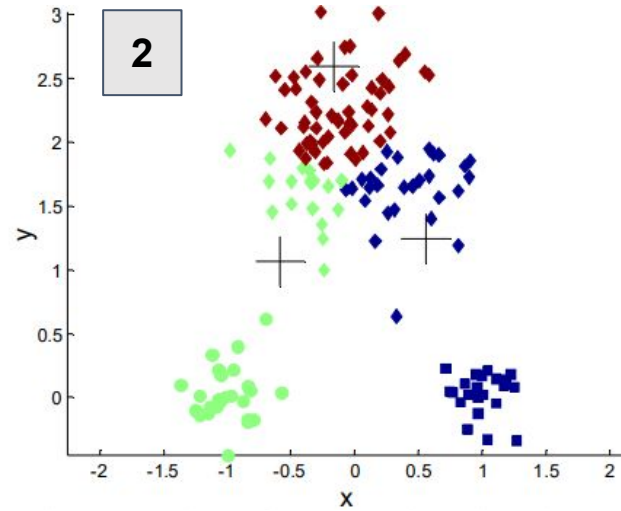
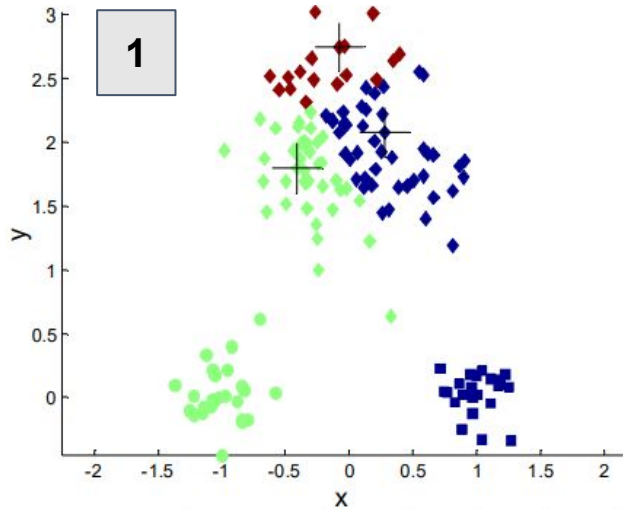




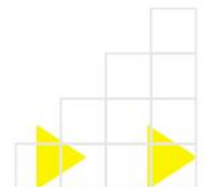
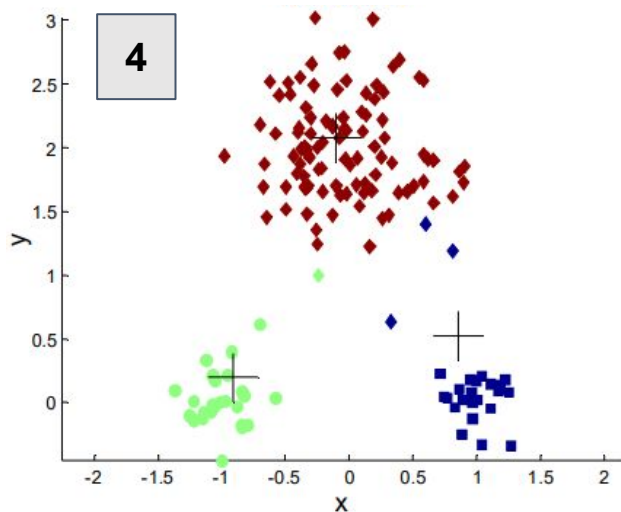
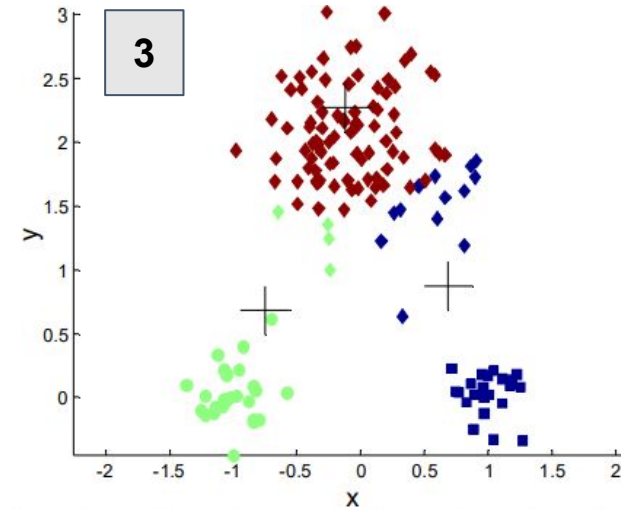
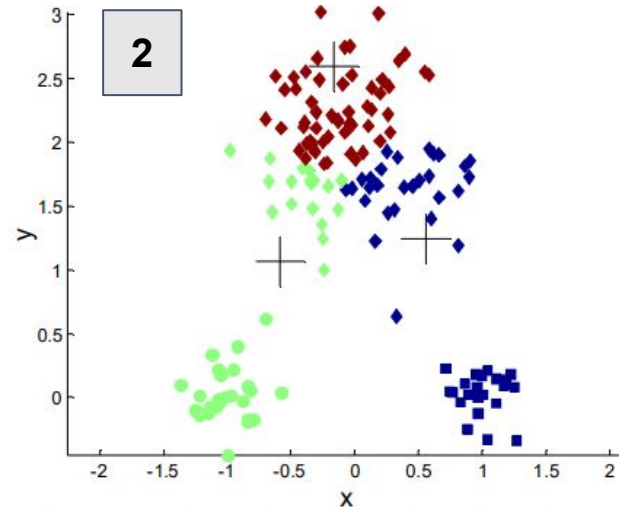
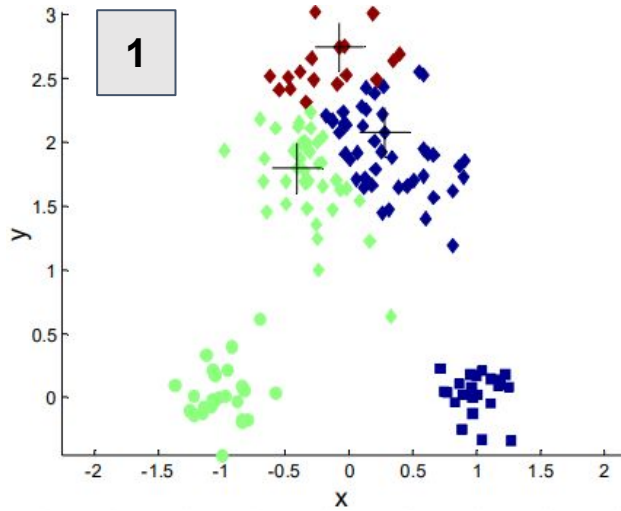
# Algoritmo *k-Means*



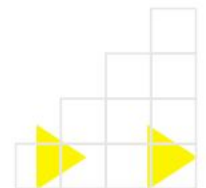
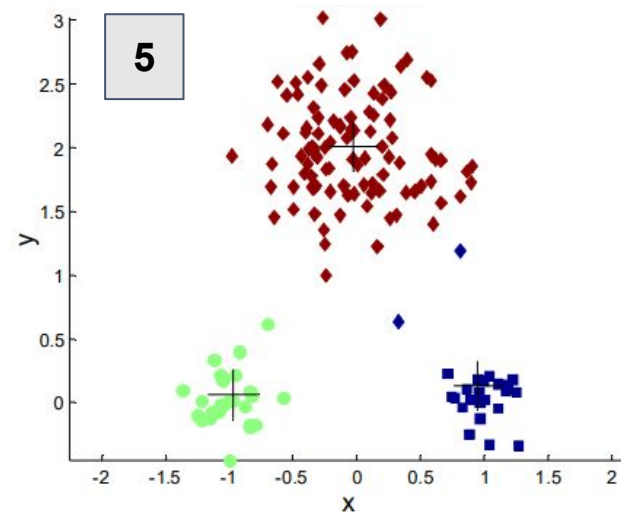
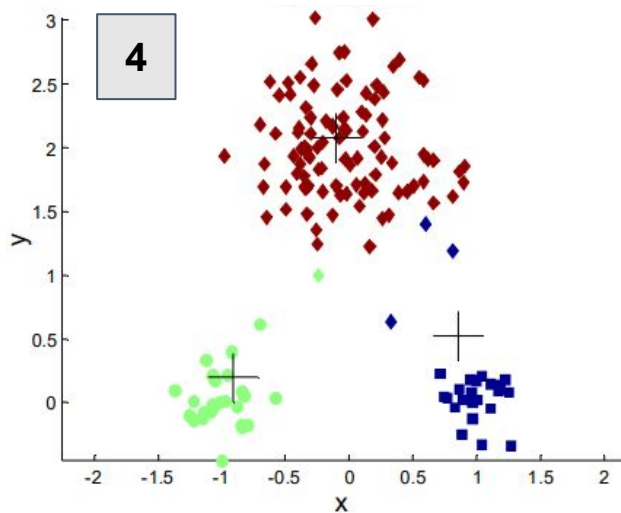
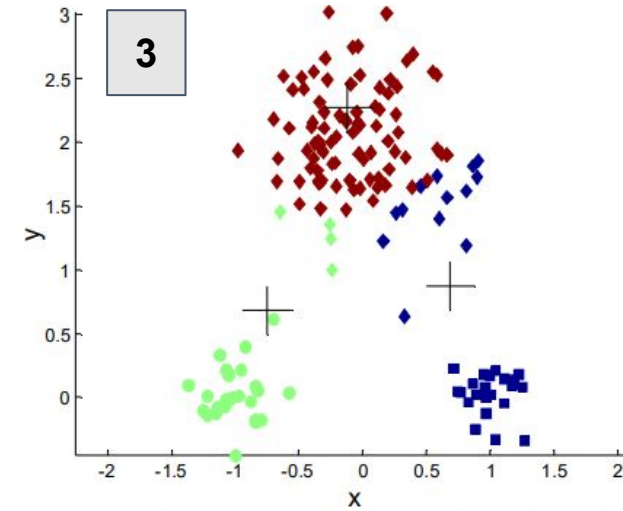
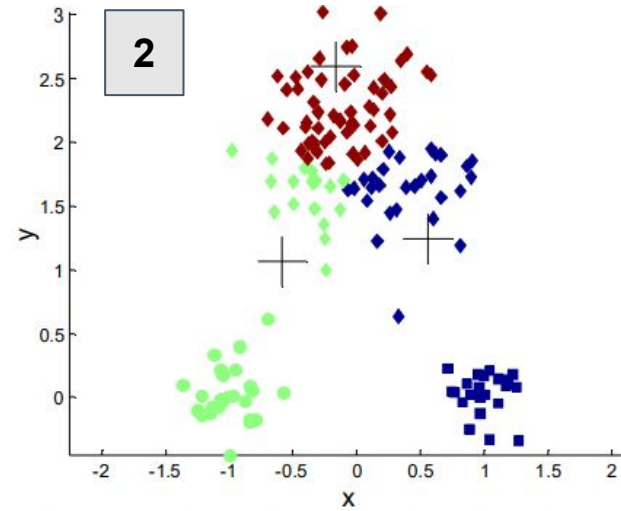
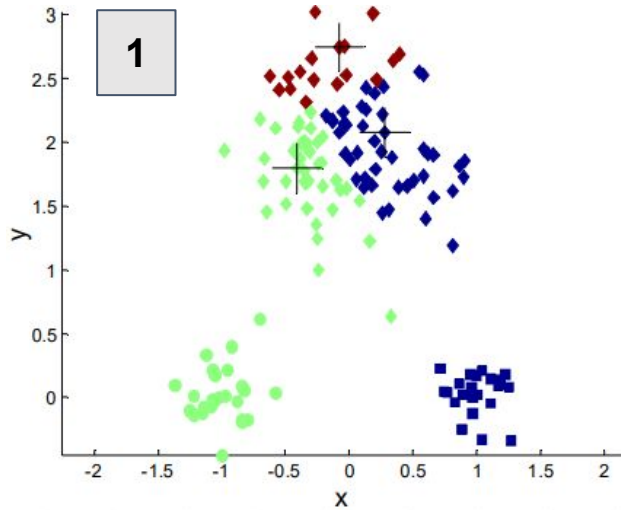
# Algoritmo *k-Means*



# Algoritmo *k*-Means

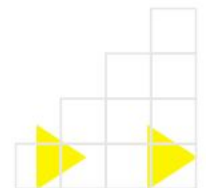
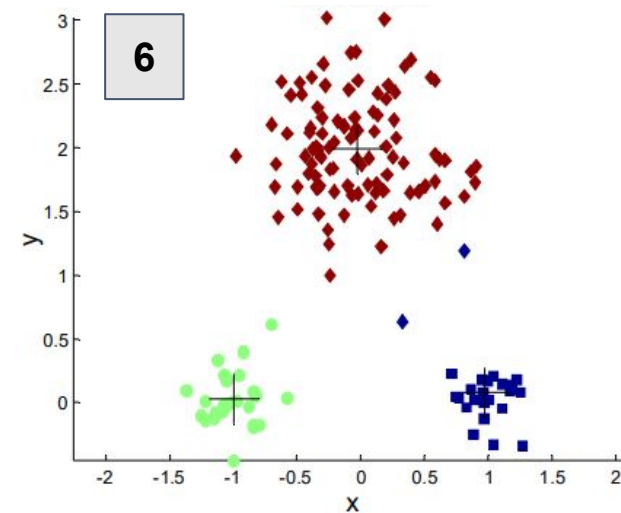
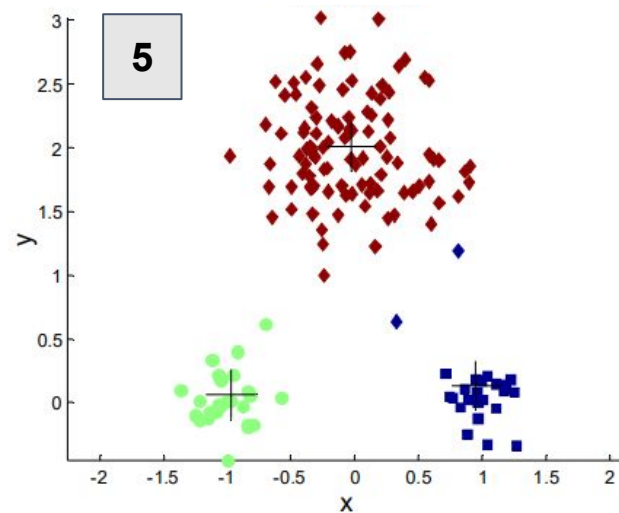
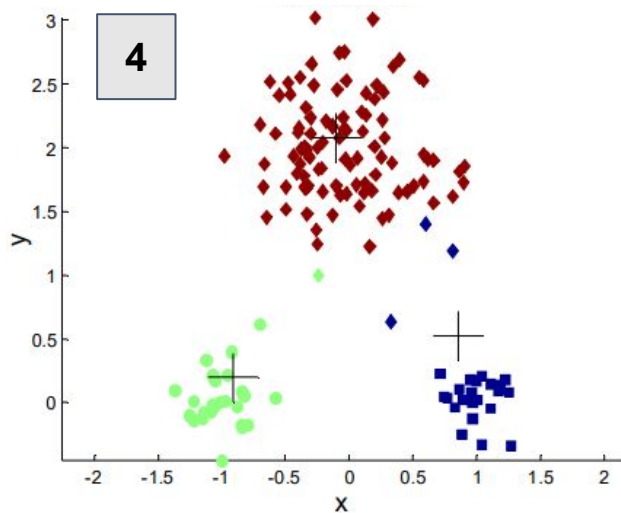
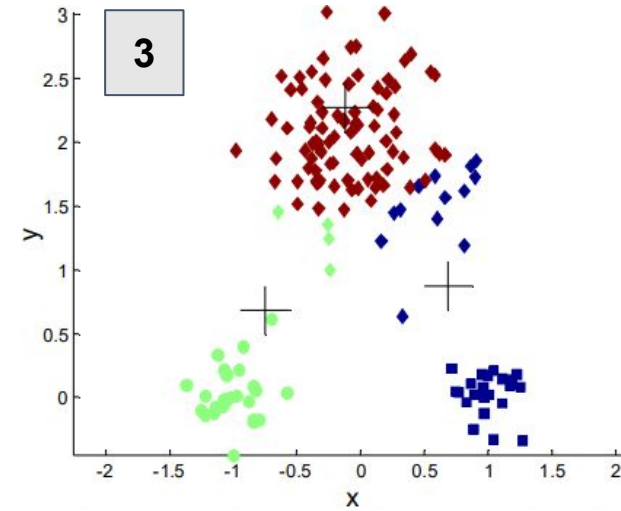
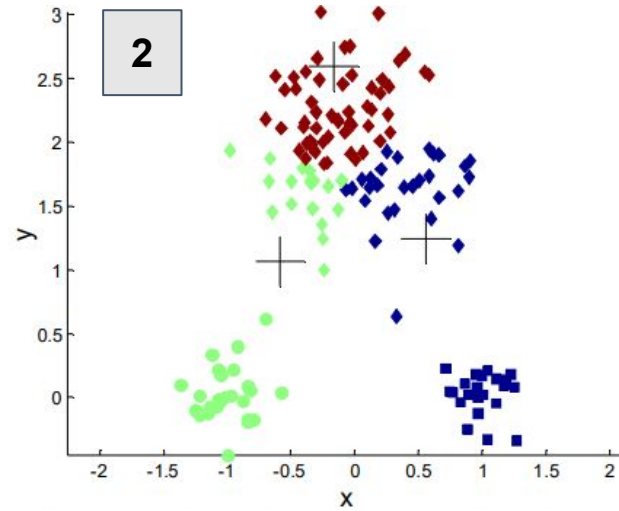
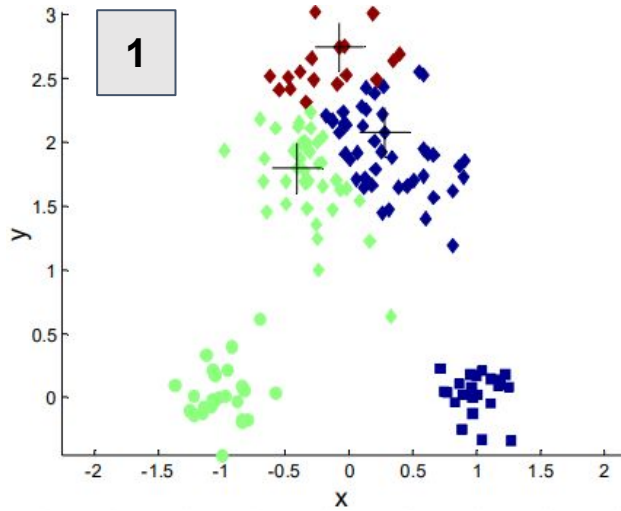


# Algoritmo *k*-Means

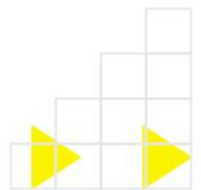
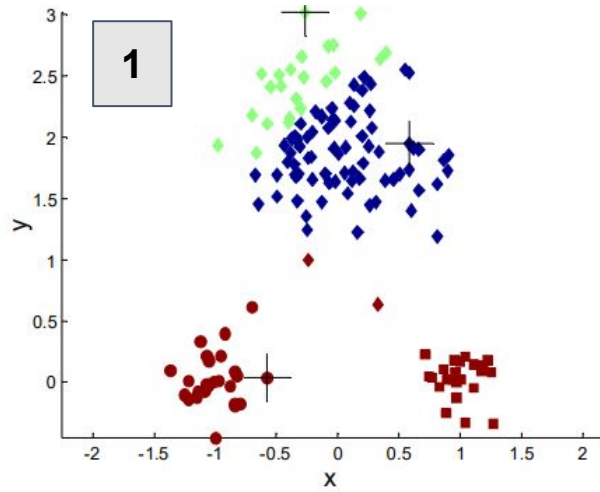




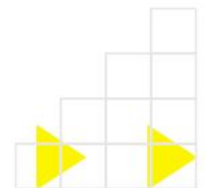
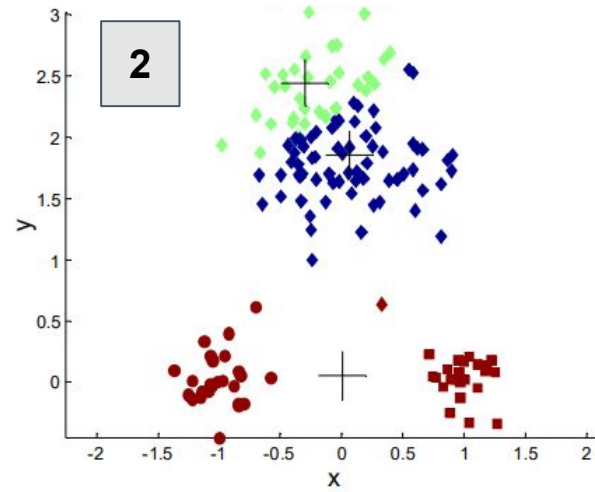
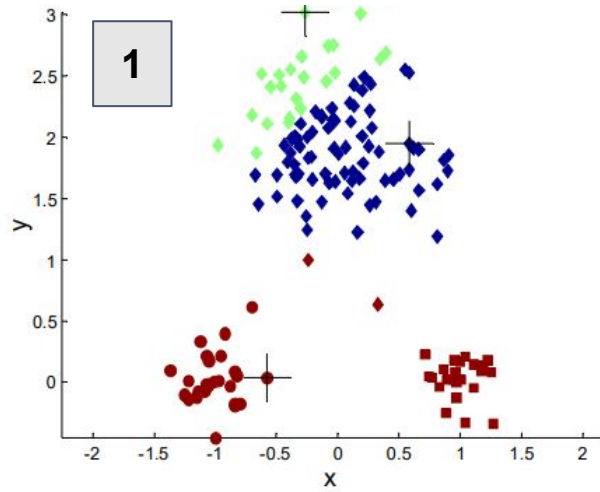
# Algoritmo *k-Means*



# Algoritmo *k-Means*

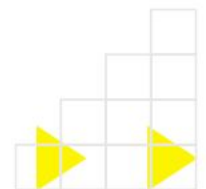
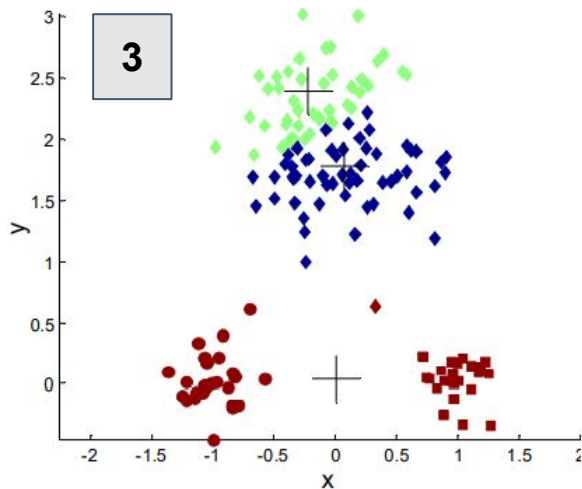
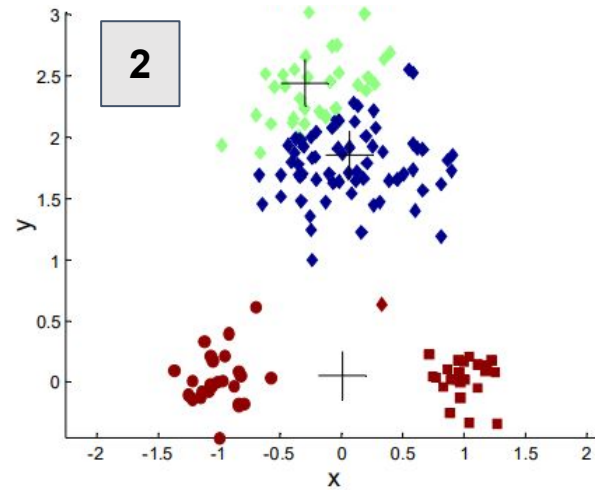
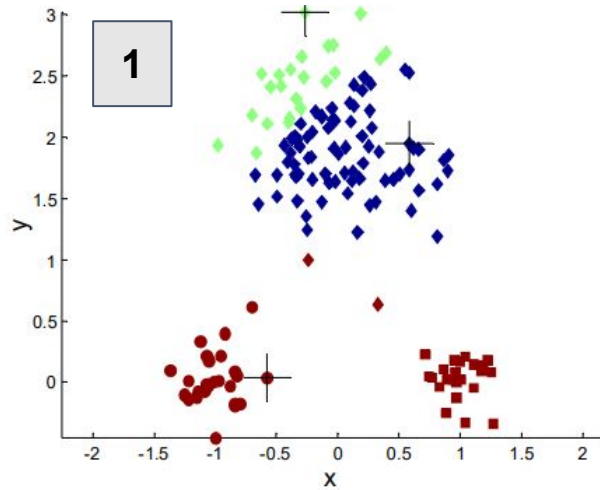


# Algoritmo *k-Means*

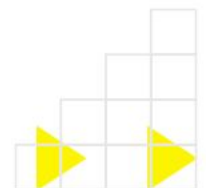
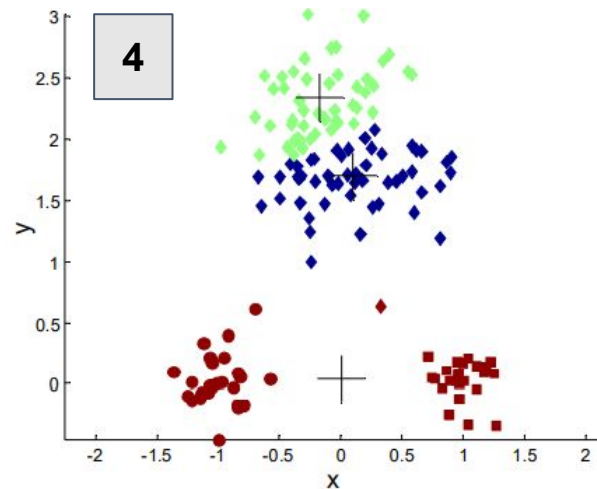
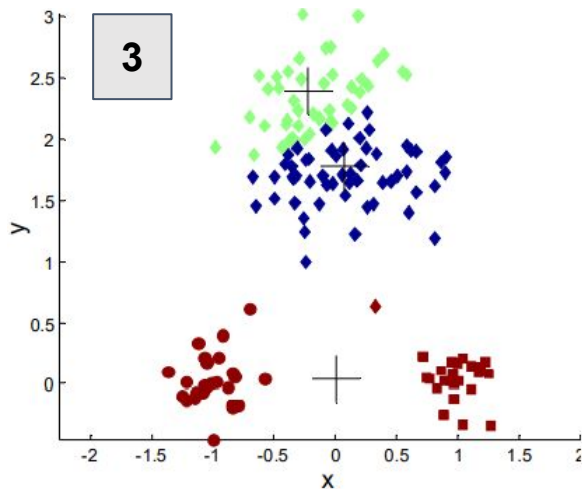
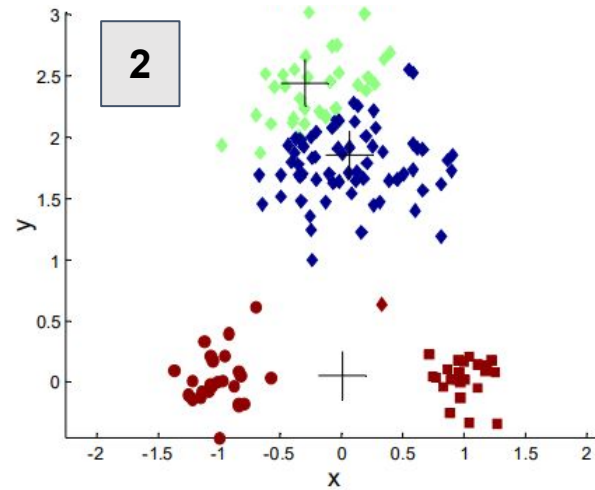
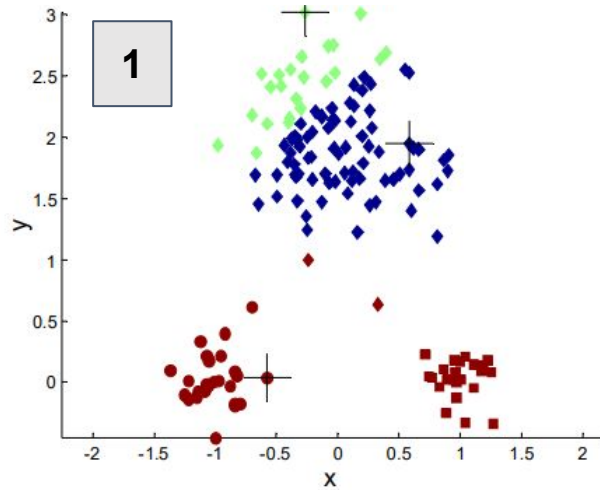




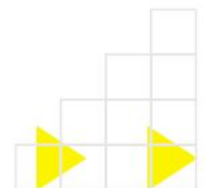
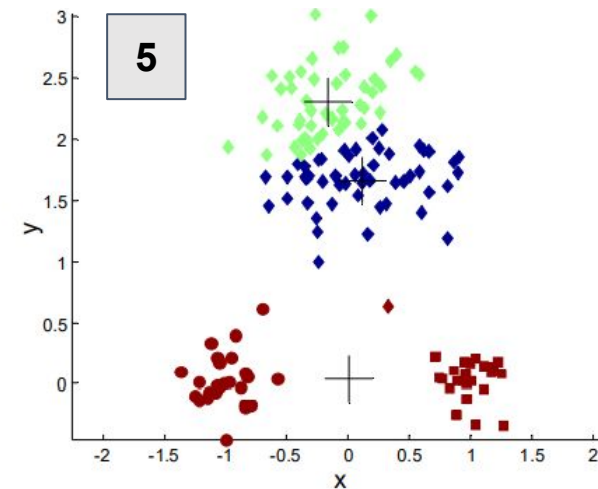
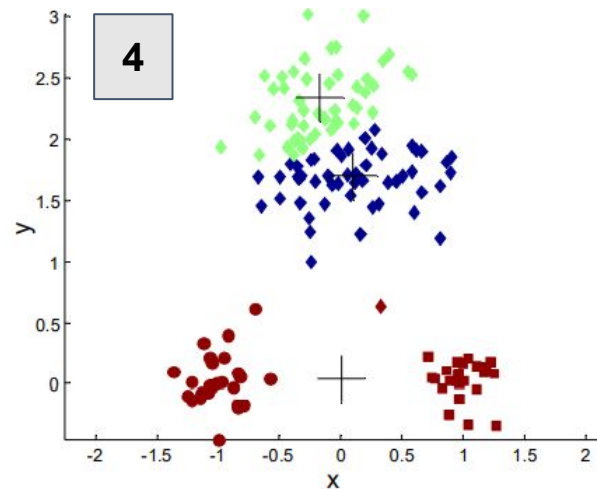
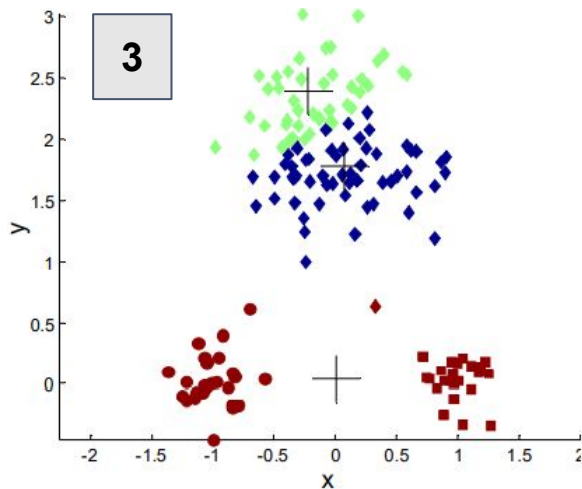
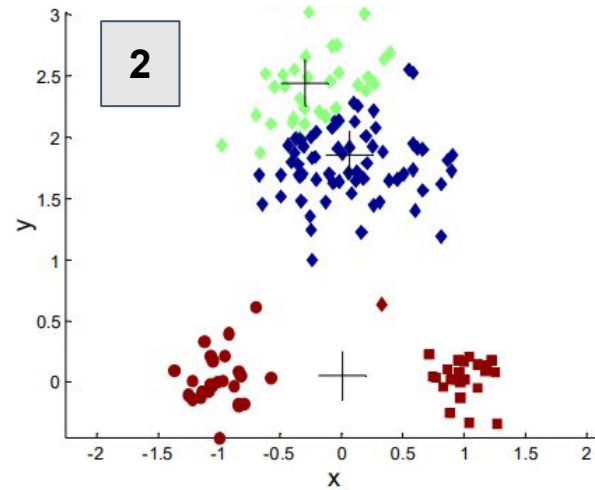
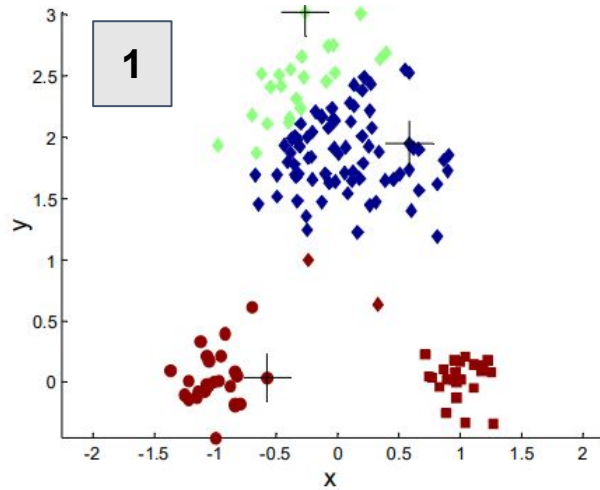
# Algoritmo *k-Means*



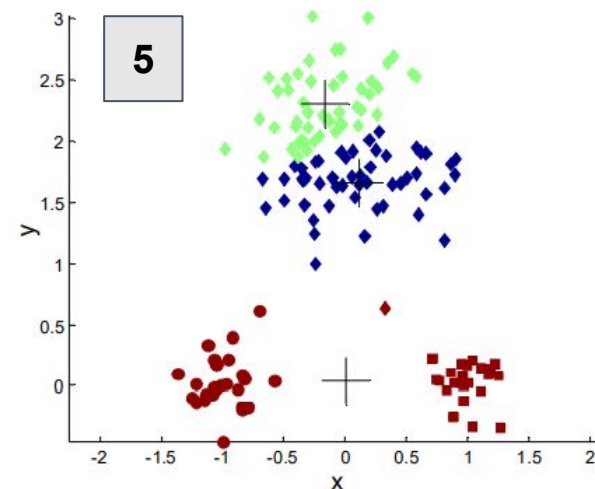
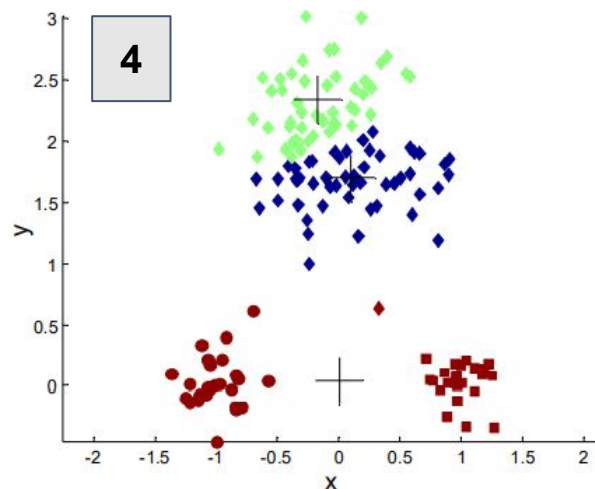
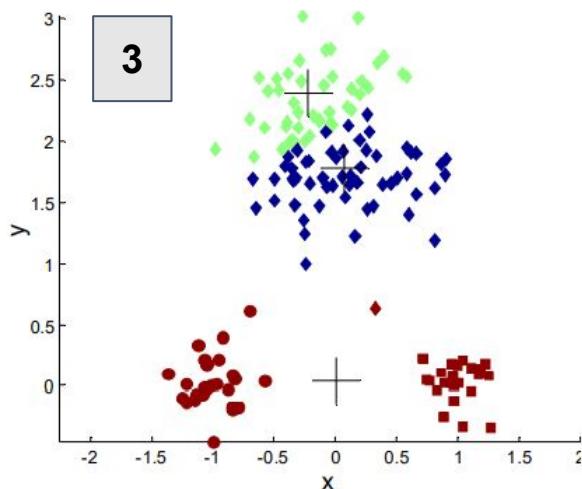
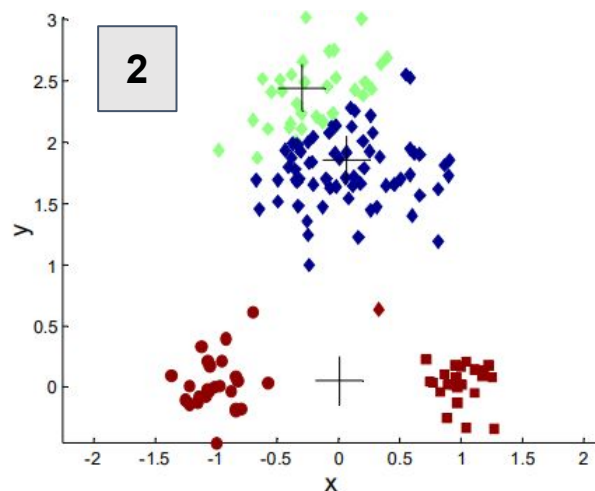
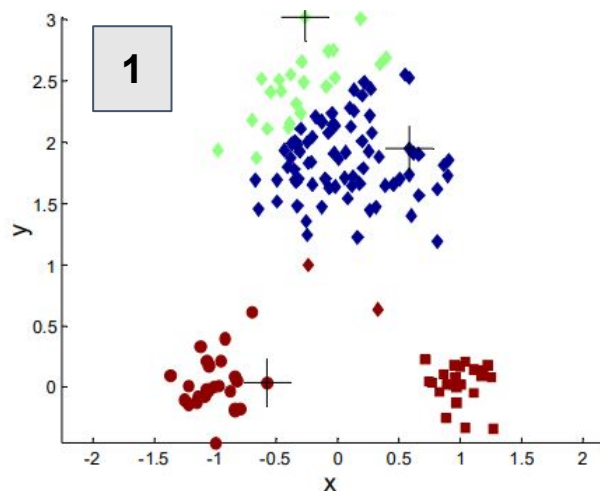
# Algoritmo *k-Means*



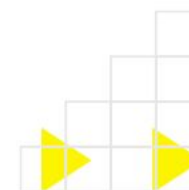
# Algoritmo *k-Means*



# Algoritmo *k-Means*



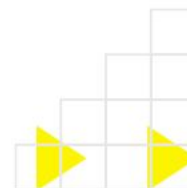
Nesta inicialização de centroides, o *k-means* obteve uma partição sub-ótima.



# Algoritmo *k-Means*

- Importância da escolha dos centroides iniciais
- Soluções comuns:
  - Múltiplas execuções do *k-means* e escolher a “melhor” solução de agrupamento (minimizar erro quadrático  $E$ )

$$E = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d^2(\mu_i, \mathbf{x})$$





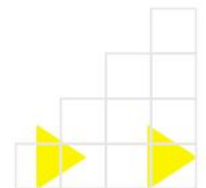
# Algoritmo *k-Means*



- Importância da escolha dos centroides iniciais
- Soluções comuns:
  - Múltiplas execuções do *k-means* e escolher a “melhor” solução de agrupamento (minimizar erro quadrático  $E$ )

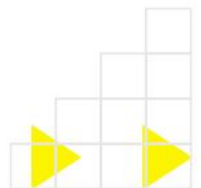
$$E = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d^2(\mu_i, \mathbf{x})$$

- Seleção “informada” dos centroides:
  - Garantir que sejam distantes entre si
  - Analista pode indicar centroides considerando sua experiência sobre o domínio dos dados



# Algoritmo *k-Means*

- Limitações do *k-Means*
  - *Outliers*
  - Clusters de tamanhos muito diferentes
  - Clusters de densidades muito diferentes
  - Clusters de formatos não globulares

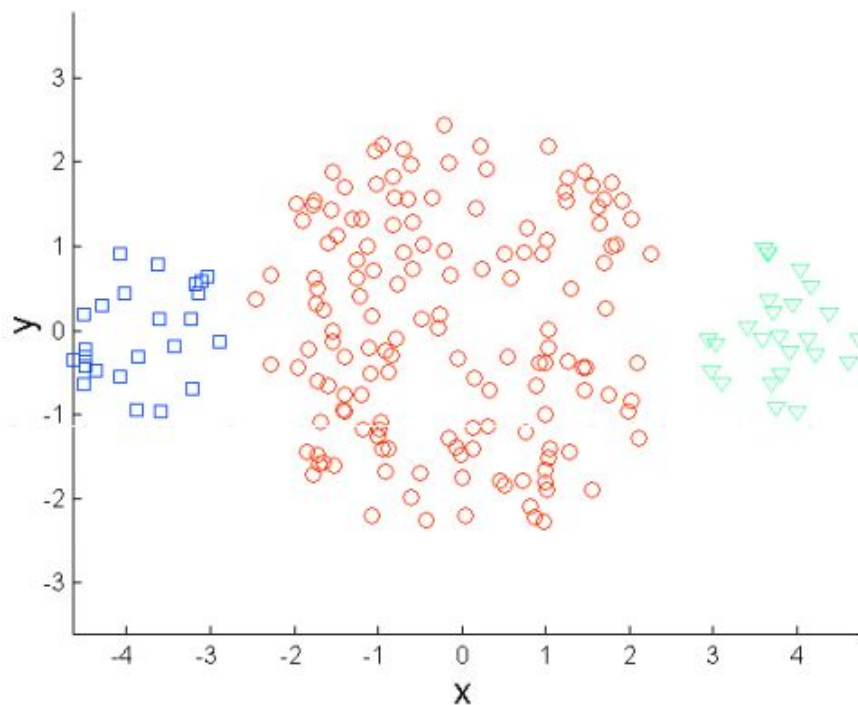




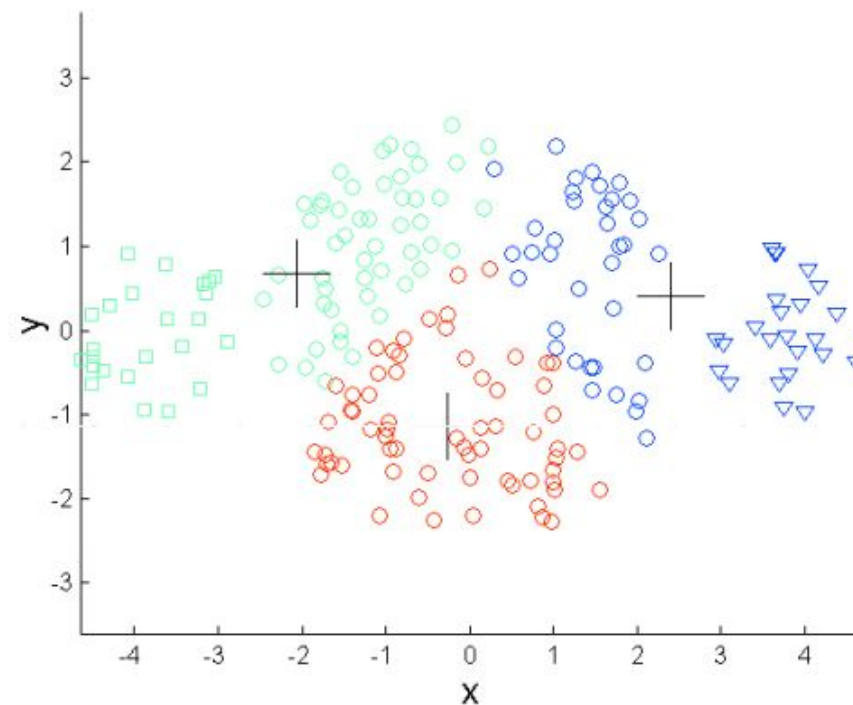
# Algoritmo *k-Means*



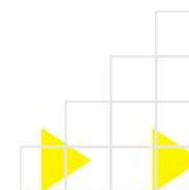
- Limitações do *k-Means*:
  - Clusters de tamanhos muito diferentes



Clusters esperados



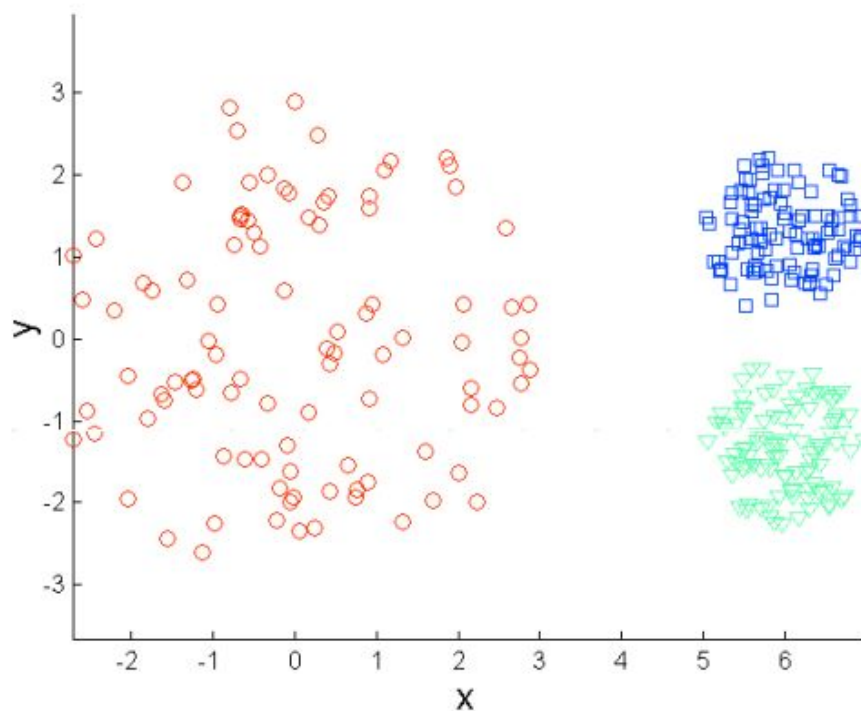
Clusters obtidos



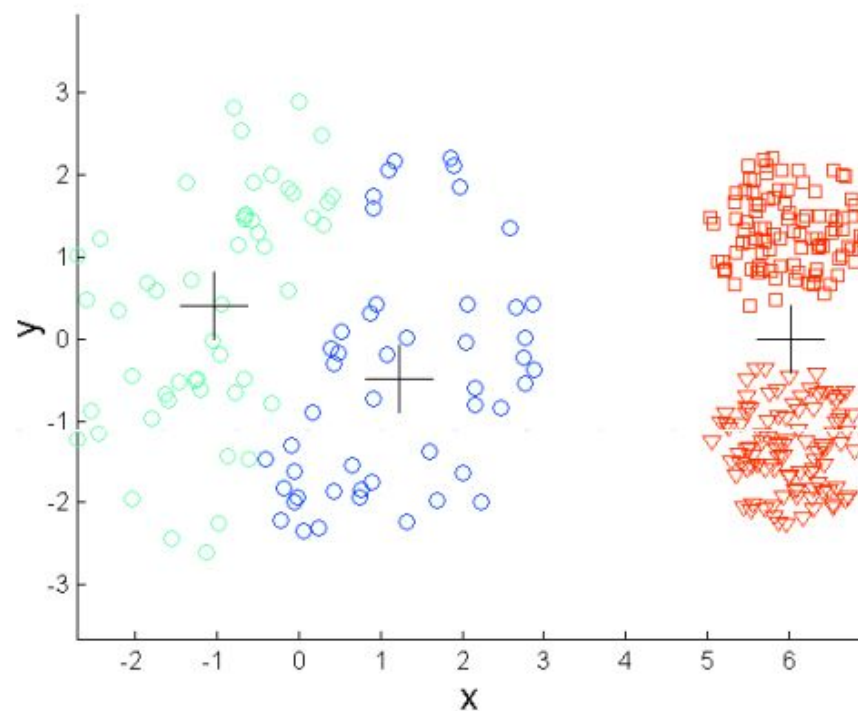
# Algoritmo *k-Means*



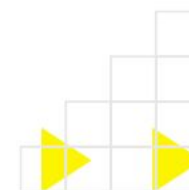
- Limitações do *k-Means*:
  - Clusters de densidades muito diferentes



Clusters esperados



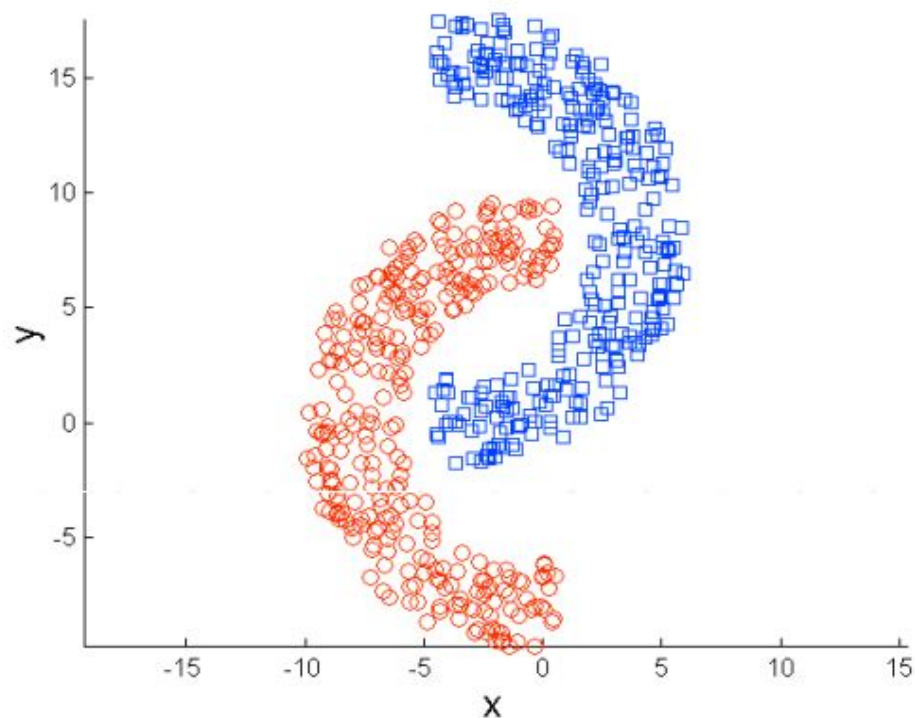
Clusters obtidos



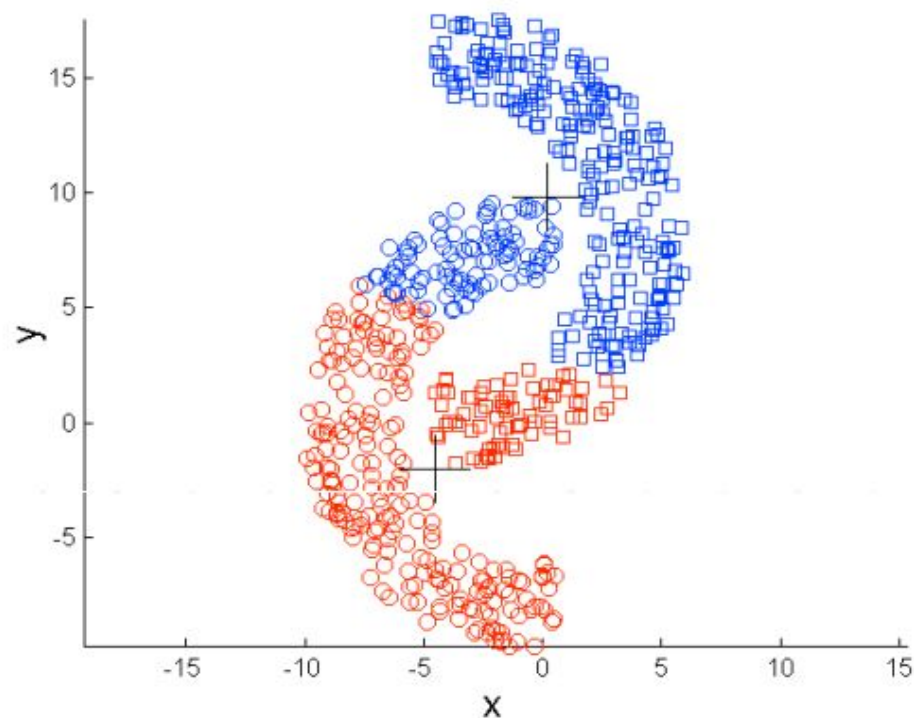
# Algoritmo *k-Means*



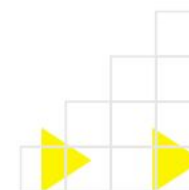
- Limitações do *k-Means*:
  - Clusters de formatos não globulares



Clusters esperados



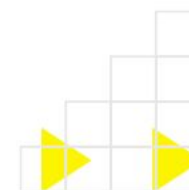
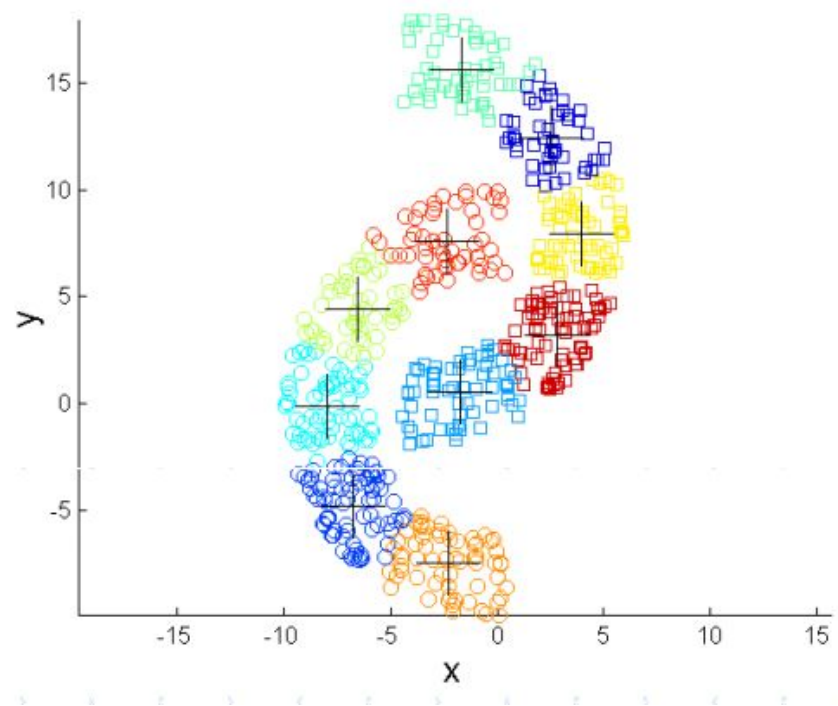
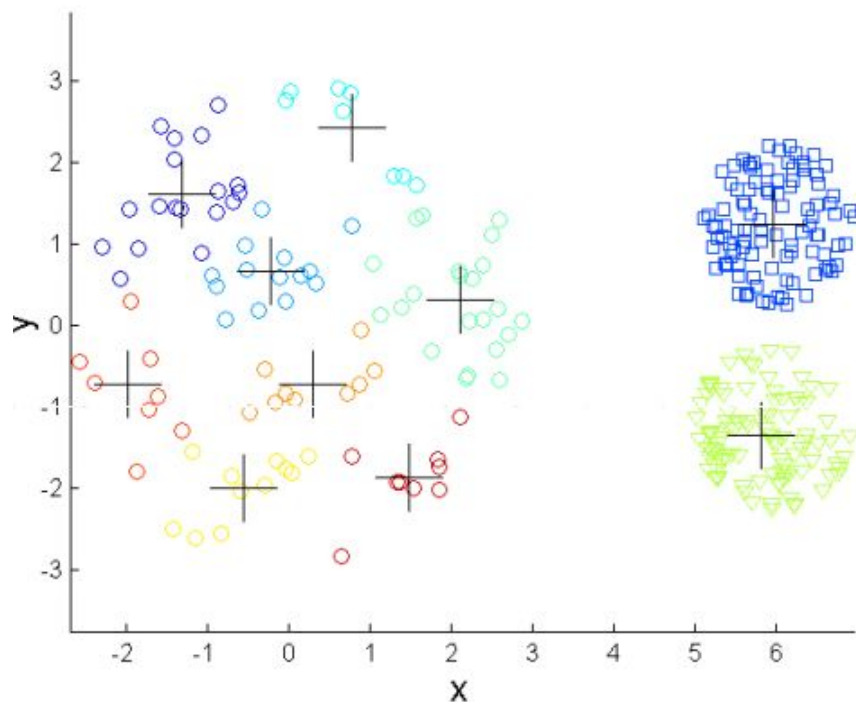
Clusters obtidos



# Algoritmo *k-Means*



- Mitigando as limitações do *k-Means*:
  - Podemos aumentar o número de *clusters*





# Algoritmos relacionados



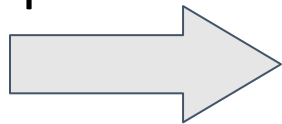
- **Algoritmo k-Medóides:**

- Similar ao k-Means
- São utilizados medoides no lugar de centroides
  - Medoides são objetos reais do conjunto de dados
  - Medoides representam objetos centrais do *cluster*
  - Não são calculados vetores médios
- Podemos usar apenas a matriz de dissimilaridades

Matriz atributo-valor

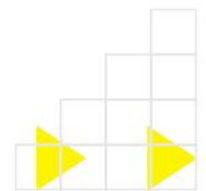
$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}$$

Aplicar medida  
de proximidade



Matriz dissimilaridades

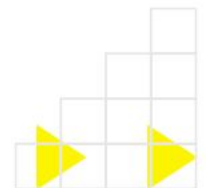
$$\mathbf{M} = \begin{bmatrix} d_{11} & d_{12} & \cdots & d_{1n} \\ d_{21} & d_{22} & \cdots & d_{2n} \\ \vdots & & \ddots & \vdots \\ d_{n1} & d_{n2} & \cdots & d_{nn} \end{bmatrix}$$



# Algoritmos relacionados



- **Algoritmo *bisecting K-Means*:**
  - Usa o *k-Means* para produzir um agrupamento hierárquico
  - Inicialmente, agrupar os dados com  $k=2$
  - Escolher o maior cluster e repetir o agrupamento com  $k=2$
  - Repetir até que obter a quantidade desejada de cluster



# Agrupamento Hierárquico



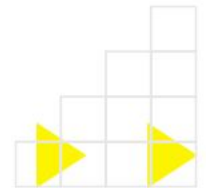
- Dois métodos clássicos para agrupamento hierárquico

## 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*

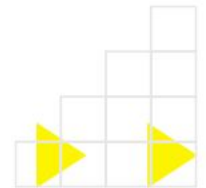
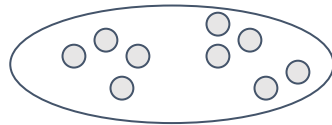




# Algoritmos relacionados



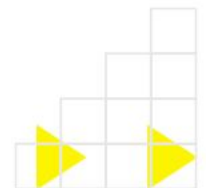
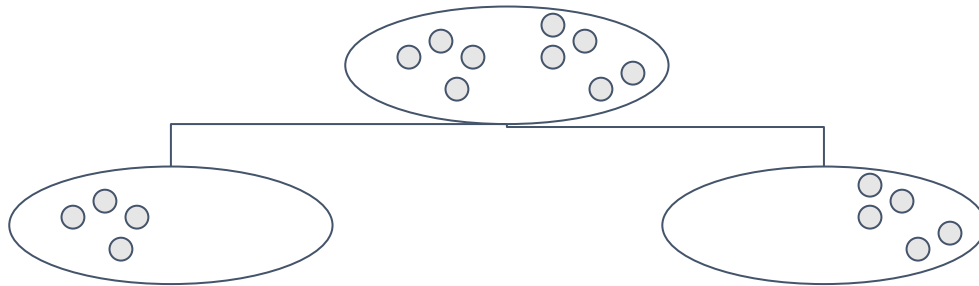
- **Algoritmo *bisecting K-Means*:**
  - Usa o *k-Means* para produzir um agrupamento hierárquico
  - Inicialmente, agrupar os dados com  $k=2$
  - Escolher o maior cluster e repetir o agrupamento com  $k=2$
  - Repetir até que obter a quantidade desejada de cluster



# Algoritmos relacionados



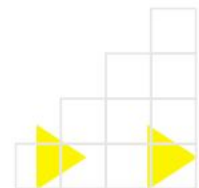
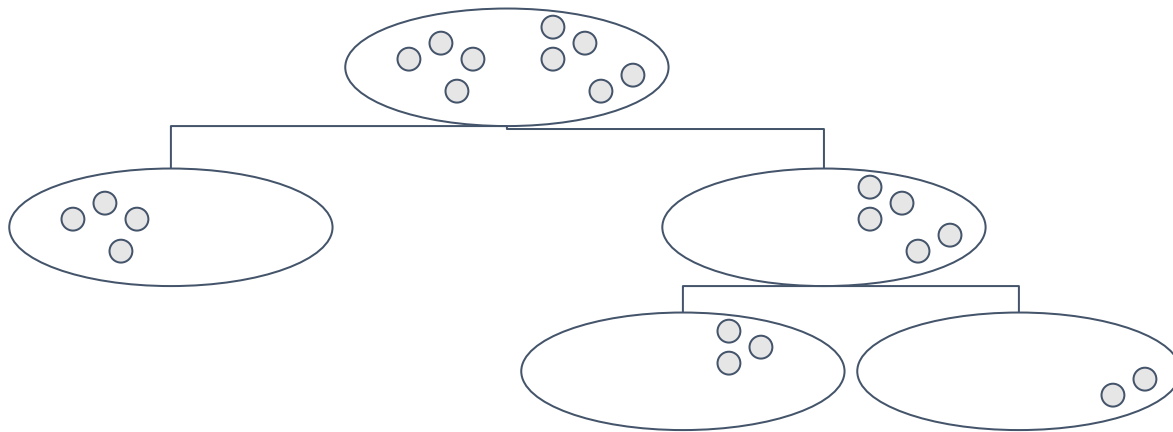
- **Algoritmo *bisecting K-Means*:**
  - Usa o *k-Means* para produzir um agrupamento hierárquico
  - Inicialmente, agrupar os dados com  $k=2$
  - Escolher o maior cluster e repetir o agrupamento com  $k=2$
  - Repetir até que obter a quantidade desejada de cluster



# Algoritmos relacionados



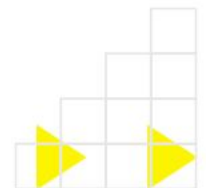
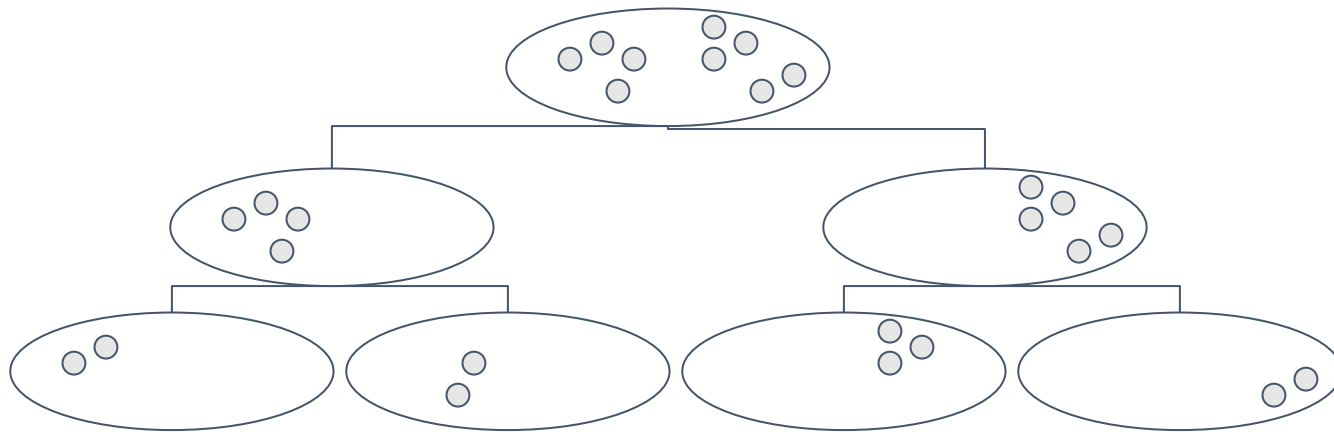
- **Algoritmo *bisecting K-Means*:**
  - Usa o *k-Means* para produzir um agrupamento hierárquico
  - Inicialmente, agrupar os dados com  $k=2$
  - Escolher o maior cluster e repetir o agrupamento com  $k=2$
  - Repetir até que obter a quantidade desejada de cluster



# Algoritmos relacionados



- **Algoritmo *bisecting K-Means*:**
  - Usa o *k-Means* para produzir um agrupamento hierárquico
  - Inicialmente, agrupar os dados com  $k=2$
  - Escolher o maior cluster e repetir o agrupamento com  $k=2$
  - Repetir até que obter a quantidade desejada de cluster

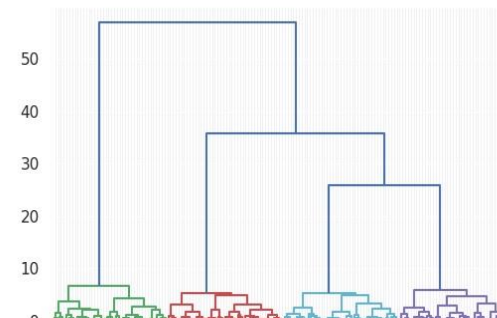
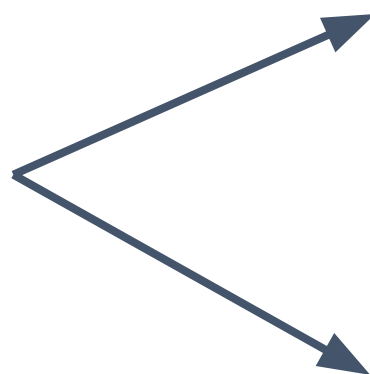
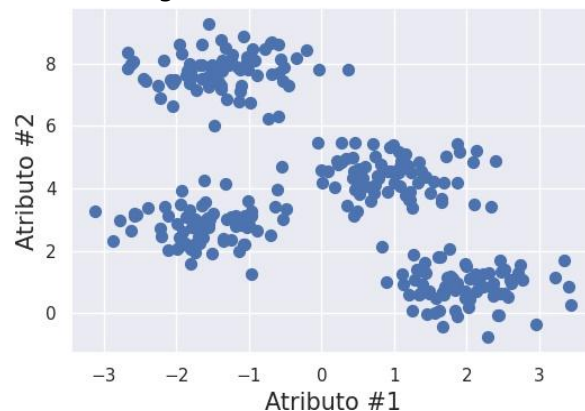


# Métodos para Agrupamento de Dados

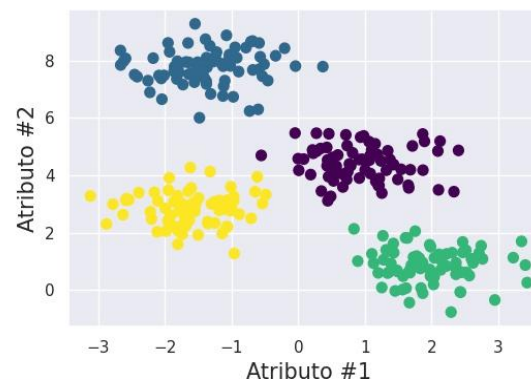


- Estudamos diferentes métodos e algoritmos
- Qual solução de agrupamento escolher?
- Qual o número apropriado de clusters para meus dados?

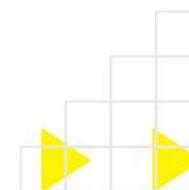
**Conjunto de Dados**



**Agrupamento Hierárquico**



**Agrupamento Particional**

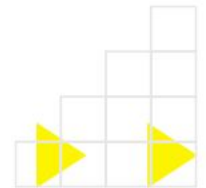


# Métodos para Agrupamento de Dados

- Estudamos diferentes métodos e algoritmos
- Qual solução de agrupamento escolher?
- Qual o número apropriado de clusters para meus dados?

## Próxima aula

Validação de Agrupamentos



# 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.

## Agradecimentos:

Notas de aula do curso de Análise de Agrupamentos, Prof. Eduardo Hruschka. Programa de Pós-Graduação do ICMC/USP. 2012.

