

CLUSTERING

K-MEANS

OUTUBRO 2024

Introdução

01

O que é K-means?

É um algoritmo de aprendizado não supervisionado utilizado para agrupar dados em diferentes clusters com base em similaridade

02

Importância

É amplamente utilizado em aplicações como segmentação de mercado, pois ajuda a identificar padrões ocultos em grandes conjuntos de dados

03

Como o funciona?

Atribui pontos de dados a clusters e ajustando as posições desses clusters até que a diferença entre os grupos seja minimizada

04

Quando usar?

K-means é eficaz quando temos grandes quantidades de dados não rotulados e queremos encontrar grupos naturais nos dados

Formula matemática

K-means é baseado no conceito de minimizar a soma das distâncias quadradas entre os pontos de dados e o centroide (o centro do cluster).

Fórmula da Função de Custo (Inércia)

$$J = \sum_{i=1}^k \sum_{x \in C_i} ||x - \mu_i||^2$$

J : A função de custo que queremos minimizar.

k : O número de clusters.

C_i : O i -ésimo cluster.

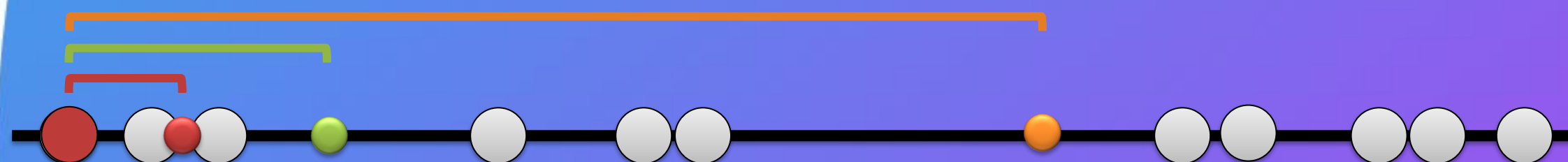
x : Um ponto de dado pertencente ao cluster C_i .

μ_i : O centroide do cluster C_i .

$||x - \mu_i||^2$: A distância euclidiana quadrada entre o ponto x e o centroide μ_i .

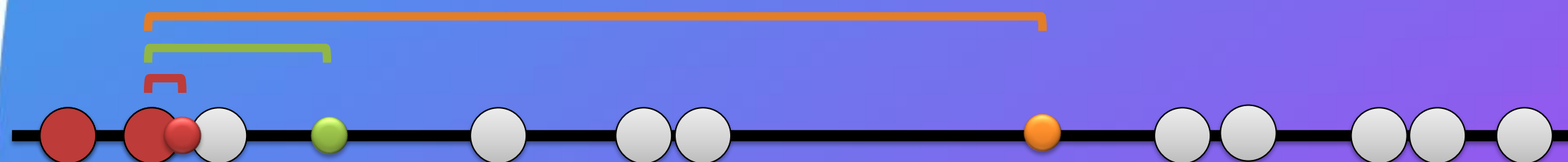
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



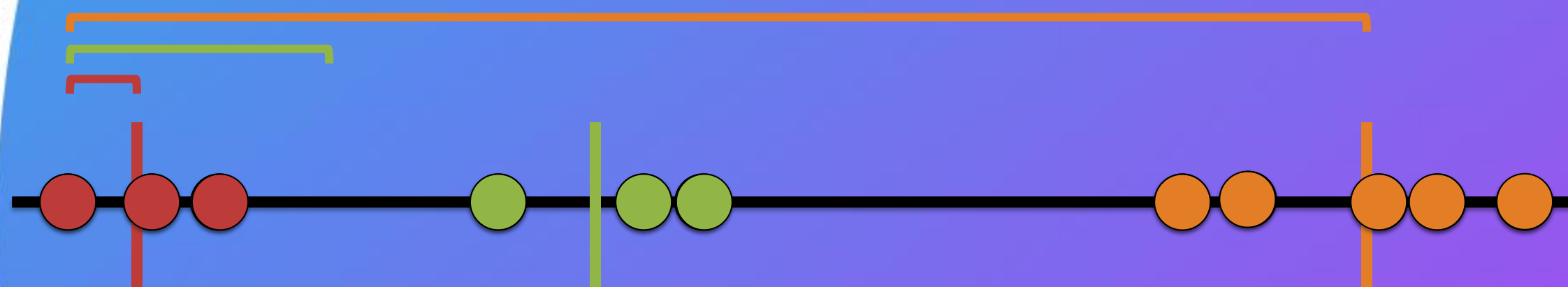
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



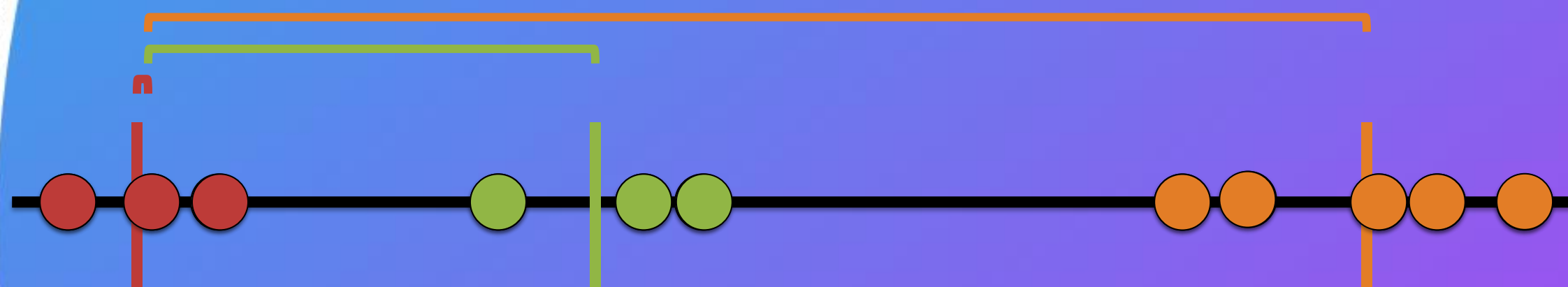
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



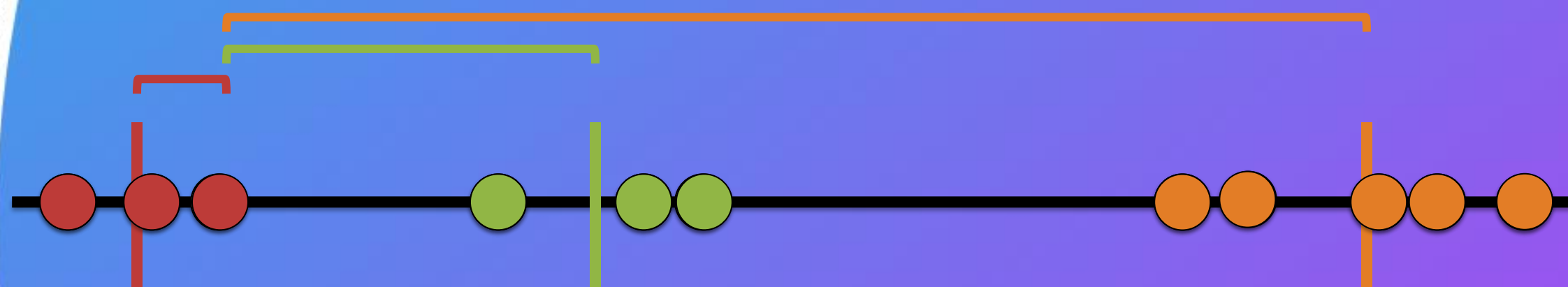
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



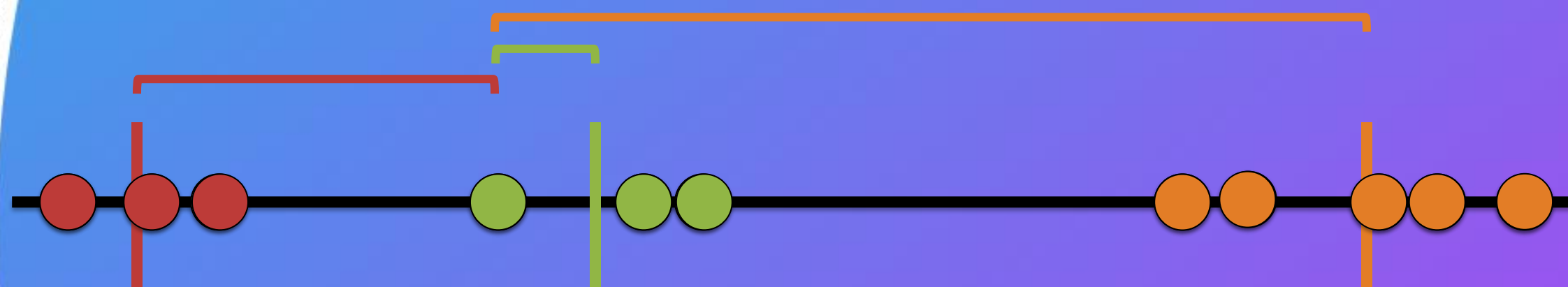
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



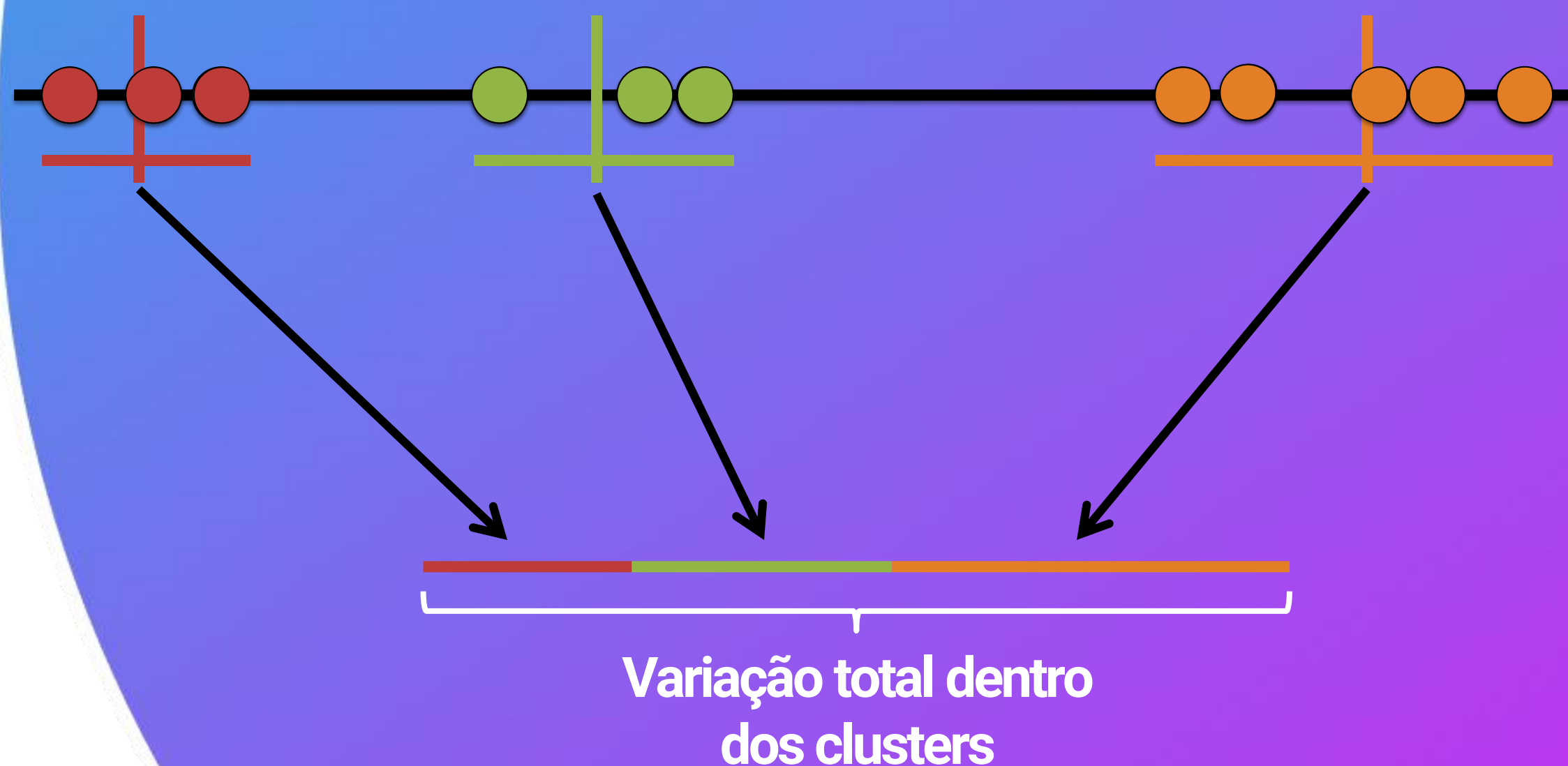
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centroide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centroide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



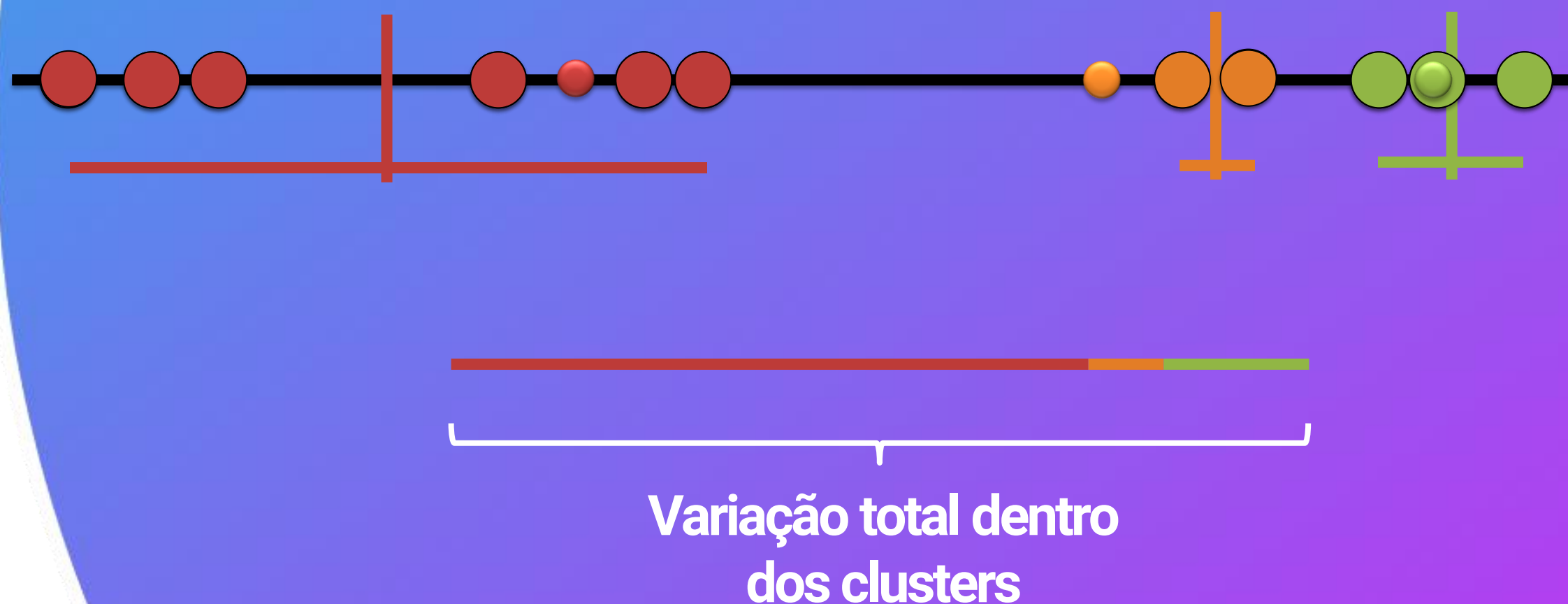
Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centróide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.



Funcionamento

O algoritmo K-means começa escolhendo k centroides aleatórios. Cada ponto de dado é atribuído ao centroide mais próximo. Em seguida, os centroides são recalculados como a média dos pontos de cada cluster. Os pontos são reatribuídos aos novos centroides e o processo se repete até que os centroides se estabilizem ou o número máximo de iterações seja atingido. O resultado final são k clusters com pontos de dados similares.

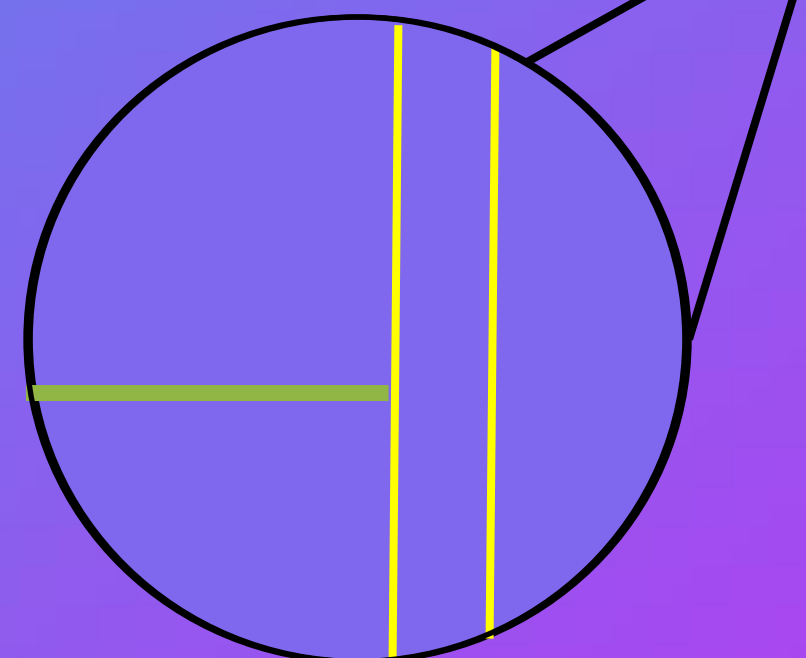
1° tentativa de cluster:



2° tentativa de cluster:

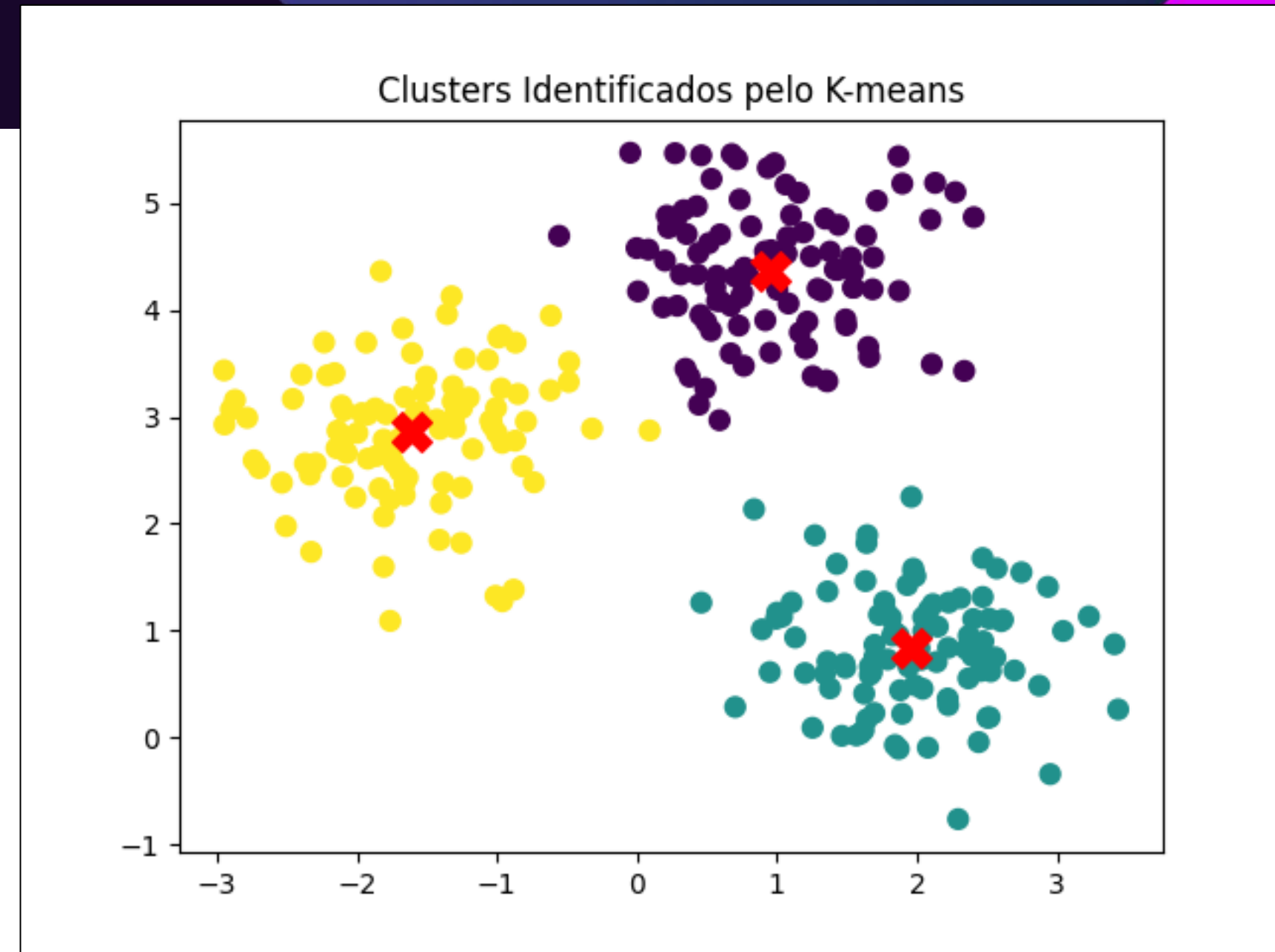


3° tentativa de cluster:

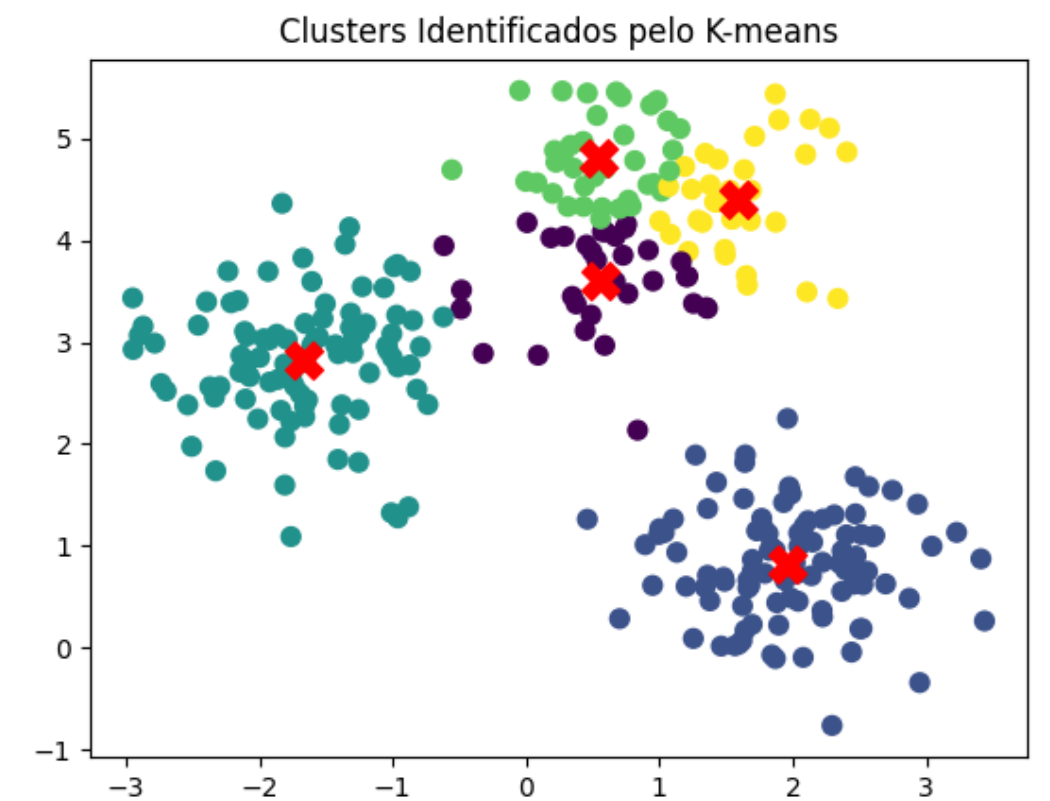
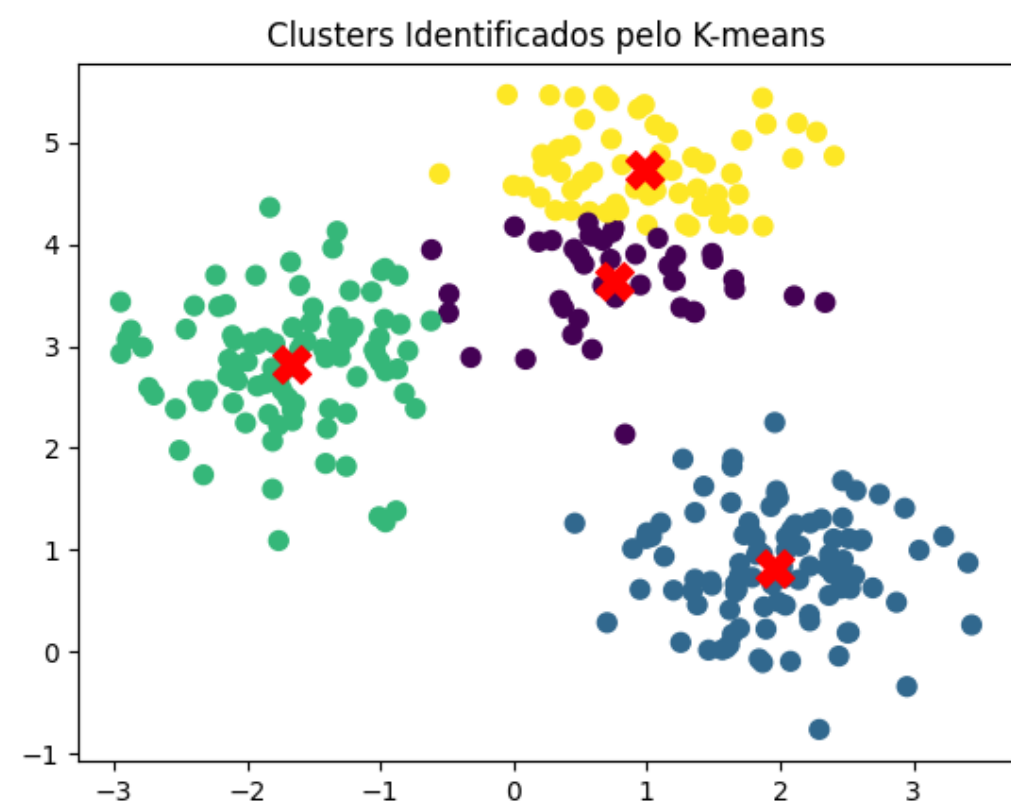
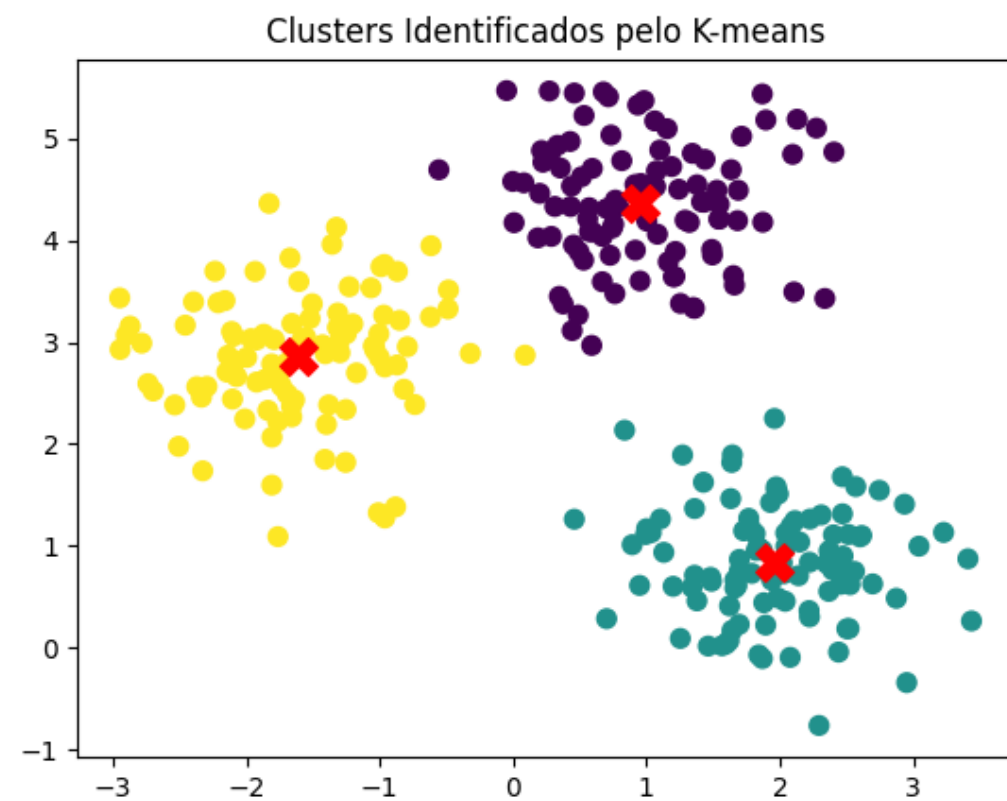


2º Período > PMC105A - Inteligência Artificial e Aprendizagem de Máquina > apresentação.py > ...

```
1  # Importando bibliotecas necessárias
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.cluster import KMeans
5  from sklearn.datasets import make_blobs
6
7  # Gerando dados fictícios
8  X, y = make_blobs(n_samples=300, centers=3, cluster_std=0.60, random_state=0)
9
10 # Visualizando os dados gerados
11 plt.scatter(X[:, 0], X[:, 1], s=50)
12 plt.title("Dados Gerados")
13 plt.show()
14
15 # Aplicando K-means com k = 3 (número de clusters)
16 kmeans = KMeans(n_clusters=3)
17 kmeans.fit(X)
18
19 # Obtendo os centroides e os rótulos (clusters) para cada ponto de dado
20 centroides = kmeans.cluster_centers_
21 rótulos = kmeans.labels_
22
23 # Visualizando os clusters com seus centroides
24 plt.scatter(X[:, 0], X[:, 1], c=rótulos, s=50, cmap='viridis')
25
26 # Plotando os centroides
27 plt.scatter(centroides[:, 0], centroides[:, 1], s=200, c='red', marker='x')
28 plt.title("Clusters Identificados pelo K-means")
29 plt.show()
```



Como saber o melhor valor de k ?

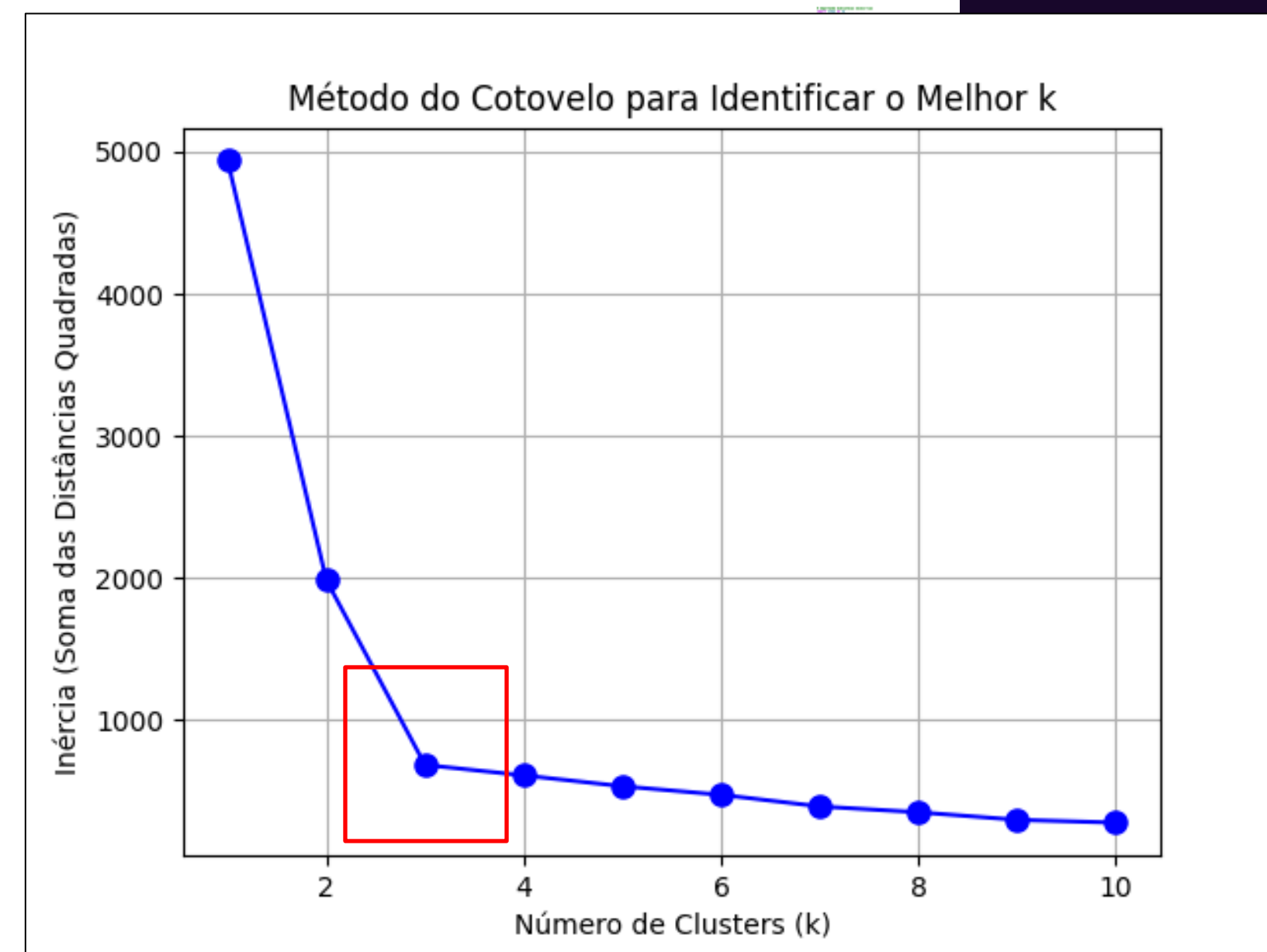


Como saber o melhor valor de k ?

Método do Cotovelo

2º Período > PMC105A - Inteligência Artificial e Aprendizagem de Máquina > apresentação.py > ...

```
1 # Importando bibliotecas necessárias
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 from sklearn.datasets import make_blobs
6
7 # Gerando dados fictícios (mesmo procedimento)
8 X, y = make_blobs(n_samples=1000, centers=3, cluster_std=0.60, random_state=0)
9
10 # Lista para armazenar os valores de inércia para diferentes k's
11 inercias = []
12
13 # Testando diferentes valores de k (de 1 a 10)
14 k_values = range(1, 11)
15 for k in k_values:
16     kmeans = KMeans(n_clusters=k, random_state=0)
17     kmeans.fit(X)
18     # Guardando o valor da inércia (soma das distâncias quadradas)
19     inercias.append(kmeans.inertia_)
20
21 # Plotando o gráfico do Método do Cotovelo
22 plt.plot(k_values, inercias, 'bo-', markersize=8)
23 plt.xlabel('Número de Clusters (k)')
24 plt.ylabel('Inércia (Soma das Distâncias Quadradas)')
25 plt.title('Método do Cotovelo para Identificar o Melhor k')
26 plt.grid(True)
27 plt.show()
```

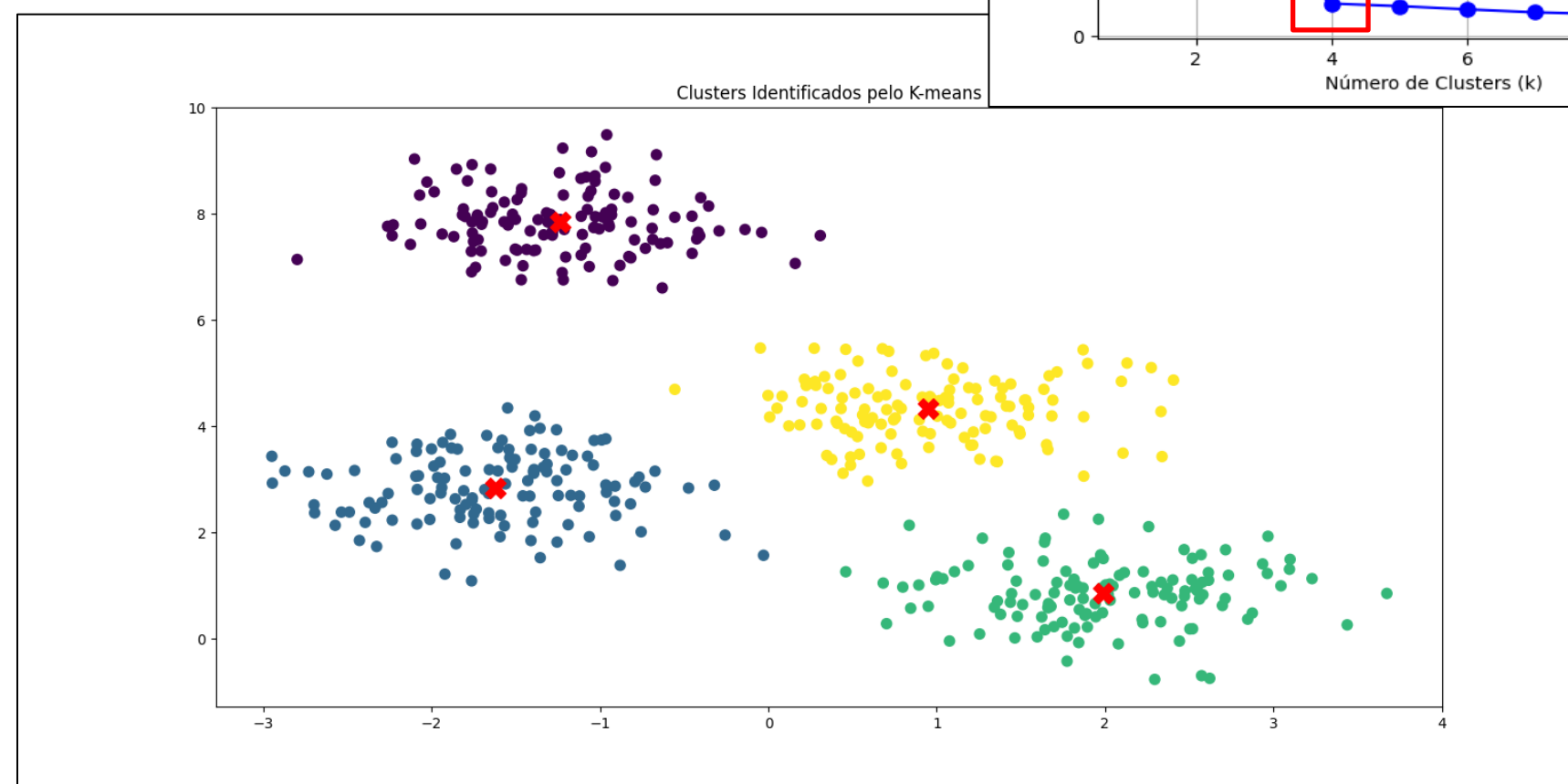
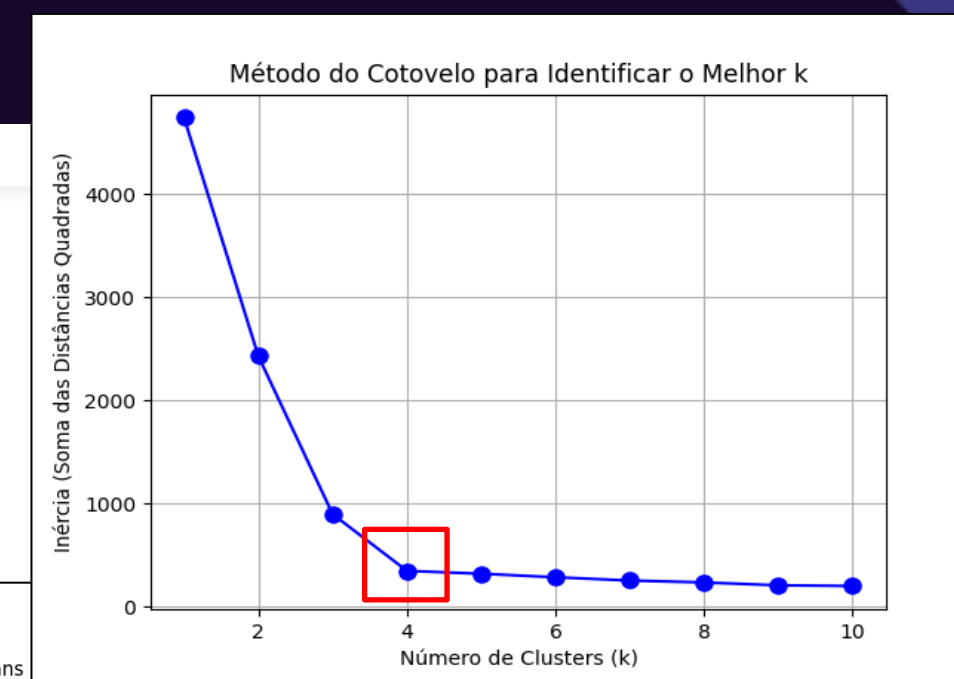


Como saber o melhor valor de k ?

Método do Cotovelo

2º Período > PMC105A - Inteligência Artificial e Aprendizagem de Máquina > k-means cotovelo.py > ...

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import KMeans
4 from sklearn.datasets import make_blobs
5
6 X, y = make_blobs(n_samples=500, centers=4, cluster_std=0.60, random_state=0)
7
8 inercias = []
9
10 # Testando diferentes valores de k (de 1 a 10)
11 k_values = range(1, 11)
12 for k in k_values:
13     kmeans = KMeans(n_clusters=k, random_state=0)
14     kmeans.fit(X)
15     # Guardando o valor da inércia (soma das distâncias quadradas)
16     inercias.append(kmeans.inertia_)
17
18 plt.plot(k_values, inercias, 'bo-', markersize=8)
19 plt.xlabel('Número de Clusters (k)')
20 plt.ylabel('Inércia (Soma das Distâncias Quadradas)')
21 plt.title('Método do Cotovelo para Identificar o Melhor k')
22 plt.grid(True)
23 plt.show()
```



stall.leonardo@pucpr.edu.br

OBRIGADO