

# Homework sobre agrupamento usando o algoritmo K-means

1) Usando a distância Euclidiana, realize o agrupamento dos 8 vetores bidimensionais abaixo a partir do algoritmo K-means. Com o mesmo, encontre  $K = 2$  centróides (médias) para representar estes dados e mostre estas centróides em um gráfico Cartesiano:

$X = \begin{bmatrix} -3, & 4 \\ 2, & 0 \\ 4, & -4 \\ 4, & 3 \\ 0, & 2 \\ 2, & 1 \\ 1, & -1 \\ -5, & 3 \end{bmatrix}$

1. Usando a distância Euclidiana, realize o agrupamento dos 8 vetores bidimensionais abaixo a partir do algoritmo K-means. Com o mesmo, encontre  $K=2$  centróides para representar estes dados e mostre estes centróides em um gráfico Cartesiano:

$$X = [-3, 4] \rightarrow P_1$$

$$2, 0 \rightarrow P_2$$

$$4, -4 \rightarrow P_3$$

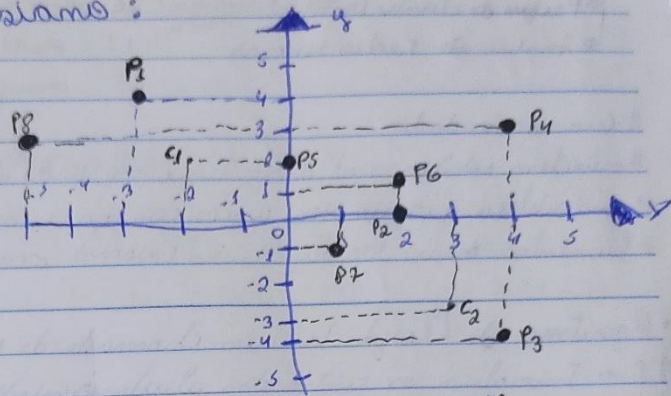
$$4, 3 \rightarrow P_4$$

$$0, 2 \rightarrow P_5$$

$$2, 1 \rightarrow P_6$$

$$1, -1 \rightarrow P_7$$

$$-5, 3 \rightarrow P_8$$



Como nas literaturas os centróides iniciais não são escolhidos aleatoriamente, vou escolher dois pontos aleatórios para representar os centróides.  $C_1(-2, 2)$ ;  $C_2(3, -3)$

# Agora faremos um agrupamento inicial para determinar quais pontos estão mais próximos dos centróides.

# Para  $C_1$

$$C_1 \rightarrow P_1 = 2.23 \rightarrow C_1$$

$$C_1 \rightarrow P_2 = 4.47$$

$$C_1 \rightarrow P_3 = 8.48$$

$$C_1 \rightarrow P_4 = 6.08$$

$$C_1 \rightarrow P_5 = 2.00 \rightarrow C_1$$

$$C_1 \rightarrow P_6 = 4.12 \rightarrow C_1$$

$$C_1 \rightarrow P_7 = 4.24$$

$$C_1 \rightarrow P_8 = 3.16 \rightarrow C_1$$

# Para  $C_2$

$$C_2 \rightarrow P_1 = 9.21$$

$$C_2 \rightarrow P_2 = 3.16 \rightarrow C_2$$

$$C_2 \rightarrow P_3 = 9.43 \rightarrow C_2$$

$$C_2 \rightarrow P_4 = 6.08 \rightarrow C_2$$

$$C_2 \rightarrow P_5 = 5.83$$

$$C_2 \rightarrow P_6 = 4.12$$

$$C_2 \rightarrow P_7 = 2.82 \rightarrow C_2$$

$$C_2 \rightarrow P_8 = 10$$

# Agrupado os dados, vamos calcular os novos centróides a partir da média dos pontos dentro de cada cluster.

Para C1

$(-3, 4) P1$

$(0, 2) P5$

$(2, 1) P6$

$(-5, 3) P8$

$$\text{média}_x = (-3 + 0 + 2 - 5) / 4$$

$$\text{média}_y = (4 + 2 + 1 + 3) / 4$$

$$\text{Novo centroide } C1 = (-1.5, 2.5)$$

Para C2

$(2, 0) P2$

$(4, -4) P3$

$(4, 3) P4$

$(1, -1) P7$

$$\text{média}_x = (2 + 4 + 4 + 1) / 4$$

$$\text{média}_y = (0 + (-4) + 3 - 1) / 4$$

$$\text{Novo } C2 = (2.75, -0.5)$$

# Agora temos que recalcular as distâncias com base nesses novos centroides. Para verificar se há mudança na proximidade dos pontos em relação a esses novos centroides.

# Para C1

$$C1 \rightarrow P1 = 2.12 \Rightarrow P1$$

$$C1 \rightarrow P2 = 4.3$$

$$C1 \rightarrow P3 = 8.51$$

$$C1 \rightarrow P4 = 5.52$$

$$C1 \rightarrow P5 = 3.58 \Rightarrow P5$$

$$C1 \rightarrow P6 = 3.80$$

$$C1 \rightarrow P7 = 4.30$$

$$C1 \rightarrow P8 = 3.93 \Rightarrow P8$$

# Para C2

$$C2 \rightarrow P1 = 7.3$$

$$C2 \rightarrow P2 = 0.9 \Rightarrow P2$$

$$C2 \rightarrow P3 = 3.71 \Rightarrow P3$$

$$C2 \rightarrow P4 = 3.71 \Rightarrow P4$$

$$C2 \rightarrow P5 = 3.71$$

$$C2 \rightarrow P6 = 1.67 \Rightarrow P6$$

$$C2 \rightarrow P7 = 1.82 \Rightarrow P7$$

$$C2 \rightarrow P8 = 8.5$$

# Como temos o acréscimo de P6 ao cluster de C2, temos que recalcular a média dos centroides

Para C1:

$P1(-3, 4)$

$P5(0, 2)$

$P8(-5, 3)$

$$\text{média}_x = (-3 + 0 - 5) / 3$$

$$\text{média}_y = (4 + 2 + 3) / 3$$

$$\text{Novo centroide } C1 = (-2.66, 3)$$

Para C2:

$P2(2, 0)$

$P3(4, -4)$

$P4(4, 3)$

$P6(2, 1)$

$P7(1, -1)$

$$\text{média}_x = (2 + 4 + 4 + 2 + 1) / 5$$

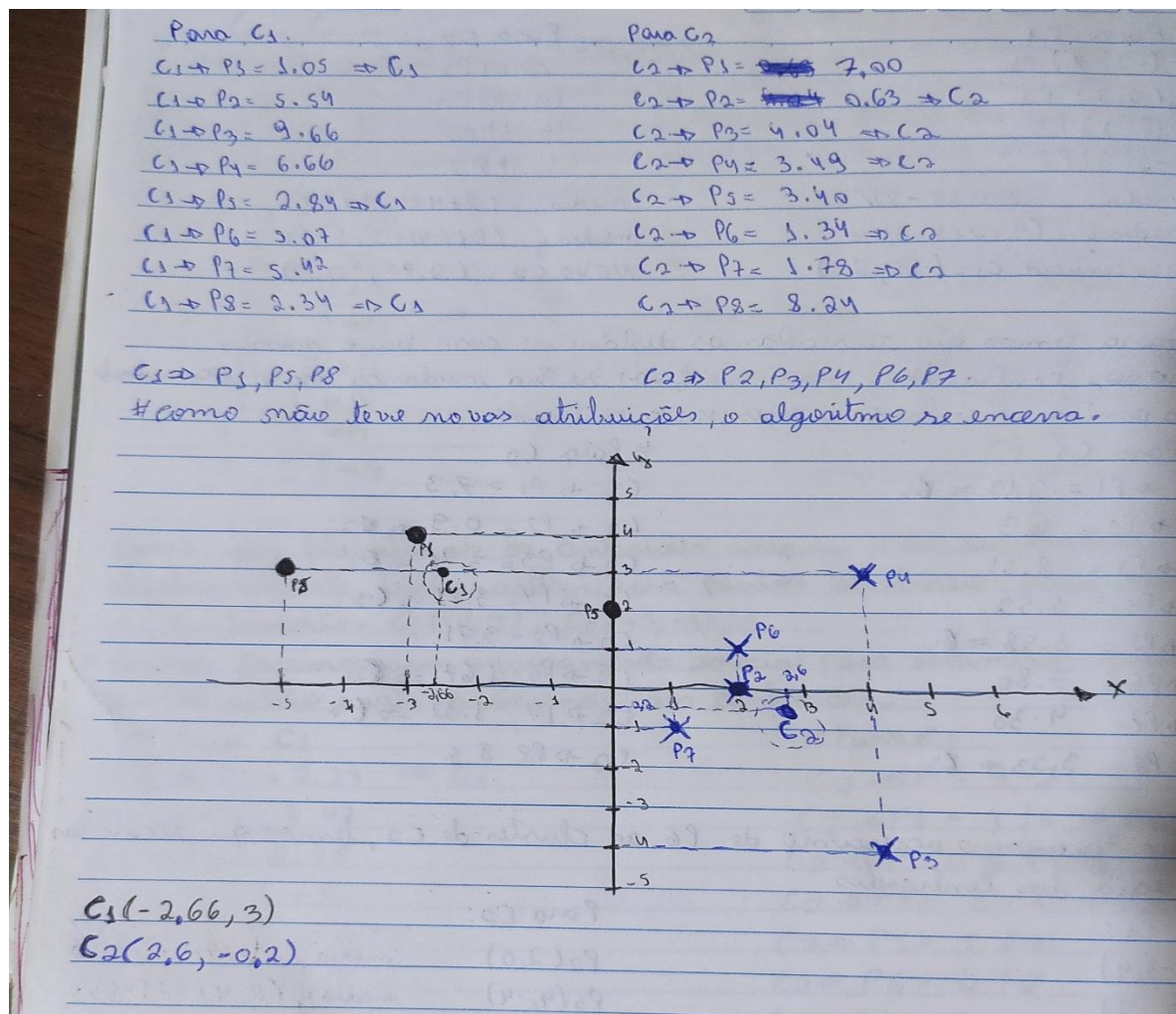
$$\text{média}_y = (0 - 4 + 3 + 1 - 1) / 5$$

$$C2 = (2.6, -0.2)$$

# Agora, temos que recalcular as distâncias com base nesses novos centroides para verificar se é adicionado algum novo ponto a algum cluster.

SÃO DOMINGOS





2) Um problema sério do uso do KNN na prática é quando o conjunto de treino é muito grande, e fazer a busca pelo vizinho mais próximo, ou mesmo armazenar o conjunto de treino em memória RAM pode ser proibitivo. Uma solução para isso é usar, antes do KNN, um algoritmo de agrupamento como o K-means ([https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)) para diminuir o conjunto de dados que o KNN irá manipular. Note que o KNN é um algoritmo de aprendizado supervisionado. Já algoritmos de agrupamento como o K-means são tipicamente “não-supervisionados” (não exigem rótulo para cada exemplo). Mas neste problema usaremos o K-means de forma independente para cada classe (rótulo) do conjunto de treino: executaremos um K-means apenas para os exemplos com rótulo “0” e outro K-means para os exemplos com rótulo “1”. Usando esses dois agrupamentos, converta seu conjunto de treino em um novo conjunto com apenas dois vetores (chamados “centróides”), um representando a classe “0” e outro a classe “1”. Depois o KNN deve ser usado com base neste conjunto de centróides para classificar o seu conjunto de teste.

Observe que o K-means é um algoritmo iterativo. Ele é implementado no Matlab / Octave como kmeans e no Scikit-learn: <https://scikit->

[learn.org/stable/modules/generated/sklearn.cluster.KMeans.html](http://learn.org/stable/modules/generated/sklearn.cluster.KMeans.html). Mas nós não precisaremos usar estas implementações pois assumiremos que a distância é Euclideana e queremos apenas 1 centróide. Daí esta centróide pode simplesmente ser obtida como a média aritmética dos vetores com os quais estamos lidando, e não precisamos de iterações extras do K-means.

Por exemplo, assuma que há 2 classes em um dado conjunto de treino com 6 exemplos no total:

```
1000,0,2,0
0,1,-4,0
2,2,2,0
-6000,-4,-1,1
0,-1,-4,1
-3000,-6,-1,1
```

O vetor de entrada  $x$  possui 3 elementos, e os organizando para as classes 0 e 1, e retirando a média aritmética (o que equivale a executar o K-means neste caso), encontra-se as centróides desejadas. No Octave isso pode ser feito com os comandos (note que não usamos o label "0" e "1" ao compor  $X_0$  e  $X_1$ ):

```
X0=[1000,0,2
    0,1,-4
    2,2,2]
X1=[-6000,-4,-1
    0,-1,-4
    -3000,-6,-1]
centroide_0 = mean(X0)
centroide_1 = mean(X1)
```

que gera:

```
centroide_0 = [ 334    1    0]
centroide_1 = [-3000.0000  -3.6667  -2.0000]
```

Daí você comporia o equivalente a um novo conjunto de treino para o KNN com estes dois vetores:

```
334, 1, 0, 0
-3000.0000, -3.6667, -2.0000, 1
```

e classificaria seu conjunto de teste usando KNN com  $K=1$  para estes dois vetores.

Observação: como você é curioso, caso queira aprender a usar o Kmeans no Octave (Matlab e Python é parecido), segue um exemplo abaixo (mas não é requerido o executar):

Carregue primeiro o pacote que contém o kmeans

pkg load statistics

e execute os comandos:

```
K=1;
```

```
[~,centroide_0] = kmeans (X0, K, 'distance','sqeuclidean')
```

```
[~,centroide_1] = kmeans (X1, K, 'distance','sqeuclidean')
```

2º) Dados o conjunto de treino e teste.

set - train		
X	Y	label
3000	0	0
3000	-3	0
-1000	2	0
1000	0	0
-6000	-2	1
0	0	1
-3000	-5	1
-2000	-6	1

set - test		
X	Y	label
-4000	-2	0
4000	3	0
0	0	1
-3000	-1	1
-5000	-6	1

# Normalizando os dados usando o min-max scaling

$$D = \frac{\text{dado} - \text{min}}{\text{max} - \text{min}}$$

set - train			set - test		
X	Y	label	X	Y	label
1	0.75	0	P1 0.22	0.5	0
1	0.375	0	P2 1.11	1.125	0
0.55	1	0	P3 0.66	0.75	1
0.77	0.75	0	P4 0.33	0.625	1
0	0.5	1	P5 0.11	0	1
0.66	0.75	1			
0.33	0.125	1			
0.44	0	1			

# 1º Passo dividir o conjunto de treino em label 0 e label 1 e encontrar seus respectivos Centroides.

# label 0

$$X_0 = [1, 0.75]$$

$$1, 0.375 \Rightarrow \text{Centroid}_0 = (0.83, 0.718)$$

$$0.55, 1$$

$$0.77, 0.75]$$

# label 1

$$X_1 [0, 0.5]$$

$$0.66, 0.75 \Rightarrow \text{Centroid}_1 = (0.3575, 0.625)$$

$$0.33, 0.125$$

$$0.44, 0]$$



D S T Q Q S S

$C_0(0.83, 0.718)$        $C_1(0.3575, 0.343)$

# utilizando o K-means, com  $K=2$ , para o teste distância euclidiana

# Para $C_0$	# Para $C_1$
$C_0 \rightarrow P_1 = 0.64$	$C_1 \rightarrow P_3 = 0.2 \rightarrow C_1$
$C_0 \rightarrow P_2 = 0.49 \rightarrow C_0$	$C_1 \rightarrow P_2 = 1.08$
$C_0 \rightarrow P_3 = 0.19 \rightarrow C_0$	$C_1 \rightarrow P_4 = 0.5$
$C_0 \rightarrow P_4 = 0.5$	$C_1 \rightarrow P_5 = 0.28 \rightarrow C_1$
$C_0 \rightarrow P_5 = 1.01$	$C_1 \rightarrow P_5 = 0.42 \rightarrow P_5$ fica no cluster $C_1$

# verificando os resultados

$P_1 \rightarrow 0$	# $C_0$ representa o cluster com classe 0
$P_2 \rightarrow 0$	# $C_1$ representa o cluster com classe 1.
$P_3 \rightarrow 1$	então
$P_4 \rightarrow 1$	$C_0 \rightarrow P_2, P_3 \rightarrow 1$ erro, 1 acerto
$P_5 \rightarrow 1$	$C_1 \rightarrow P_1, P_4, P_5 \rightarrow 1$ erro, 2 acertos

Taxa de erro: 40%  
 60% de acerto.

### 3) Quando o K-means foi inventado e quais foram suas primeiras aplicações? O que são K-medoides?

R=O termo "k-means" foi usado pela primeira vez por James MacQueen em 1967. O algoritmo padrão foi proposto pela primeira vez por Stuart Lloyd do Bell Labs em 1957 como uma técnica para modulação por código de pulso. O k-means clustering é um método de quantização vetorial, originário do processamento de sinais, que visa particionar  $n$  observações em  $k$  clusters em que cada observação pertence ao cluster com a média mais próxima (cluster centers ou cluster centroid), servindo como um protótipo de o aglomerado.

O que são medoides?

Medoids são objetos representativos de um conjunto de dados ou um cluster dentro de um conjunto de dados cuja soma das distâncias para outros objetos no cluster é mínima. Lembre-se que para K-Means, estávamos trabalhando com centroides. A relação entre centroides e medoides é semelhante à relação entre médias e medianas. Medoides e medianas sempre serão uma das observações nos dados, enquanto isso não é necessariamente o caso de centroides e médias.



O K-Means é sensível a valores discrepantes, mas o K-Medoids atenuou essa sensibilidade por não depender de centroides

4) Como as árvores K-d são usadas no scikit-learn para acelerar o agrupamento K-means?

R=

Ao usar a árvore para acelerar a etapa de atribuição, as informações resumidas coletadas nos permitem podar rapidamente centróides candidatos de grupos de pontos mantidos em nós internos. Quando um grupo de pontos tem apenas um candidato a centróide restante, determinamos que o candidato é o centróide mais próximo para todo o grupo e não precisa mais continuar percorrendo o subárvore. Isso elimina muitos cálculos de distância e comparações em comparação com outras implementações de k-means. Como a árvore é construída apenas uma vez, o algoritmo é melhor otimizado desde que a sobrecarga de tempo de construção da árvore seja menor do que o tempo economizado usando a árvore.