



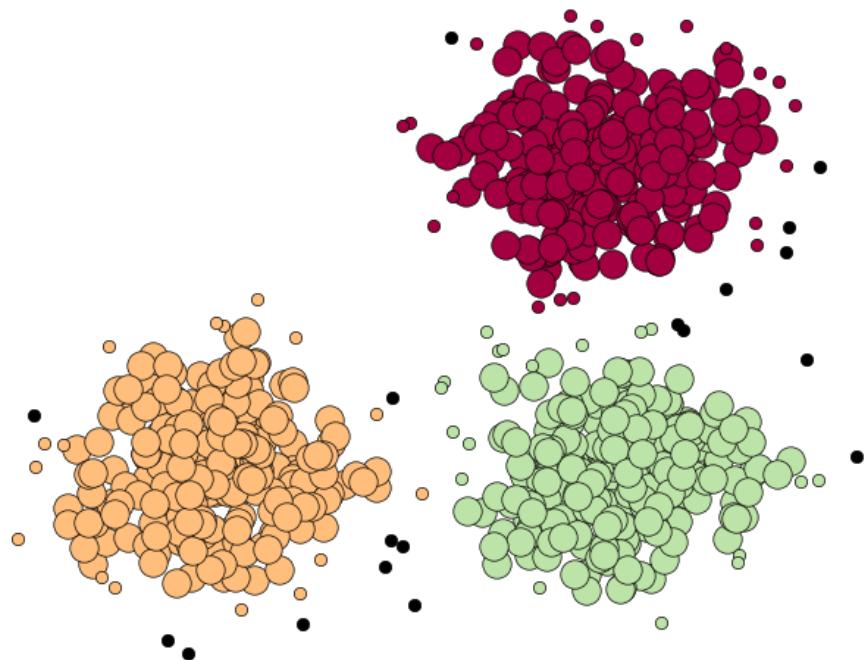
Universidade do Minho

4 de Abril de 2020

---

## O Método *Elbow*

Clusters particionais com dados numéricos



---

Bruno Jácome, A89515  
Carolina Barros, A84950  
Dinis Gomes, A87993  
Joana Gouveia, A85650

João Silva, A84617  
Jorge Gonçalves, A84133  
Pedro Peixoto, A89602

# Índice

<b>1</b>	<b>Introdução</b>	<b>6</b>
<b>2</b>	<b>Clusters</b>	<b>7</b>
2.1	O que são? . . . . .	7
2.2	Clustering . . . . .	7
2.2.1	Clustering na História . . . . .	8
2.3	Clusters Particionais . . . . .	8
2.3.1	O que são? . . . . .	8
<b>3</b>	<b>Centróides</b>	<b>9</b>
3.1	O que são? . . . . .	9
3.2	Relação entre centróide e média . . . . .	9
3.2.1	Semelhanças . . . . .	9
3.2.2	Diferenças . . . . .	9
3.2.3	Exemplo . . . . .	10
3.3	Como se determinam? . . . . .	10
<b>4</b>	<b>Os Algoritmos de Clusters Particionais</b>	<b>11</b>
4.1	Representação dos Dados . . . . .	11
4.2	K-Means . . . . .	11
4.2.1	O que é? . . . . .	11
4.2.2	Diagrama de Voronoi . . . . .	11
4.2.3	Restrições . . . . .	12
4.2.4	Determinação do K-Means . . . . .	12
4.3	K-Medoids . . . . .	13
4.3.1	Medóide . . . . .	13
4.3.2	O que é? . . . . .	13

4.4	Diferença entre K-Means e K-Medoids . . . . .	14
4.4.1	A nível de sensibilidade . . . . .	14
4.4.2	A nível de Centróide . . . . .	15
4.4.3	A nível atributos . . . . .	15
5	<b>O Método Elbow</b>	16
5.1	O que é? . . . . .	16
5.2	Pré-aplicação . . . . .	16
5.3	WCSS . . . . .	17
5.3.1	Ilusão de solução ótima . . . . .	17
5.3.2	Então, qual é a solução ótima ao problema? . . . . .	18
5.4	Aplicação . . . . .	18
6	<b>Demonstrações ilustrativas</b>	20
6.1	Exemplo de código: . . . . .	20
6.2	<i>K-means</i> com K = 2 . . . . .	21
6.3	<i>K-means</i> com K = 3 . . . . .	22
6.4	<i>The Elbow Method</i> . . . . .	23
7	<b>Conclusões</b>	24

# Listas de Ilustrações

2.1	Figura ilustrativa do agrupamento de clusters. . . . .	7
2.2	Clustering. . . . .	7
2.3	Cluster Particional . . . . .	8
4.2	<i>PAM</i> com <b>k = 3</b> . . . . .	14
5.1	<i>WCSS</i> mínimo . . . . .	17
5.2	<i>WCSS</i> máximo para um grande conjunto de <i>data points</i> . . . . .	18
5.3	Solução ótima: pequeno valor de <i>WCSS</i> e de número de clusters . . . . .	18
5.4	Valor de <i>wcss</i> <i>versus</i> número de <i>clusters</i> . . . . .	19
6.1	<i>k-means clustering</i> com <b>k = 2</b> . . . . .	21
6.2	<i>k-means clustering</i> com <b>k = 3</b> . . . . .	22
6.3	Procedimentos do <i>Elbow method</i> . . . . .	23

# Tabelas

# 1 | Introdução

Este trabalho foi realizado no âmbito da Unidade Curricular de Matemática das Coisas e tem como objetivo primordial o estudo do *Clusters* particionais com dados numéricos (centróide) através do *The Elbow Method*.

O presente relatório divide-se essencialmente em 4 partes. Primeiramente, no Capítulo 2, será feita uma contextualização do assunto, apresentan-se a definição de clusters no geral e, mais em concreto, de clusters particionais.

Seguidamente, no Capítulo 3, será descrito o conceito de centróides, bem como outros aspectos relevantes relativos.

Depois, no Capítulo 4, serão abordados alguns algoritmos de clusters particionais, com as apresentações das suas aplicações mais práticas.

No capítulo seguinte, a parte teórica será aplicada em exemplos mais práticos, de forma a melhor entendermos a aplicação dos tópicos referidos nos capítulos anteriores.

Para finalizar, expor-se-á uma breve conclusão do trabalho apontando-se os aspectos mais enriquecedores para o nosso conhecimento.

# 2 | Clusters

## 2.1 O que são?

Um **cluster** é um conjunto de objetos similares entre si e dissimilares em relação a objetos noutros clusters. A análise de clusters ou o seu conceito, é um procedimento humano normal, muitas vezes usado de forma inconsciente. [6][7]

Muito cedo nas escolas, os alunos aprendem a classificar e agrupar, por exemplo distinguir entre gatos e cães, entre animais e planta, progredindo num refinamento de classificação que tem subjacente teorias de *clustering*. A análise de clusters é usada em inúmeras aplicações, tais como no reconhecimento de padrões (*machine learning*), processamento de imagem e pesquisa de mercado.

## 2.2 Clustering

O *clustering* é o conjunto de técnicas de prospeção de dados, isto é, exames minuciosos e metódicos, que fazem agrupamentos automáticos de dados segundo o seu grau de semelhança. Normalmente o usuário do sistema deve escolher a priori o número de grupos a serem detetados. Alguns algoritmos mais sofisticados pedem apenas o nú-

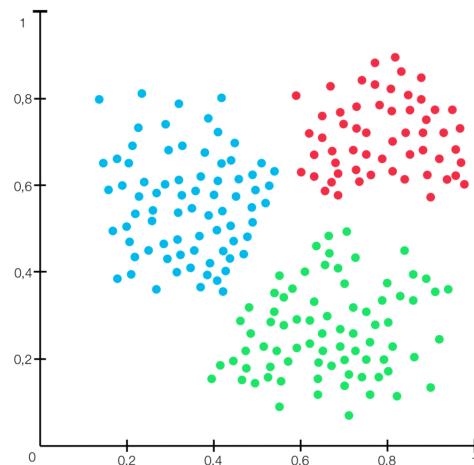


Figura 2.1: Figura ilustrativa do agrupamento de clusters.

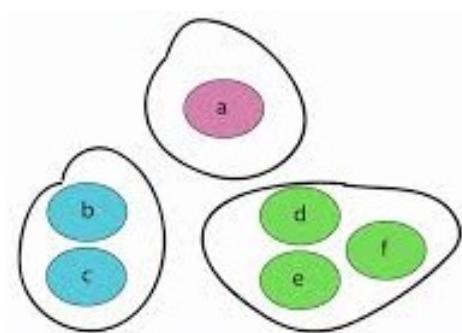


Figura 2.2: Clustering.

mero mínimo e outros tem a capacidade de subdividir um grupo em dois. Existem vários tipos de agrupamentos, mas o que será analisado com mais detalhe serão os **particionais**.

### 2.2.1 Clustering na História

O primeiro registo publicado sobre um método de clustering foi feito em 1948, com o trabalho de *SORENSEN* (1948) sobre o Método Hierárquico de Ligação Completa. Desde então mais de uma centena de algoritmos distintos de clustering já foram definidos.

## 2.3 Clusters Particionais

### 2.3.1 O que são?

Cluster particional define-se, especificamente, pelo facto de ao utilizar o agrupamento particional estar a dividir objetos de dados em subconjuntos sem sobrepor grupos, o que leva a que cada dado esteja exatamente num subconjunto.

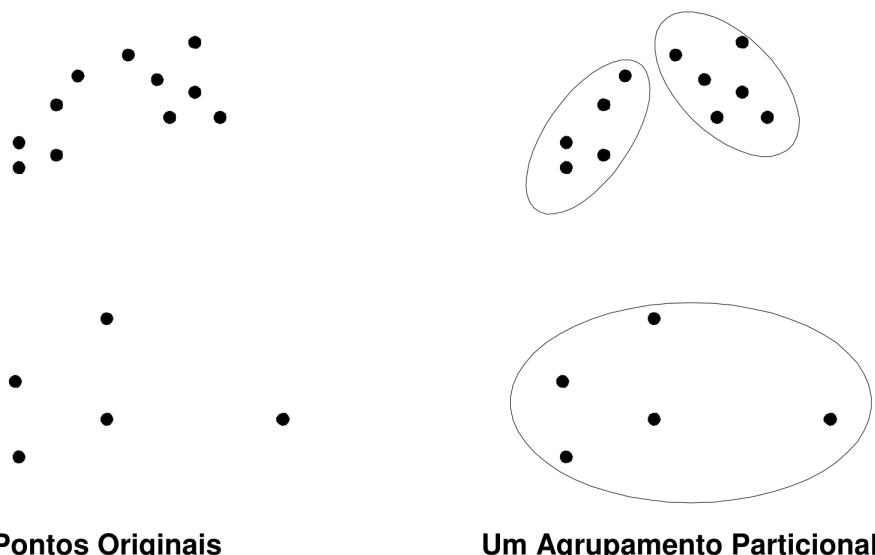


Figura 2.3: Cluster Particional

# 3 | Centróides

## 3.1 O que são?

Um **centróide** é o ponto que representa o *centro* de todos os pontos pertencentes a um cluster. No que diz respeito aos modelos centróides, a noção de similaridade deriva da proximidade dos pontos com o centróide do *cluster*.

Além disso, os centróides são obtidos através de operações algébricas (somas e multiplicações por escalares) e, em regra, estes não pertencem à base de dados. Logo, são uma mera interpretação de resultados e que dependem maioritariamente da definição de proximidade entre dois objetos de estudo.

## 3.2 Relação entre centróide e média

### 3.2.1 Semelhanças

... a *média* de um cluster é o mesmo que o centróide, contudo o termo **centróide** é mais preciso quando se estuda *multivariate data*, isto é, dados multivariados.

Um centróide é às vezes denominado de **centro de massa** ou **barycenter**(centro de gravidade), baseado na sua interpretação física. Assim como a média, a localização do centróide **minimiza a *sum-squared distance* entre os outros pontos**.

### 3.2.2 Diferenças

Há, no entanto, uma diferença entre **distância de centróide** e **distância média** quando se comparam clusters. A distância de centróide entre dois quaisquer clusters A e B é simplesmente a distância entre o centróide de A e o centróide de B. Já a distância média é calculada encontrando-se a distância média entre todos os pares de pontos de cada cluster.

Estes dois cálculos são duas métricas possíveis para calcular a distância entre dois clusters, mas existem mais métodos.<sup>[8]</sup>

$$\text{dist}(A, B) = \frac{\sum_{ij} \text{dist}(a_i, b_j)}{\#A \times \#B}, \quad \forall a_i \in A, b_j \in B$$

Métrica de Clusters: Distância média

$$\text{dist}(A, B) = \text{dist}\left(\frac{\sum_i a_i}{\#A}, \frac{\sum_i b_i}{\#B}\right), \quad \forall a_i \in A, b_i \in B$$

Métrica de Clusters: Distância entre centróides

### 3.2.3 Exemplo

## 3.3 Como se determinam?

# 4 | Os Algoritmos de Clusters Particionais

O Cluster Particional tem dois algoritmos: o K-means e o K-medoids. Estes algoritmos tem as suas diferenças. Uma delas é o facto de K-means temos a soma máxima das distâncias, em K-medoids temos a soma mínima nas distâncias.

## 4.1 Representação dos Dados

- K representa o número e clusters
- $(m^1, \dots, m^k, \dots, m^K)$  K pontos distintos de D
- $(x^1, \dots, x^n, \dots, x^N)$  representa a base de dados
- d representa uma função distância

## 4.2 K-Means

### 4.2.1 O que é?

O *K-Means* é um algoritmo de *clustering* bastante comum e popular usado por numerosos investigadores em todo o mundo. Este tem por objetivo por em partes  $n$  observações dentro de  $k$  clusters, onde cada observação está dentro do cluster com que tem a média mais próxima, usando o Diagrama de Voronoi. Nestes modelos, os números de clusters necessários no final ( $n$ ) têm de ser mencionados com antecedência, o que torna importante o conhecimento prévio do conjunto de dados.

### 4.2.2 Diagrama de Voronoi

O diagrama de Voronoi relaciona-se com o algoritmo K-means pelo facto de que o diagrama é uma parte do conjunto de dados com alguns pontos centrais, que se

denominam de centroides. Estes centroides não pertencem a base de dados, e um centroide é a localização (pode ser real ou imaginária) do centro de um cluster.

### 4.2.3 Restrições

Uma restrição que este algoritmo tem é o facto de apenas funcionar com atributos quantitativos, necessita de fazer operações algébricas, como somas e multiplicações por escalar, que dará origem uma matriz que é a “matriz da partição”. A nível de pontos que se encontram fora da curva, tem que se ter cuidado devido ao facto de os mesmos poderem facilmente influenciar o valor da média e levar a mesma a alterar-se.

### 4.2.4 Determinação do K-Means

**1º Passo:** Temos de determinar os clusters que estão associados a M, e para isso temos de ter cuidado com os algoritmos fora do grupo, para não se influenciar a mesma e para isso podemos utilizar um número mediano.

$$M = (m^1, \dots, m^k, \dots, m^K) \rightarrow P = (P^1, \dots, P^k, \dots, P^K)$$
$$P^k = x^i \in D : d(x^i, m^k) < d(x^i, m^j), j \in (1, \dots, k-1, k+1, \dots, K)$$

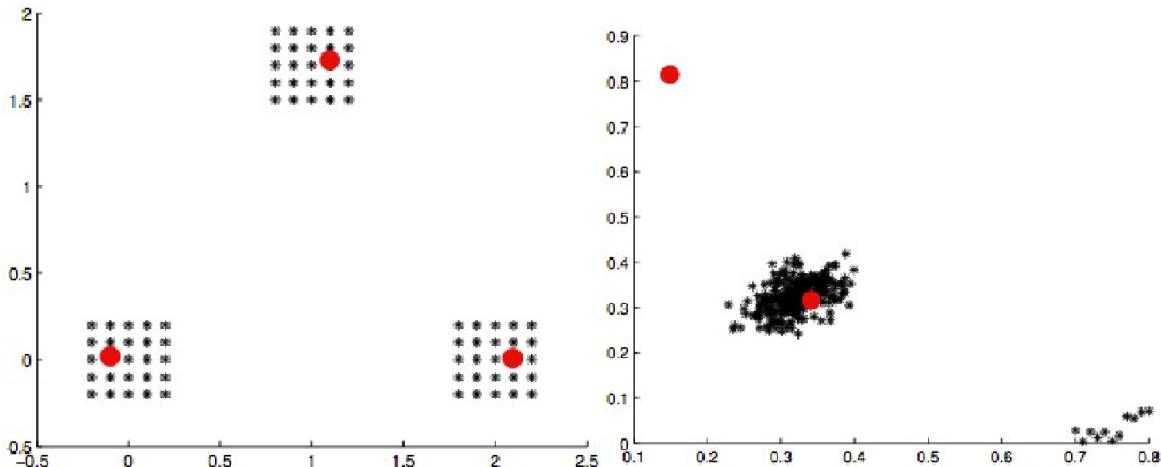
**2º Passo:** Temos de determinar os novos centroides associados a P, pegando no primeiro conjunto de dados e fazer uma seleção aleatória de k pontos de dados para se verificar onde se encontra o centro.

$$P = (p^1, \dots, p^k, \dots, p^K) \rightarrow M = (m^1, \dots, m^k, \dots, m^K)$$

$$m^k = \frac{1}{|P^k|} \sum_{x \in P^k} x$$

Fórmula do K-Means

**3º Passo:** Repetir os dois passos anteriores até se verificar que nenhum cluster muda de grupo.



(a) Gráfico bem resolvido por K-Means

(b) Gráfico mal resolvido por K-Means

## 4.3 K-Medoids

### 4.3.1 Medóide

Uma ideia semelhante é a centróide é a de **medóide**, que é o ponto de dado que é *menos parecido* de todos os outros pontos de dados.

Ao contrário do centróide, a medóide tem de ser um dos pontos originais.

### 4.3.2 O que é?

O *k-medoid* ou *partitioning around medoids* (**PAM**) são algoritmos de *clustering* reminiscentes do algoritmo de *k-means*, na medida em que ambos operam de modo particional e ambos tentam minimizar a distância entre os pontos e o centróide, dentro de um cluster.

Figura 4.2: *PAM* com  $k = 3$

## 4.4 Diferença entre K-Means e K-Medoids

### 4.4.1 A nível de sensibilidade

O *K-medoid* lida melhor com os *outliers* (pontos fora da curva) do que *K-means*, é menos sensível a eles, porque minimiza a soma das diferenças contrariamente a *k-means*, que maximiza.

#### **4.4.2 A nível de Centróide**

O centro de *k-medoids* não é o ponto médio mas sim um ponto real, porque é o objeto mais centralmente localizado do cluster, que como já referi, tem somas mínimas de distância.

#### **4.4.3 A nível atributos**

Os atributos de *K-medoids* podem ser atributos quantitativos, tal como *k-mean*, mas também podem ser atributos qualitativos, o que leva a que não exista uma necessidade e obrigação do uso de operações algébricas neste algoritmo. Estes atributos encontram-se representados na base de dados.

# 5 | O Método Elbow

## 5.1 O que é?

Uma etapa fundamental para qualquer aprendizagem não-supervisionada é determinar o número ideal de clusters segundo os quais os dados podem ser agrupados: **K**.

O **The Elbow Method** é uma heurística, uma vez que, é um método criado para encontrar soluções sobre um problema, neste caso, para determinar o número ideal de clusters no *k-means clustering*. Este método parcela o valor da função custo produzida pelos diferentes valores de **K**. Ora, isto só é possível, ignorando parte da informação com o objetivo de tornar a escolha mais fácil e rápida.

Sendo assim, não há uma resposta universal para este problema já que o número ideal de *clusters* é de alguma forma subjetivo e depende do método usado para medir as similaridades e os parâmetros usados para particionar. Portanto, em algumas situações, pode ser considerado ambíguo e pouco confiável. Nesse caso, é preferível utilizar-se outras abordagens para determinar o número de clusters.

## 5.2 Pré-aplicação

Numa fase inicial, criar um dendrograma, ou seja, um diagrama que organize as variáveis, agrupando-as de forma hierárquica ascendente - o que em termos gráficos se assemelha aos ramos de uma árvore.

Seguidamente, inspecionar o dendrograma produzido usando o cluster hierárquico para verificar se ele sugere um número específico de clusters. (Todavia, esta abordagem também é subjetiva.)

Estes métodos, apresentados a seguir, incluem métodos diretos e teste estatístico:

- **Métodos diretos:** consistem em otimizar um critério, como a somas de erros quadrados dentro do cluster (*With-in Cluster Sum of Squares*) ou a média si-

*lhhouette*. Os métodos correspondentes são denominados métodos de *Elbow* e *silhouette*, respetivamente.

- **Métodos de teste estatístico:** consiste em comparar evidências contra hipóteses nulas. Um exemplo é a estatística de gap.

É importante referir ainda que, a ideia básica por detrás dos métodos de particionamento, como o *k-means clustering*, é definir clusters de forma que a variação total intra-cluster, ou a soma total quadrada dentro do cluster (WSS), seja minimizada.

## 5.3 WCSS

### 5.3.1 Ilusão de solução ótima

Comumente, considera-se que para obter a solução ótima de número de clusters, deve-se obter o mínimo de **WCSS**.

Porque será então isto um erro?

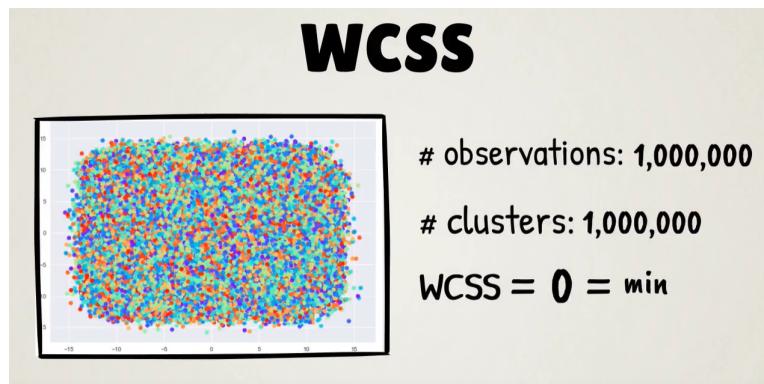


Figura 5.1: WCSS mínimo

Para um valor mínimo de *WCSS*, a solução ótima de número de cluster é igual ao número total de *data points* e a noção de cluster acaba por perder o seu propósito, acabando por ser uma solução trivial ao problema.

Por oposição,

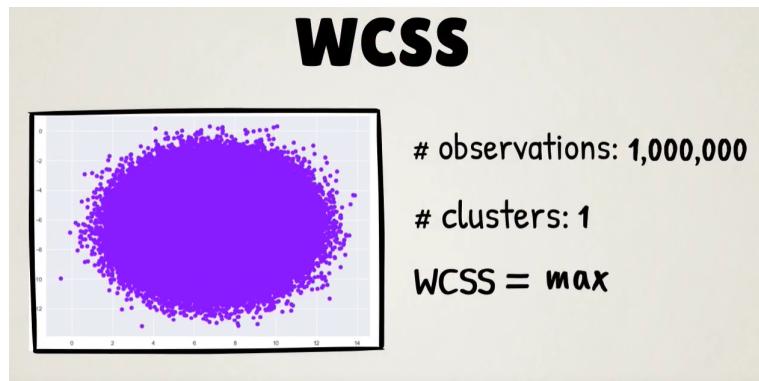


Figura 5.2: WCSS máximo para um grande conjunto de *data points*

Para um valor máximo de *WCSS*, a solução ótima seria apenas um cluster, isto é óbvio uma vez que a soma quadrada dentro dos clusters só poderia ser máxima se este contivesse todos os pontos.

### 5.3.2 Então, qual é a solução ótima ao problema?

Em boa verdade, não há. Como foi dito anteriormente, o número de clusters ótimo depende de vários fatores, inclusivé do objetivo de cada "clusterização"específica.

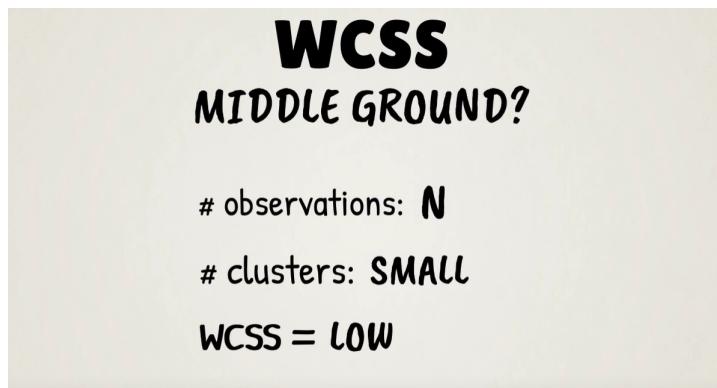


Figura 5.3: Solução ótima: pequeno valor de WCSS e de número de clusters

Contudo, com a ajuda do **método elbow**, é possível obter um resultado ótimo de **equilíbrio** entre o **número de cluster** e **wcss**.

## 5.4 Aplicação

O método de Elbow considera o WCSS total como uma função do número de clusters: deve-se escolher um número de clusters para que a adição de outro cluster não melhore muito mais o WCSS total.

O número ótimo de clusters pode ser obtido da seguinte forma:

1. Calcular o algoritmo de *clustering*, por exemplo, *k-means clustering*, para diferentes valores de  $k$ .
2. Para cada  $k$ , calcular o WCSS;
3. Representar, graficamente, o WCSS em função do  $k$ ;
4. Localizar, no gráfico, a curva com uma aparência de **cotovelo (elbow)**, geralmente considerado o indicador do número ótimo de *clusters*.

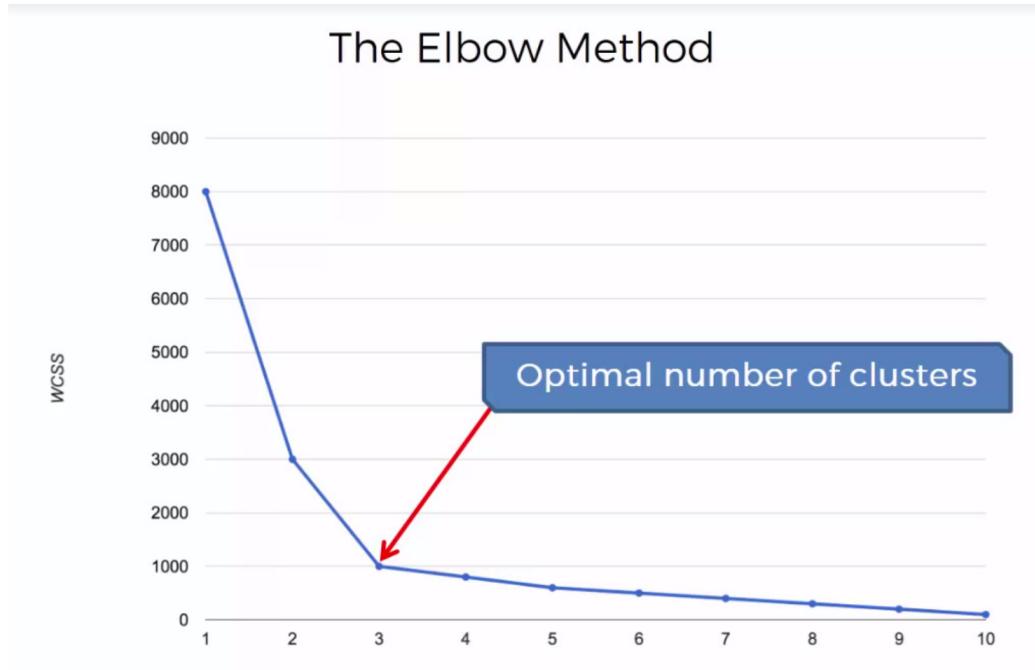


Figura 5.4: Valor de wcss *versus* número de *clusters*

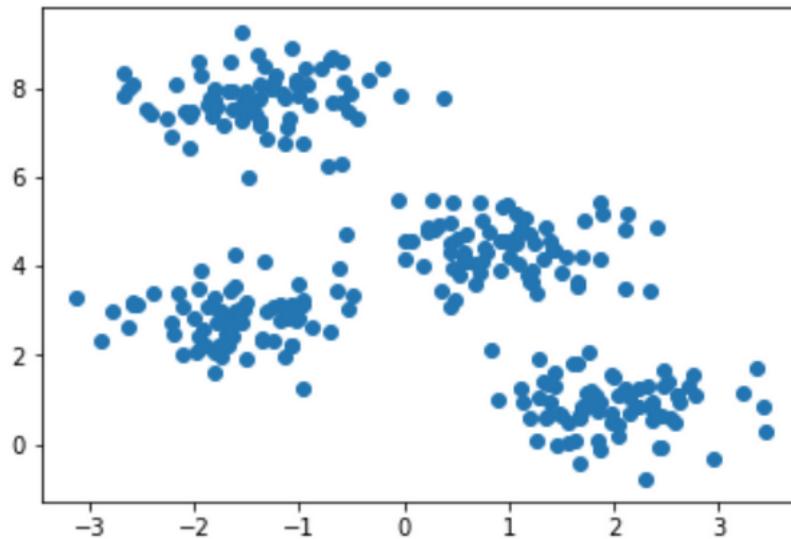
# 6 | Demonstrações ilustrativas

Mais fundamental que apenas expor a teoria por detrás dos conceitos abordados, será mesmo mostrar, de forma ilustrativa, o resultado da aplicação deste teoria. Então, com a intenção de facilitar a assimilação do estudo, seguem uma série de animações sobre alguns dos processos descritos.

## 6.1 Exemplo de código:

Para aplicar o algoritmo, precisamos de primeiro criar alguns conjuntos aleatórios de pontos e distribui-los com algum espaçamento.

```
points = make_blobs(n_samples=300, centers=4, cluster_std=0.60, random_state=0);
2 points.scatter(distance=1.5);
```



De seguida vamos aplicar kmeans aos nossos pontos. Vamos aplicar a função várias vezes, para numeros de clusters desde 1 até 9 e vamos guardar o valor de WCSS de cada resultado.

```
1 int wcss[10];
2
3 for(int i=1; i<10; i+=1) {
4     kmeans = points.KMeans(n_clusters=i, init="k-means++", max_iter=300, n_init=10,
5         random_state=0);
6     wcss[i] = kmeans.getWCSS();
7 }
```

Prestemos agora atenção à seguinte animação que pretende mostrar a aplicação do *k-means clustering* considerando **k = 2**.<sup>[9]</sup>

## 6.2 *K-means* com K = 2

Figura 6.1: *k-means clustering* com **k = 2**

### 6.3 *K-means* com $K = 3$

Analogamente, aplicando ao mesmo conjunto de *data points* agora com  $\mathbf{k = 3}^{[9]}$ :

Figura 6.2: *k-means clustering* com  $\mathbf{k = 3}$

Ora, o procedimento será idêntico para qualquer  $k$  maior que estes valores. No entanto, o  $k$  é escolhido aleatoriamente. O que aconterá se aplicar-mos o **método de Elbow** ?

## 6.4 *The Elbow Method*

Assim como na figura 5.4, uma representação gráfica e sua intrepertação, a partir de uma codificação (em *python*), do *elbow method*, utilizando o *k-means clustering* e o *wcss*.<sup>[10]</sup>

Figura 6.3: Procedimentos do *Elbow method*

# 7 | Conclusões

# Bibliografia

[1] *What is “Within cluster sum of squares by cluster” in K-means*

<https://discuss.analyticsvidhya.com/t/what-is-within-cluster-sum-of-squares-by-cluster-in-k-means/2706>

[2] *Elbow Method,*

<https://www.scikit-yb.org/en/latest/api/cluster/elbow.html>

[3] *Determining the optimal number of clusters,*

<https://www.datanovia.com/en/lessons/determining-the-optimal-number-of-clusters-3-must-know-methods/#elbow-method>

[4] *Finding the optimal number of clusters for K-Means through Elbow method using a mathematical approach compared to graphical approach,*

<https://www.linkedin.com/pulse/finding-optimal-number-clusters-k-means-through-elbow-asanka-perera>

[5] Lachi, Ricardo Luís & Rocha, Heloísa Vieira da. Fevereiro 2005. *Aspectos básicos de clustering: conceitos e técnicas*. (Brasil).

[6] [https://www.maxwell.vrac.puc-rio.br/24787/24787\\_5.PDF](https://www.maxwell.vrac.puc-rio.br/24787/24787_5.PDF)

[7] <http://www.dei.isep.ipp.pt/~paf/proj/Julho2003/Clustering.pdf>

[8] *Hierarchical Clustering 3: single-link vs. complete-link*

<https://www.youtube.com/watch?v=VMyXc3SiEqs>

[9] *K Means Clustering: Pros and Cons of K Means Clustering*

<https://www.youtube.com/watch?v=YIGta1P1mv0>

[10] *How to Choose the Number of Clusters / Advanced Statistical Methods - K-Means Clustering*

<https://www.youtube.com/watch?v=SCA07-7Xe6Q>