

K-Means Clustering

Dr. Junya Michanan

What is K-Means Clustering

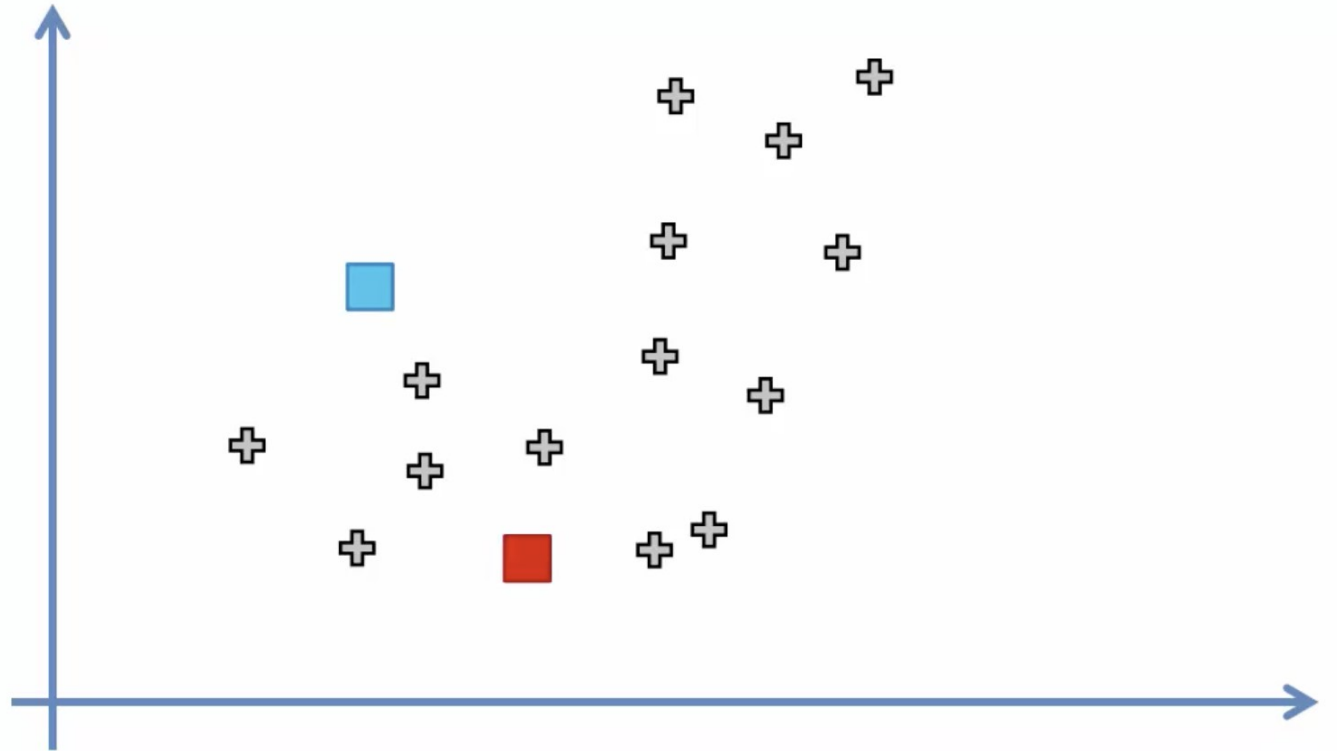
- K-Means Clustering is an unsupervised machine learning algorithm.
- Is an UN-supervised machine learning algorithm.
- K-Means attempts to classify data without having first been trained with labeled data.
- Once the algorithm has been run and the groups are defined, any new data can be easily assigned to the most relevant group.

K-Means Clustering Applications

- customer profiling
- market segmentation
- computer vision
- search engines
- astronomy

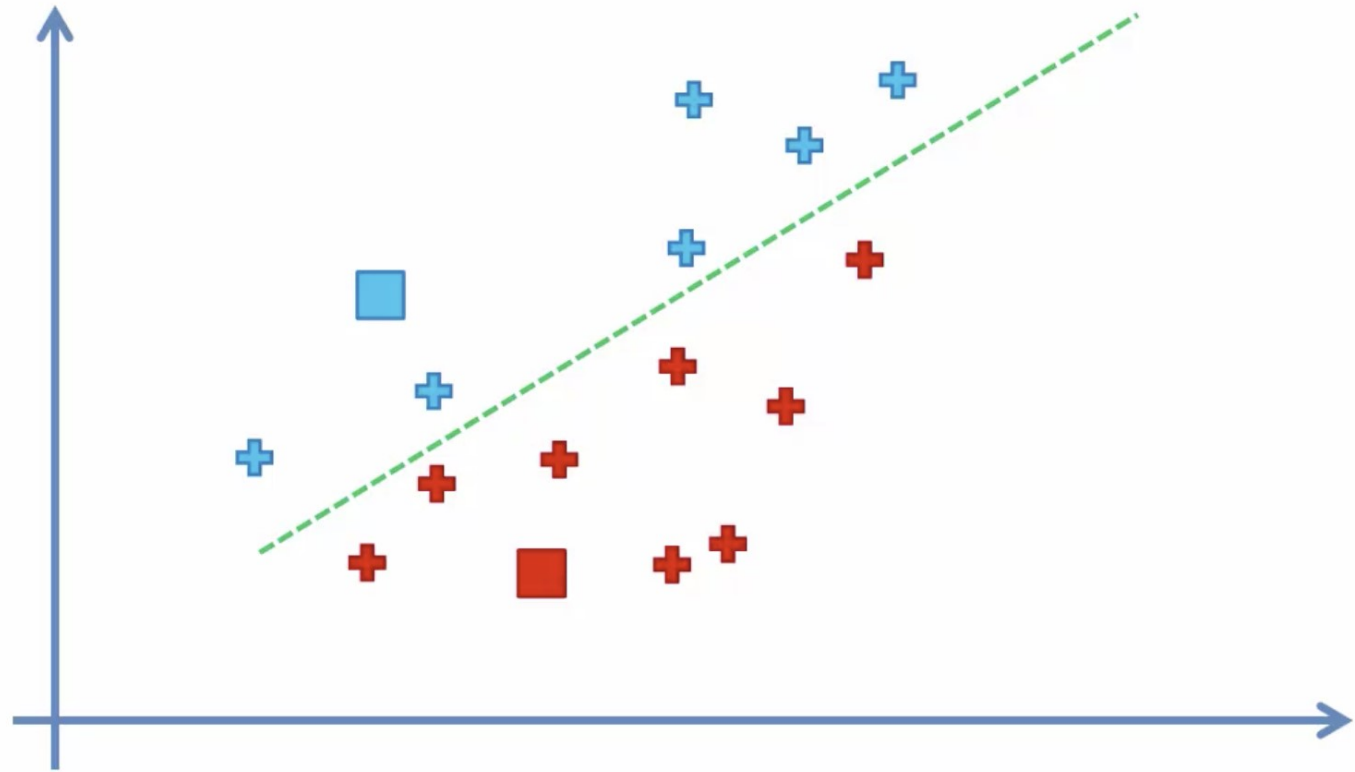
How it works

1. Select K (i.e. 2) random points as cluster centers called centroids



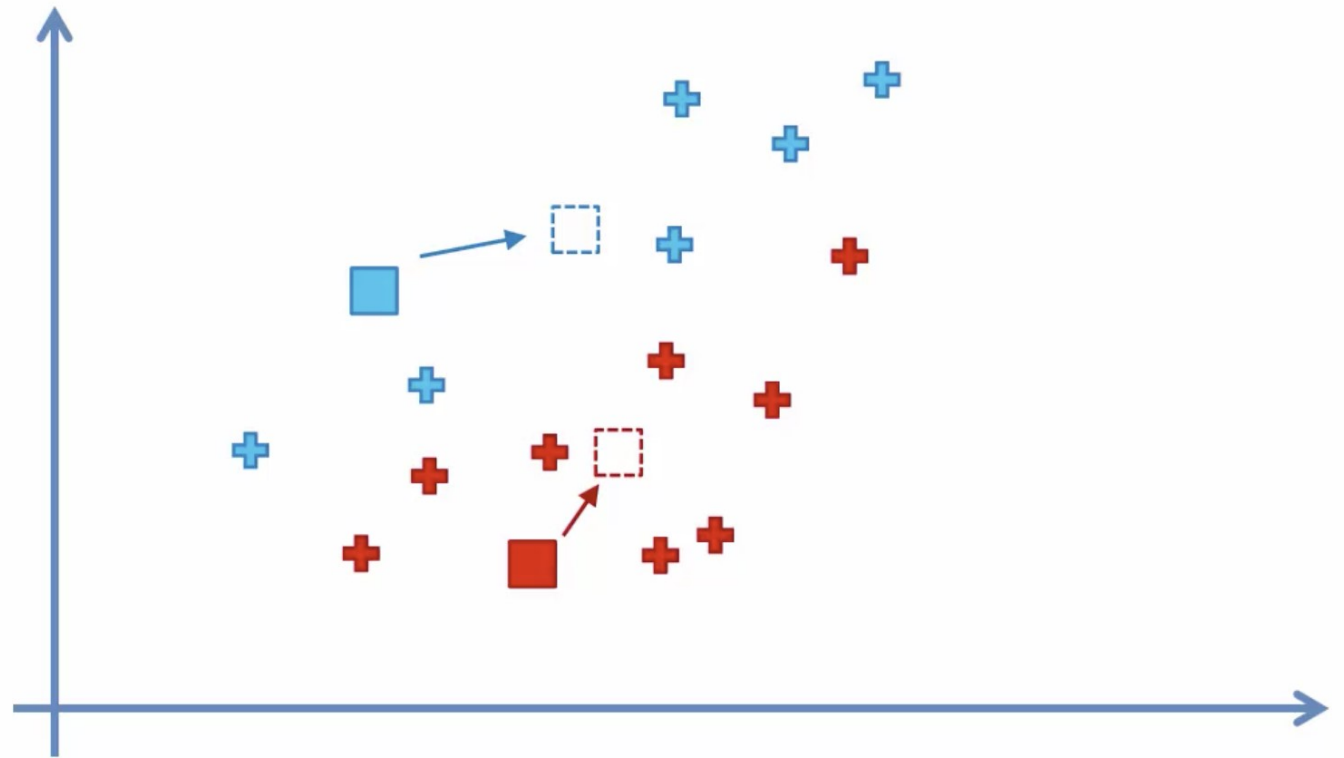
How it works

2. Assign each data point to the closest cluster by calculating its distance with respect to each centroid



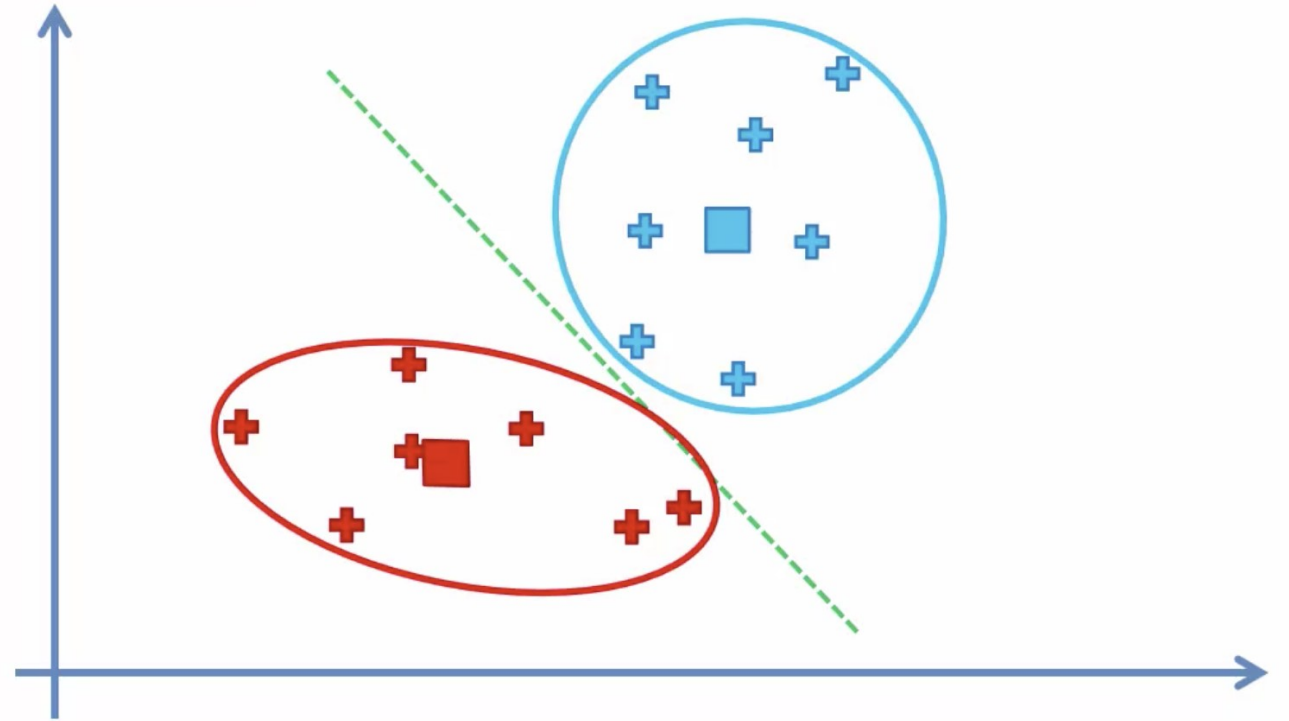
How it works

3. Determine the new cluster center by computing the average of the assigned points



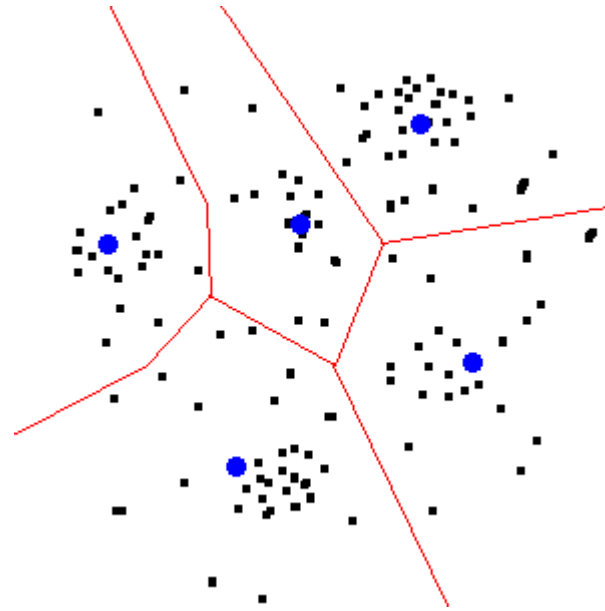
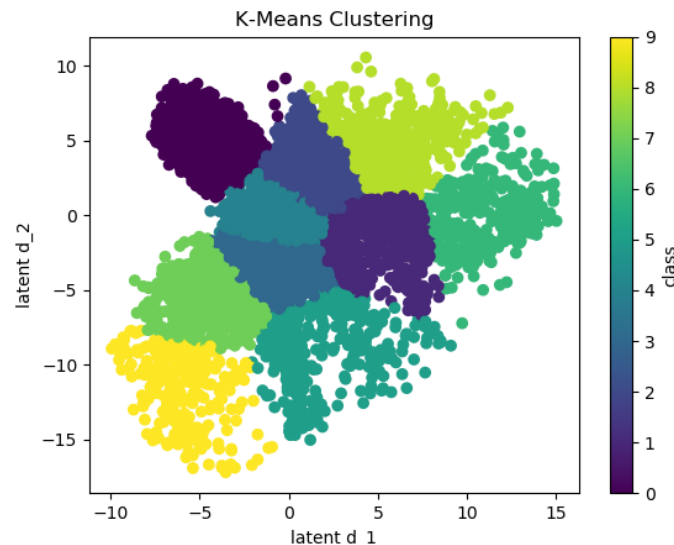
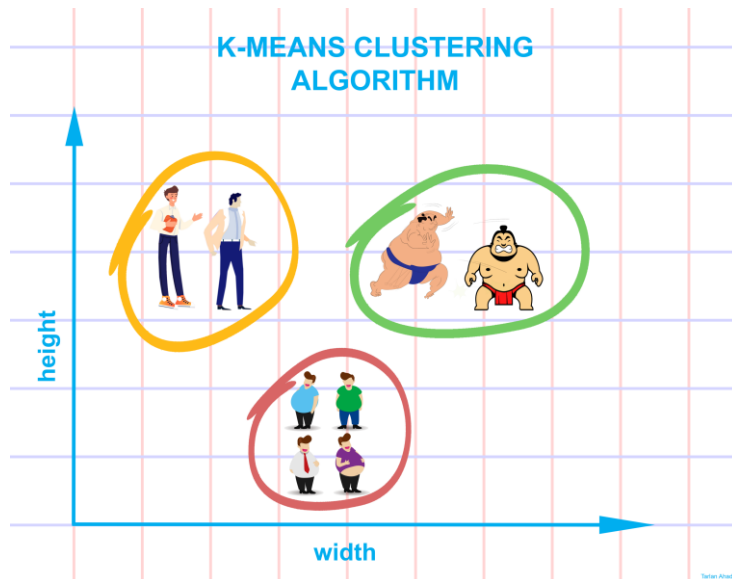
How it works

4. Repeat steps 2 and 3 until none of the cluster assignments change

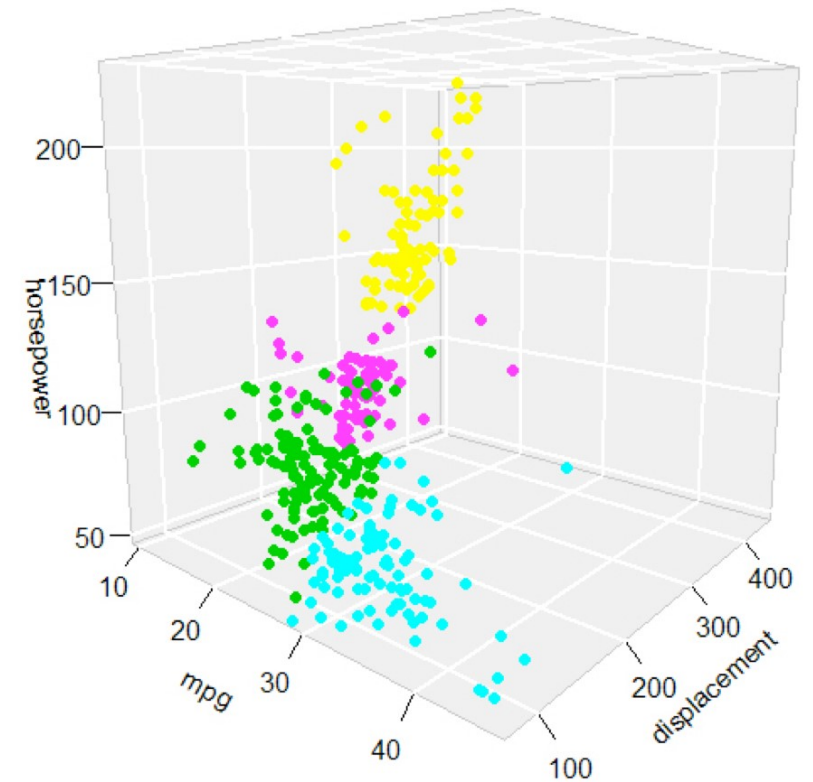


Examples

- Often the data you'll be working with will have multiple dimensions making it difficult to visual.



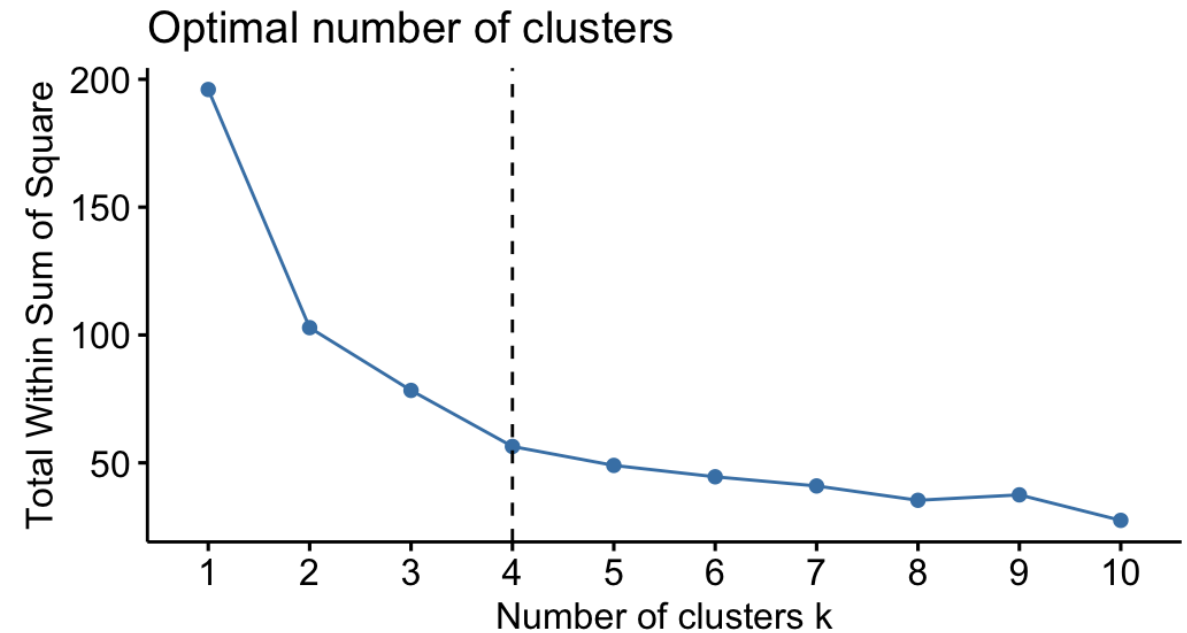
Clustering of Horsepower, MPG, and Displacement



Choosing the right number of clusters

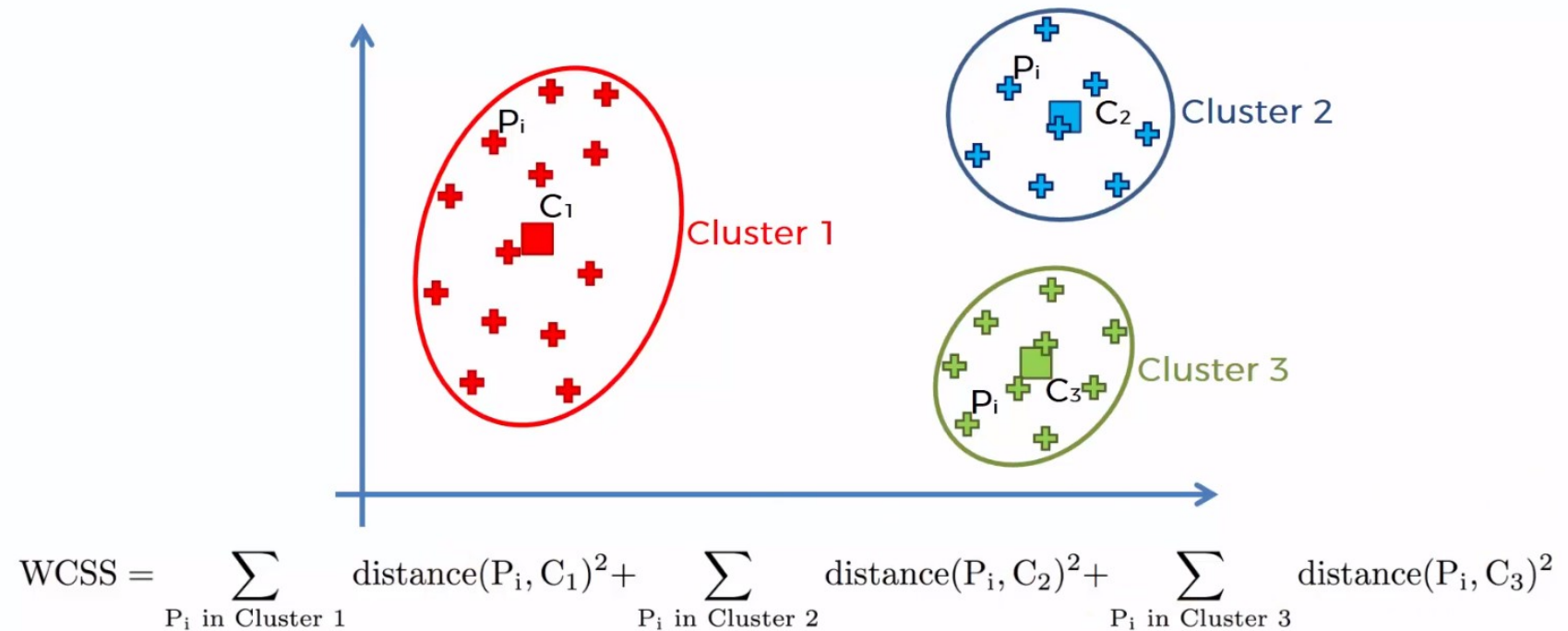
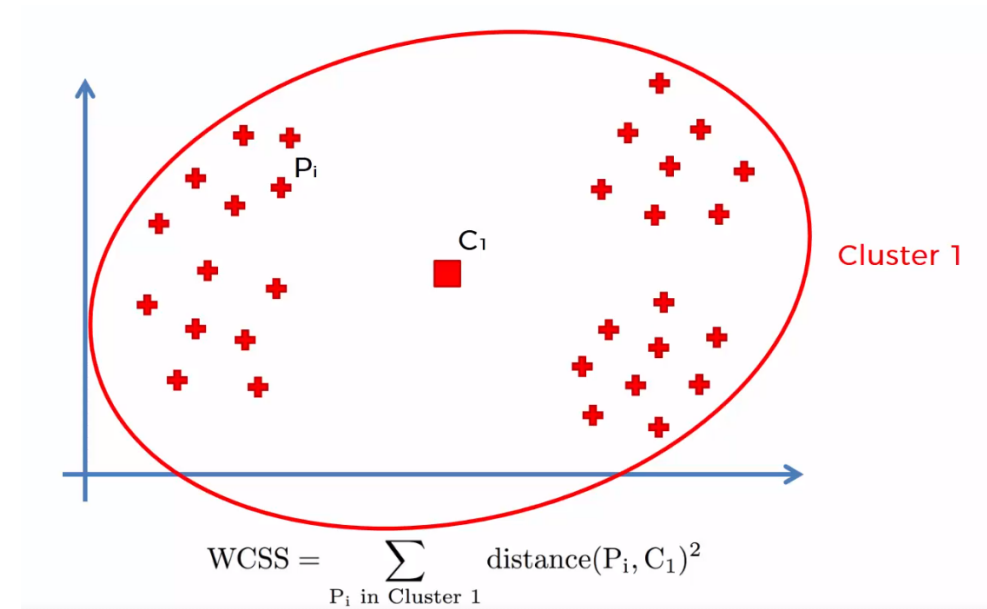
- We graph the relationship between the number of clusters and Within Cluster Sum of Squares (WCSS) then we select the number of clusters where the change in WCSS begins to level off (elbow method).
- WCSS is defined as the sum of the squared distance between each member of the cluster and its centroid.

$$WSS = \sum_{i=1}^m (x_i - c_i)^2$$



Elbow method

WCSS Calculations



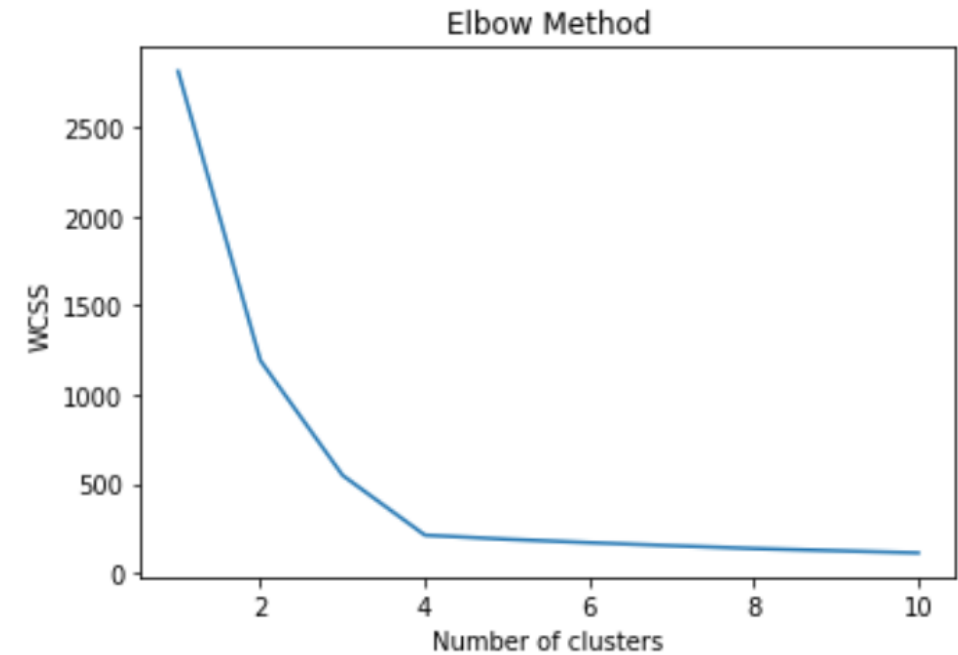
Code

- we'll generate our own data using the `make_blobs` function from the `sklearn.datasets` module. The `centers` parameter specifies the number of clusters.

```
|from matplotlib import pyplot as plt  
|from sklearn.datasets.samples_generator import make_blobs  
|from sklearn.cluster import KMeans
```

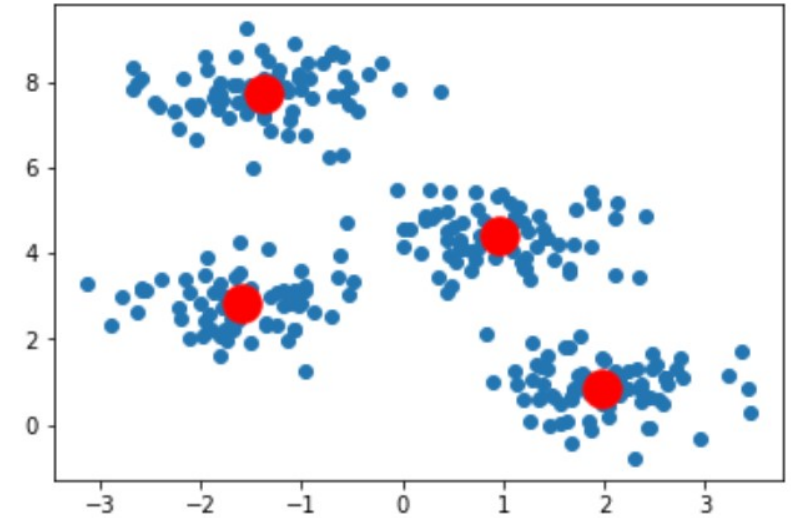
- Elbow method:
 - To get the values used in the graph, we train multiple models using a different number of clusters and storing the value of the
 - inertia_ property (WCSS) every time.

```
wcss = []  
for i in range(1, 11):  
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)  
    kmeans.fit(X)  
    wcss.append(kmeans.inertia_)  
plt.plot(range(1, 11), wcss)  
plt.title('Elbow Method')  
plt.xlabel('Number of clusters')  
plt.ylabel('WCSS')  
plt.show()
```



Code

- K-means clustering with $K = 4$ (the optimal k number)
- k-means++ ensures that you get don't fall into the random initialization trap.



```
kmeans = KMeans(n_clusters=4, init='k-means++', max_iter=300, n_init=10, random_state=0)
pred_y = kmeans.fit_predict(X)

print(pred_y)

plt.scatter(X[:,0], X[:,1])
plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s=300, c='red')
plt.show()
```

pred_y

```
=====y data=====
[1 3 0 3 1 1 2 0 3 3 2 3 0 3 1 0 0 1 2 2 1 1 0 2 2 0 1 0 2 0 3 3 0 3 3 3
 3 2 1 0 2 0 0 2 2 3 2 3 1 2 1 3 1 1 2 3 2 3 1 3 0 3 2 2 2 3 1 3 2 0 2 3 2
 2 3 2 0 1 3 1 0 1 1 3 0 1 0 3 3 0 1 3 2 2 0 1 1 0 2 3 1 3 1 0 1 1 0 3 0 2
 2 1 3 1 0 3 1 1 0 2 1 2 1 1 1 1 2 1 2 3 2 2 1 3 2 2 3 0 3 3 2 0 2 0 2 3 0
 3 3 3 0 3 0 1 2 3 2 1 0 3 0 0 1 0 2 2 0 1 0 0 3 1 0 2 3 1 1 0 2 1 0 2 2 0
 0 0 0 1 3 0 2 0 0 2 2 2 0 2 3 0 2 1 2 0 3 2 3 0 3 0 2 0 0 3 2 2 1 1 0 3 1
 1 2 1 2 0 3 3 0 0 3 0 1 2 0 1 2 3 2 1 0 1 3 3 3 3 2 2 3 0 2 1 0 2 2 2 1 1
 3 0 0 2 1 3 2 0 3 0 1 1 2 2 0 1 1 1 0 3 3 1 1 0 1 1 1 3 2 3 0 1 1 3 3 3 1
 1 0 3 2]
=====y predict data=====
[0 2 1 2 0 0 3 1 2 2 3 2 1 2 0 1 1 0 3 3 0 0 1 3 3 1 0 1 3 1 2 2 1 2 2 2
 2 3 0 1 3 1 1 3 3 2 3 2 0 3 0 2 0 0 3 2 3 2 0 2 1 2 3 3 3 2 0 2 3 1 3 2 3
 3 2 3 1 0 2 0 1 0 0 2 1 0 1 2 2 1 0 2 3 3 1 0 0 1 3 2 0 2 0 1 0 0 1 2 1 3
 3 0 2 0 1 2 0 0 1 3 0 3 0 0 0 0 3 0 3 2 3 3 0 2 3 3 2 1 2 2 3 1 3 1 3 2 1
 2 2 2 1 2 1 0 3 2 3 0 1 2 1 1 0 1 3 3 1 0 1 1 2 0 1 3 2 0 0 1 3 0 1 3 3 1
 1 1 1 0 2 1 3 1 1 3 3 3 1 3 2 1 3 0 3 1 2 3 2 1 2 1 3 1 1 2 3 3 0 0 1 2 0
 0 3 0 3 1 2 2 1 1 2 1 0 3 1 0 3 2 3 0 1 0 2 2 2 2 3 3 2 1 3 0 1 3 3 3 0 0
 2 1 1 3 0 2 3 1 2 1 0 0 3 3 1 0 0 0 1 2 2 0 0 1 0 0 0 2 3 2 1 0 0 2 2 2 0
 0 1 2 3]
```

K-Means Clustering 2

- Class Exercises