

# k-means

---

## Theory

---

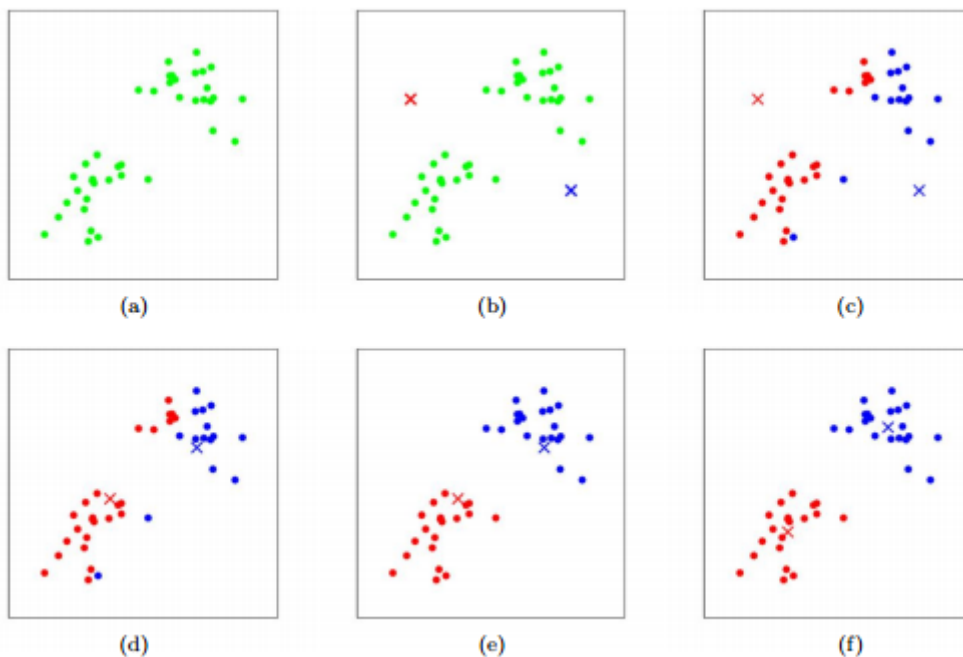
K-means is an unsupervised machine learning technique used to identify clusters of data objects in a dataset. A point is considered to be in a particular cluster if it is closer to that cluster's centroid than any other centroid.

K-means requires only a few steps. The first step is to randomly select  $k$  centroids, where  $k$  is equal to the number of clusters you choose. Centroids are data points representing the center of a cluster.

The second steps called expectation-maximization. The expectation step assigns each data point to its nearest centroid. Then, the maximization step computes the mean of all the points for each cluster and sets the new centroid. Here's the k-means algorithm looks like:

The quality of the cluster assignments is determined by computing the sum of the squared error (SSE) after the centroids converge, or match the previous iteration's assignment. The SSE is defined as the sum of the squared Euclidean distances of each point to its closest centroid. Since this is a measure of error, the objective of k-means is to try to minimize this value.

K-means can be described as follows:



1. Specify the number  $k$  of clusters to assign.
2. Randomly initialize  $k$  centroids.
3. Repeat
  - Assign each point to its closest centroid.
  - Compute the new centroid(mean) of each cluster.
4. Until The centroid position do not change

## Implement

---

## Import Packages And Dataset

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from sklearn.cluster import KMeans
4 # from sklearn import datasets
5 from sklearn.datasets import load_iris
```

```
1 iris = load_iris()
2 # Take the last two dimensions in the feature
3 x = iris.data[:,2:]
4 # drawer
5 plt.scatter(x[:, 0], x[:, 1], c = "blue", marker='+', label='label1')
6 plt.xlabel('petal length')
7 plt.ylabel('petal width')
8 plt.legend(loc=2)
9 plt.show()
```

## Start process

```
1 estimator = KMeans(n_clusters=3)
2 estimator.fit(x)
3 label_pred = estimator.labels_
```

## Visualization

```
1 x0 = x[label_pred == 0]
2 x1 = x[label_pred == 1]
3 x2 = x[label_pred == 2]
4 plt.scatter(x0[:, 0], x0[:, 1], c = "red", marker='o', label='label10')
5 plt.scatter(x1[:, 0], x1[:, 1], c = "green", marker='*', label='label11')
6 plt.scatter(x2[:, 0], x2[:, 1], c = "blue", marker='+', label='label12')
7 plt.xlabel('petal length')
8 plt.ylabel('petal width')
9 plt.legend(loc=2)
10 plt.show()
```

## Practice

Please imitate it to perform on the following dataset. The result should be visualized like above

<https://archive.ics.uci.edu/ml/datasets/seeds>