# 30-K_Means_Clustering

October 20, 2024

## 1 K-Means Clustering

**K-Means Clustering** is one of the simplest and most widely used unsupervised learning algorithms for partitioning a dataset into distinct groups or clusters. In K-Means, each cluster is represented by a centroid, and each data point is assigned to the cluster with the closest centroid.

The key idea behind K-Means is to partition data into k clusters, where k is a predefined number. The algorithm works iteratively to assign data points to one of the k clusters based on the similarity (often measured by Euclidean distance).

### 1.0.1 Why do we use K-Means Clustering?

K-Means Clustering is highly popular due to:

- **Simplicity**: Easy to understand and implement.

- **Speed**: Efficient for clustering large datasets.

- **Applicability**: Useful for a variety of tasks such as customer segmentation, image compression, and document clustering.

- **Interpretability**: Results are easy to interpret, as each data point belongs to exactly one cluster.

### 1.0.2 How does K-Means Clustering work?

The algorithm works in the following steps:

Initialization:

Choose the number of clusters k. Randomly initialize k centroids (cluster centers) in the feature space. Assignment (Step 1):

For each data point in the dataset, calculate its distance to each of the k centroids. Assign each data point to the nearest centroid (based on minimum distance). Update Centroids (Step 2):

After the assignment step, recompute the centroids by calculating the mean of all data points assigned to each cluster. The new centroid will be the new center of the cluster. Repeat:

Repeat the Assignment and Update steps iteratively until the centroids do not change significantly, or until a predefined number of iterations is reached. This is called convergence. Output:

The algorithm produces k clusters, each with a centroid and a set of assigned data points.

### 1.0.3   When to use K-Means Clustering?

K-Means is best suited for problems where:

You know the number of clusters (k): K-Means requires you to specify k in advance. Clusters are spherical and equally sized: K-Means assumes clusters are compact, and all of approximately equal size. There are no significant outliers: K-Means is sensitive to outliers, as they can heavily skew the centroid calculation. Common use cases include:

Customer Segmentation: Group customers with similar buying behaviors. Image Compression: Reduce image size by clustering similar pixels. Document Clustering: Group similar documents based on their content.

```python
[1]: from sklearn.cluster import KMeans
     import numpy as np
     import matplotlib.pyplot as plt

     X = np.array([[1, 2], [1, 4], [1, 0],
                   [4, 2], [4, 4], [4, 0]])

     kmeans = KMeans(n_clusters=2, random_state=0)

     kmeans.fit(X)
     labels = kmeans.predict(X)

     centroids = kmeans.cluster_centers_

     # Visualize the clusters
     plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')
     plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='x')
     plt.title("K-Means Clustering (k=2)")
     plt.show()
```
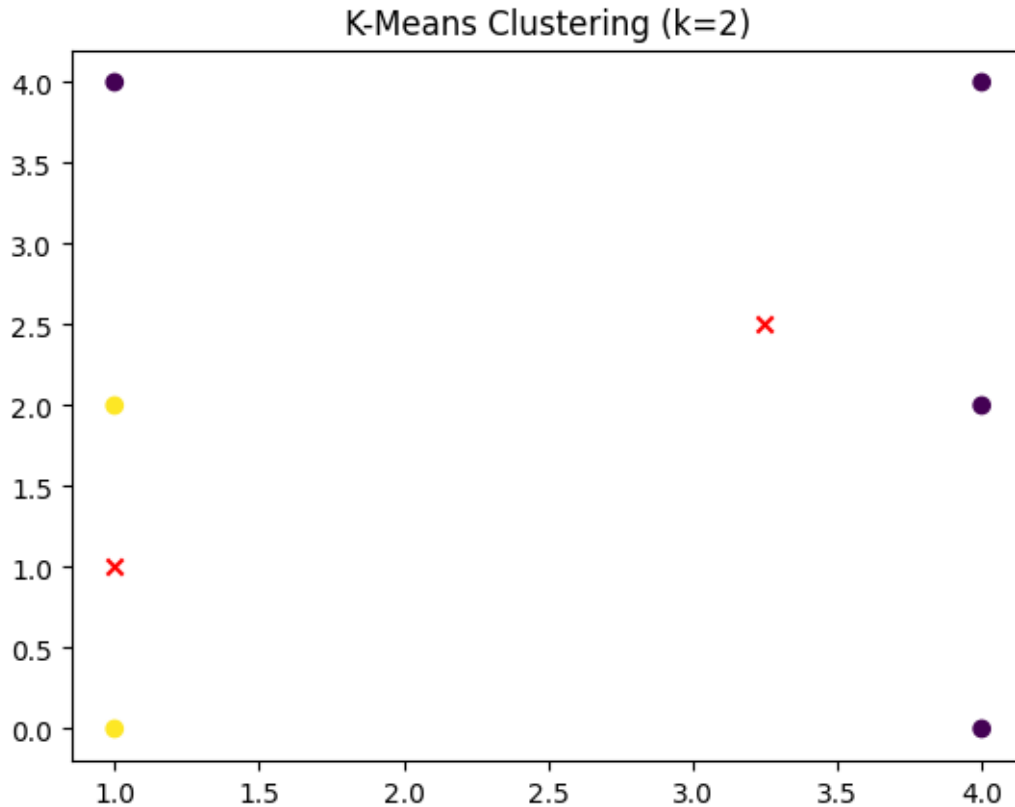
### 1.0.4 Advantages of K-Means Clustering:

Scalability: Can handle large datasets efficiently. Easy to implement: Straightforward in both understanding and coding. Fast convergence: The iterative approach usually converges quickly.

### 1.0.5 Disadvantages of K-Means Clustering:

Choosing k: The number of clusters (k) must be pre-specified, and determining the optimal k can be challenging. Sensitive to initialization: Different random initializations can lead to different clustering results (though you can mitigate this by running the algorithm multiple times or using the k-means++ initialization method). Sensitive to outliers: Outliers can significantly affect cluster centroids. Assumption of spherical clusters: K-Means works best when clusters are of equal size and spherical in shape.

### 1.0.6 Who uses K-Means Clustering?

- **Marketers**: To segment customers based on purchasing patterns.

- **Biologists**: To classify genes with similar characteristics.

- **Image Analysts**: To compress images or classify regions in an image.

- **Text Miners**: To group similar documents or articles.