

## Assignment B-6

Problem Statement:

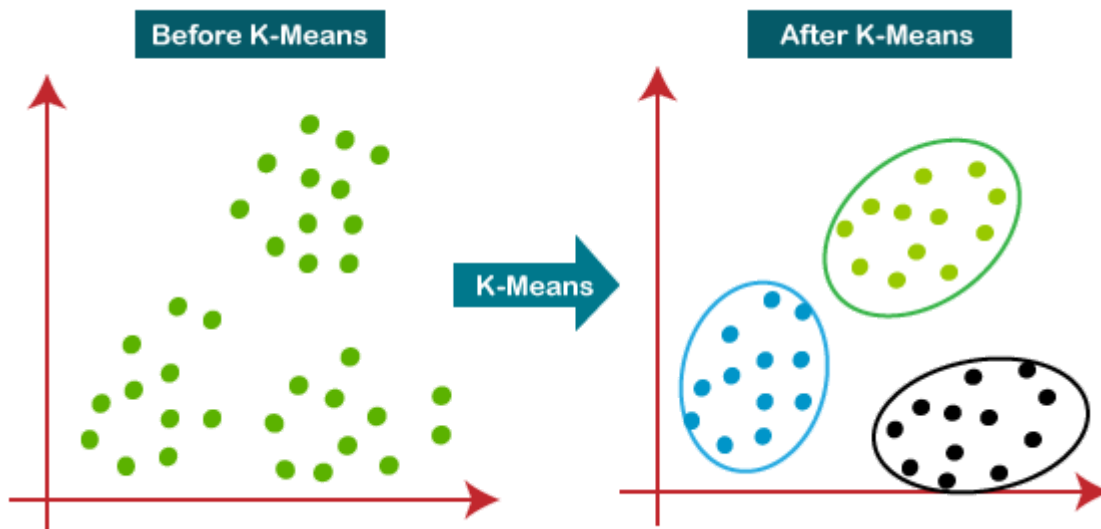
Implement K-Means clustering/ hierarchical clustering on sales\_data\_sample.csv dataset. Determine the number of clusters using the elbow method.

### K-Means Clustering Algorithm

K-Means Clustering is an [Unsupervised Learning algorithm](#), which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if  $K=2$ , there will be two clusters, and for  $K=3$ , there will be three clusters, and so on. It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training. The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

The k-means [clustering](#) algorithm mainly performs two tasks:

- o Determines the best value for K center points or centroids by an iterative process.
- o Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.



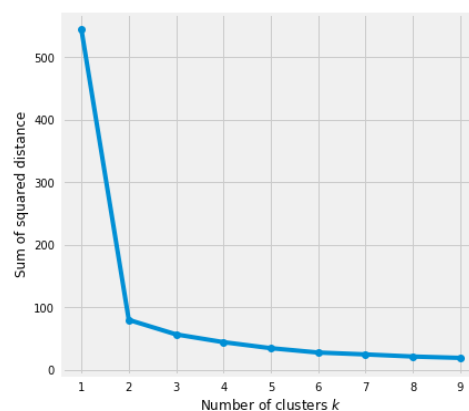
The way kmeans algorithm works is as follows:

1. Specify number of clusters  $K$ .
2. Initialize centroids by first shuffling the dataset and then randomly selecting  $K$  data points for the centroids without replacement.

3. Keep iterating until there is no change to the centroids. i.e assignment of data points to clusters isn't changing.
- Compute the sum of the squared distance between data points and all centroids.
  - Assign each data point to the closest cluster (centroid).
  - Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

### Elbow Method

Elbow method gives us an idea on what a good  $k$  number of clusters would be based on the sum of squared distance (SSE) between data points and their assigned clusters' centroids. We pick  $k$  at the spot where SSE starts to flatten out and forming an elbow.



Database Used:

Dataset link : <https://www.kaggle.com/datasets/kyanyoga/sample-sales-data>

<https://www.kaggle.com/code/victorngeno/mall-customers/data>

Python : Colab, spider or similar platform -----

YT Ref: <https://www.youtube.com/watch?v=H27rlggXlSk&t=35s>

Code (As attached) & Graphs ( wherever applicable) -----

Metrics used for performance measurement: \_\_\_\_\_

### Conclusion

Kmeans clustering is one of the most popular clustering algorithms and usually the first thing practitioners apply when solving clustering tasks to get an idea of the structure of the dataset. The goal of kmeans is to group data points into distinct non-overlapping subgroups. It does a very good job when the clusters have a kind of spherical shapes. However, it suffers as the geometric shapes of clusters deviates from spherical shapes. It also doesn't learn the number of clusters from the data and requires it to be pre-defined.