

Experiment – 7: K-means Clustering Algorithm

Date: _____

1. **Aim:** Write a program to segment the image using k-mean clustering algorithm.
2. **Requirements:** Python
3. **Pre-Experiment Exercise**

3.1 Brief Theory

k -means clustering is a method of vector quantization, originally from signal processing, that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean (cluster centers or cluster centroid), serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells. k -means clustering minimizes within-cluster variances (squared Euclidean distances), but not regular Euclidean distances, which would be the more difficult Weber problem: the mean optimizes squared errors, whereas only the geometric median minimizes Euclidean distances. For instance, better Euclidean solutions can be found using k -medians.

The problem is computationally difficult (NP-hard); however, efficient heuristic algorithms converge quickly to a local optimum. These are usually similar to the expectation-maximization algorithm for mixtures of Gaussian distributions via an iterative refinement approach employed by both k -means and *Gaussian mixture modeling*. They both use cluster centres to model the data; however, k -means clustering tends to find clusters of comparable spatial extent, while the Gaussian mixture model allows clusters to have different shapes.

The unsupervised k -means algorithm has a loose relationship to the k -nearest neighbor classifier, a popular supervised machine learning technique for classification that is often confused with k -means due to the name. Applying the 1-nearest neighbor classifier to the cluster centers obtained by k -means classifies new data into the existing clusters. This is known as nearest centroid classifier or Rocchio algorithm.

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k ($\leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares (WCSS) (i.e. variance). Formally, the objective is to find:

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \arg \min_{\mathbf{S}} \sum_{i=1}^k |S_i| \text{Var } S_i$$

Where, μ_i is the mean of points in S_i . This is equivalent to minimizing the pairwise squared deviations of points in the same cluster:

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \frac{1}{2|S_i|} \sum_{\mathbf{x}, \mathbf{y} \in S_i} \|\mathbf{x} - \mathbf{y}\|^2$$

The equivalence can be deduced from

$$\sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \sum_{\mathbf{x} \neq \mathbf{y} \in S_i} (\mathbf{x} - \boldsymbol{\mu}_i)^T (\boldsymbol{\mu}_i - \mathbf{y})$$

Because the total variance is constant, this is equivalent to maximizing the sum of squared deviations between points in different clusters, which follows from the law of total variance.

4. Laboratory Exercise

4.1 Algorithm:

K Means Clustering is an unsupervised learning algorithm that tries to cluster data based on their similarity. Unsupervised learning means that there is no outcome to be predicted, and the algorithm just tries to find patterns in the data. In k means clustering, we have to specify the number of clusters we want the data to be grouped into. The algorithm randomly assigns each observation to a cluster, and finds the centroid of each cluster. Then, the algorithm iterates through two steps: Reassign data points to the cluster whose centroid is closest. Calculate new centroid of each cluster. These two steps are repeated till the within cluster variation cannot be reduced any further. The within cluster variation is calculated as the sum of the Euclidean distance between the data points and their respective cluster centroids.

- Step-1:** Read the color image as an input and select the number K to decide the number of clusters.
- Step-2:** Select random K points or centroids. (It can be other from the input dataset).
- Step-3:** Assign each data point to their closest centroid, which will form the predefined K clusters.
- Step-4:** Calculate the variance and place a new centroid of each cluster.
- Step-5:** Repeat the third steps, which means reassign each data point to the new closest Centroid of each cluster.
- Step-6:** If any reassignment occurs, then go to step-4 else go to FINISH.
- Step-7:** Display the output segmented image into k class.

5. Post-Experiment Exercise

5.1 Conclusion:
