

Assignment - 6

Name - ~~B~~ Gaurang Patil

PRN - 1032221535

Roll No - 20.

T.Y B.Tech CSE.

Aim → W.A.P to implement k-means clustering.

Objective → To study k-means clustering algorithm.

Theory →

i) k-means clustering

Ans i) → k-means clustering is a method that sorts data into a set number of clusters, denoted as k . It starts by randomly placing cluster centres (centroids), then assigns each data point to the closest centroid. It also updates the centroids based on the points assigned to them and repeats this until the clusters stabilize.

ii) steps of k-means clustering Algo.

Ans ii) step 1 → select the no. of k to decide the number of clusters.

step 2 → Select random k points as centroid (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 \rightarrow Repeat the third step which means Reassign each datapoint to the new closet Centroid of each cluster

step 6 \rightarrow If any Reassignment occur, then go to step 4
also go to finish

step 7 \rightarrow The model is ready

3) Objective function of the k-means algorithm.

\rightarrow The objective function of the k-means clustering algo is to minimize the within-cluster sum of squares (WCSS)

Minimize within cluster Scatter - minimize the distance b/w data points assigned to be cluster.

Maximize between-cluster Scatter - Maximize the distance between cluster points assigned to different clusters

FAQ \rightarrow

1) How to determine the k using the Elbow method?

Ans. i) Run k-means for a Range k

ii) Calculate Inertia for each k

iii) Plot Inertia v k on a graph

iv) Pick the elbow point. The optimal k is where the curve bends.

2) Describe the initialization step in the k-means algo.

\rightarrow Random Initialization - Select k data points Randomly as the initial centroid.

3) k Means the Initialization - Select the first centroid Randomly and Subsequent Centroid are chosen based on optimizing the maximum distance from Centroid.

Importance of Initial Centroids

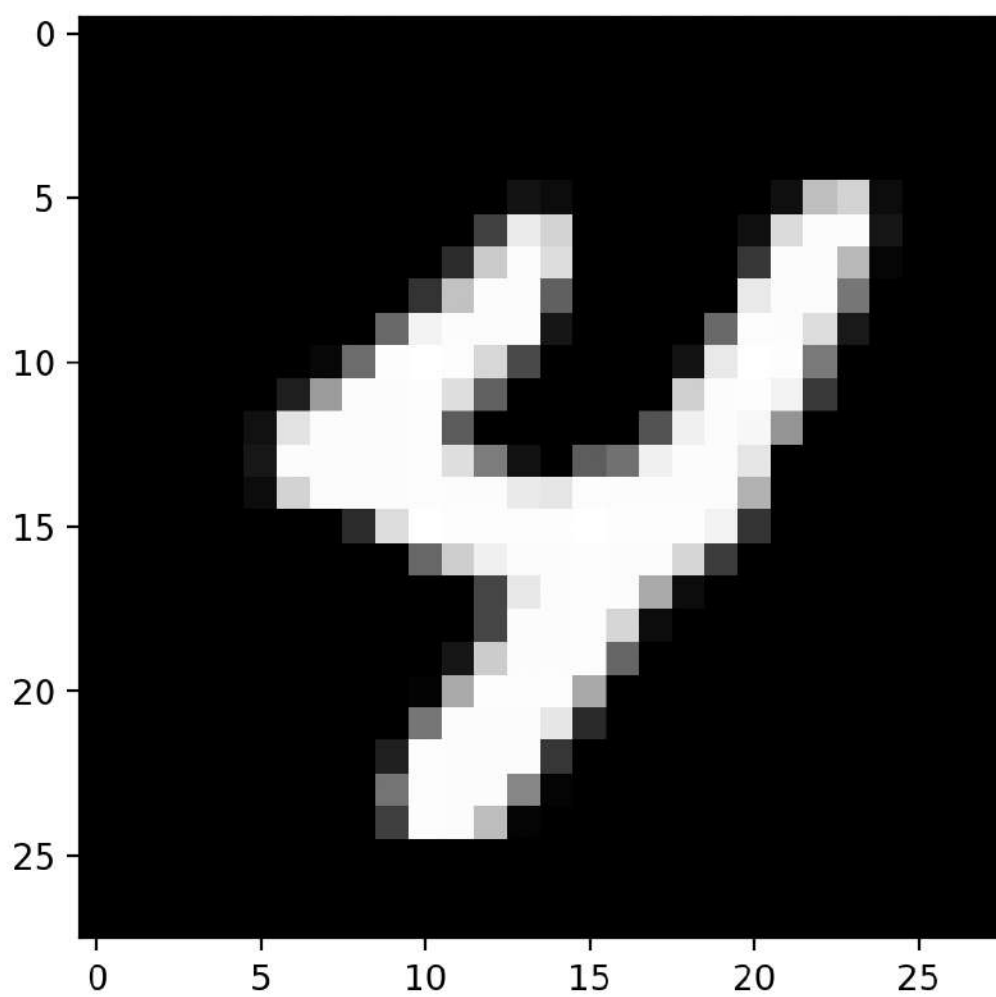
Impact on convergence: Poor initialization can lead to slower convergence.

Push of local minima: Random centroid might cause to also to get stuck in suboptimal cluster.

Consistency of Results: Different initial centroid can yield different clusters, good initialization helps provide consistent results.

~~Yogita~~
24/10/2024

Figure 1



Accuracy: 0.8262195121951219
Iteration: 350
Accuracy: 0.8283170731707317
Iteration: 360
Accuracy: 0.830780487804878
Iteration: 370
Accuracy: 0.8323658536585365
Iteration: 380
Accuracy: 0.8342682926829268
Iteration: 390
Accuracy: 0.8360731707317073
Iteration: 400
Accuracy: 0.8376829268292683
Iteration: 410
Accuracy: 0.8390731707317073
Iteration: 420
Accuracy: 0.8404390243902439
Iteration: 430
Accuracy: 0.8422926829268292
Iteration: 440
Accuracy: 0.8437317073170731
Iteration: 450
Accuracy: 0.845170731707317
Iteration: 460
Accuracy: 0.8462439024390244
Iteration: 470
Accuracy: 0.8473414634146341
Iteration: 480
Accuracy: 0.848609756097561
Iteration: 490
Accuracy: 0.8496829268292683
Prediction: [4]
Label: 4