

Assignment No : 6

Roll no : 4351

1 Title

K Means Clustering

2 Problem Statement

Implement a simple approach for k-means/ k-medoids clustering using C++.

3 Learning Objectives

1. To understand Data Mining Concepts.
2. To develop problem solving abilities using Mathematical Modelling
3. To develop time and space efficient algorithms

4 Learning Outcome

1. After successfully completing this assignment, you should be able to Understand Implement efficient design, analysis and testing of algorithmic assignments.
2. Also we will learn a very effective clustering technique, which will help us in studying further techniques

5 Theory

In statistics and machine learning, k-means clustering is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

k-means clustering is a method of vector quantization, originally from signal processing, that is popular for cluster analysis in data mining. k-means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells.

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

The algorithm has a loose relationship to the k-nearest neighbor classifier, a popular machine learning technique for classification that is often confused with k-means because of the k in the name.

6 Mathematical Model

Let S be the solution perspective of the system such that
 $S = \{ s, e, X, Y, Fme, DD, NDD, Ffriend \mid \phi s \}$

where

DD = Deterministic Data

NDD = Non - Deterministic Data = User Input

ϕs = Constraints on System S

s is start state

e is end state

$X = x_1, x_2, x_3, \dots, x_n$ — 'X' n instances of data.

$A = a_1, a_2, a_3, \dots, a_p$ — A p attributes of each instance X.

Operations Performed : initializecluster(X, Y), assigncluster(X, Y), updatecluster(X, Y)

initializecluster(X)1 = 'C' randomly selecting k cluster centres.

$C = m$ — m is the center of cluster C

assigncluster(X, C) = B'; calculate the nearest cluster center to each X_i .

$C_j = m_j$ — center of cluster C_j

if $\min(d(X_i, m_1), d(X_i, m_2), \dots, d(X_i, m_k)) = d(X_i, m_j)$ then $X_i = m_j$

$X_j = X_{j,1}, X_{j,2}, \dots, X_{j,p}$ — p attributes of X_j

$m_j = m_{j,1}, m_{j,2}, \dots, m_{j,p}$ — p attributes of m_j

$d(X_i, m_j) = \sqrt{(X_{i,1} - m_{j,1})^2 + (X_{i,2} - m_{j,2})^2 + \dots + (X_{i,p} - m_{j,p})^2}$

updatecluster(C) = C'; calculate the new cluster center. $m = \frac{1}{n} \sum_{i=1}^n X_i$ — no of instances within cluster

Y=Output

Y= C1, C2, C3, CK —O outcome of operations performed by the system k clusters where, Y1=Name

Y2 = Email

Y3 = Phone no

Success Case = Partition Successful

Failure Case = Partition Unsuccessful

7 Conclusion

Thus, after successfully completing this assignment, we understood and implemented k - mean clustering technique to partition instance data into k clusters in C++.