What is clustering? Explain in detail about k-means clustering with an example?

clustering:- "cluster analysis or clustering is the task of grouping a set of objects or data points in such a way that objects in the same group are more similar to each other than to those in other groups

(or)

clustering is the process of grouping objects into different groups which are meaningful, useful or both"

k-Means:-

A prototype based one level partitional clustering technique that attempts to find a user- specified number of cluster(k).

It is the oldest and most widely used clustering algorithm.

k-means defines a prototype as a 'centroid' Applied to objects in a continuous n-dimensional space

The Basic k-means Algorithm:-

choose k initial centroids, k-user specified and k indicates number of clusters desired

Each point is then assined to the closest centroid.

centroid of each cluster is then updated based on the assigned to the cluster

# Algorithm:-

## Algorithm Basic k-means algorithm

1. select $k$ point as initial. Centroids
2. repeat
3. Form $k$ clusters by assigning each point to its closest control.
4. Recompute the centroid of each cluster
5. Untill Centroids do not change.

## Assigning point to the closest Centroid

To assign a point into a cluster, a measure must be calculated for finding the closest of point to a Centroid.

Various Measures are:-

Euclidean $L_2$

Cosine similar. is more appropriate for documents

Manhattan $(L_1)$

Jaccard measure

| symbol | Description |
|---|---|
| $x$ | An object |
| $C_i$ | The $i$th cluster |
| $c_i$ | The Centroid of cluster $C_i$ |
| $c$ | The Centroid of all points |
| $m_i$ | The number of objects in the $i$th cluster |
| $m$ | The number of objects in the dataset. |
| $k$ | The number of clusters. |

# Centroids and objective Functions

Minimize the Squared distance of each point to its closest centroid

## Data in Euclidean space

objective Function. measures the quality of a clustering

calculate the Error of each datapoint.

SSE is formally defined as follows:

$$SSE = \sum_{i=1}^{k} \sum_{x \in C_i} dist(C_i, x)^2$$

Where dist is the standard Euclidean ($L2$) distance between two objects in Euclidean space.

The centroid that minimizes.

$$C_i = \frac{1}{m_i} \sum_{x \in C_i} x$$

To illustrate, the centroid of a cluster containing the three two-dimensional points $(4,1), (2,3)$ and $(6,2)$

## Document data:

k means is also used for document data
Document data is represented as a document-term matrix

$$Total\ Cohesion = \sum_{i=1}^{k} \sum_{x \in C_i} Cosine(x, C_i)$$

choosing __Initial Centroids__! Randomly selected initial Centroids may be poor. One technique that is commonly used to address the problem of choosing Initial Centroids is to perform multiple runs, each with a different set of randomly. chosen Initial Centroids. and then select the set of clusters with the minimum SSE