

Introduction falls into
the general category
of variance-based



**ASSIGN
BUSTER**

INTRODUCTION CLUSTERING problems arise nowadays in many different applications, such as data mining and knowledge discovery, data compression and vector quantization, and pattern recognition and pattern classification.

The notion of what constitutes a good cluster depends on the application and there are many methods for finding clusters subject to various criteria, both ad hoc and systematic. These include approaches based on splitting and merging such as ISODATA randomized approaches such as CLARA, CLARANS, methods based on neural networks, and methods designed to scale to large databases, including DBSCAN, BIRCH, etc. Surrounded by clustering formulations that are dependent on minimizing a proper objective function, possibly the most widely used and studied is k-means clustering. Given a set of n data points in real d -dimensional space, R^d , and an integer k , the problem is to determine a set of k points in R^d , called centers, so as to minimize the mean squared distance from each data point to its nearest centre. This measure is often called the squared-error distortion and this type of clustering falls into the general category of variance-based clustering. Clustering based on k-means is closely related to a number of other clustering and location problems.

These include the Euclidean k -medians in which the objective is to minimize the sum of distances to the nearest centre and the geometric k -centre problem in which the objective is to minimize the maximum distance from every point to its closest center. There are no efficient solutions known to any of these problems and some formulations are NP-hard. k-means is one of the simplest unsupervised learning algorithms that solve the well known <https://assignbuster.com/introduction-falls-into-the-general-category-of-variance-based/>

clustering problem. The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters) fixed a priori.

The main idea is to define k centers, one for each cluster. These centers should be placed in a cunning way because of different location causes different result. So, the better choice is to place them as much as possible far away from each other. The next step is to take each point belonging to a given data set and associate it to the nearest center. When no point is pending, the first step is completed and an early grouping is done.

At this point we need to re-calculate k new centroids as barycenter of the clusters resulting from the previous step. After we have these k new centroids, a new binding has to be done between the same data set points and the nearest new center. A loop has been generated. As a result of this loop we may notice that the k centers change their location step by step until no more changes are done or in other words centers do not move any more. Finally, this algorithm aims at minimizing an objective function known as squared error function given by:

Euclidean distance between x_i and v_j . c_i is the number of datapoints in i th cluster. c is the number of cluster centers.

Algorithmic steps for k -means clustering Let $X = \{x_1, x_2, x_3, \dots,$

$x_n\}$ be the set of data points and $V = \{v_1, v_2, \dots, v_c\}$ be the set of centers.

1) Randomly select 'c' cluster centers. 2) Calculate the distance between each data point and cluster centers. 3) Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers..

4) Recalculate the new cluster center using: 5) Recalculate the distance between each data point and new obtained cluster centers. 6) If no data point was reassigned then stop, otherwise repeat from step 3. For Example- Suppose that we have n sample feature vectors x_1, x_2, \dots

, x_n all from the same class, and we know that they fall into k compact clusters, $k < n$. Let m_i be the mean of the vectors in cluster i. If the clusters are well separated, we can use a minimum-distance classifier to separate them. That is, we can say that x is in cluster i if $\|x - m_i\|$ is the minimum of all the k distances.

This suggests the following procedure for finding the k means: Make initial guesses for the means m_1, m_2, \dots, m_k Until there are no changes in any mean Use the estimated means to classify the samples into clusters For i from 1 to k Replace m_i with the mean of all of the samples for cluster i end_for end_until Here is an example showing how the means m_1 and m_2 move into the centers of two clusters

Advantages 1) Fast, robust and easier to understand. 2) Relatively efficient: $O(tknd)$, where n is # objects, k is # clusters, d is # dimension of each object, and t is # iterations. Normally, $k, t, d \ll n$. 3) Gives best result when data set are distinct or well separated from each other.

Note: For more detailed figure for k-means algorithm please refer to k-means figure subpage. Disadvantages 1) The learning algorithm requires a priori specification of the number of cluster centers. 2) The use of Exclusive Assignment - If there are two highly overlapping data then k-means will not be able to resolve that there are two clusters. 3) The learning algorithm is not invariant to non-linear transformations i.

e. with different representation of data we get different results (data represented in form of cartesian co-ordinates and polar co-ordinates will give different results). 4) Euclidean distance measures can unequally weight underlying factors. 5) The learning algorithm provides the local optima of the squared error function. 6) Randomly choosing of the cluster center cannot lead us to the fruitful result. 7) Applicable only when mean is defined i.

e. fails for categorical data. 8) Unable to handle noisy data and outliers. 9) Algorithm fails for non-linear data set.