

DM assignment - 3

Q1) What is clustering? How is it different from classification?

Ans Cluster analysis groups data objects based only on information found in the data that describes the objects and their relationships. The goal is that the objects within a group be similar to one another and different from the objects in other group.

- Clustering analysis is related to other techniques that are used to divide data objects into groups. For instance, clustering can be regarded as a form of classification in that it creates a labelling of objects with class (cluster) labels.

However, it derives these labels only from the data. In contrast, classification, new unlabelled objects are assigned a class label using a model developed from objects with known class labels.

Hence, cluster analysis is sometimes referred to as unsupervised classification.

when the term classification is used without any qualification within data mining, it typically refers to supervised classification

1P8170544

Q3 Discuss the different types of clusters.

Ans Types of clusters:

- 1) well-separated
- 2) Prototype-Based
- 3) Graph-Based.
- 4) Density-Based
- 5) Shared-Property.

1) well separated: A cluster is a set of objects in which each object is closer to every other object in the cluster than to any object not in the cluster. Sometimes a threshold is used to specify that all objects in a cluster must be sufficiently close to another.

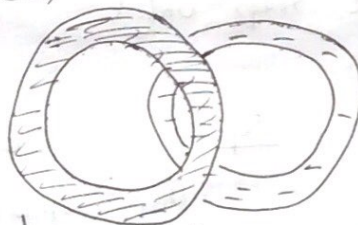
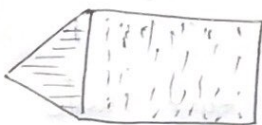
2) Prototype-Based: A cluster is a set of objects in which each object is closer to the prototype that defines the cluster than to the prototype of any cluster. For data with continuous attributes the prototype of a cluster is a centroid i.e., mean of all points in cluster.

3) Graph based: If the data is represented as a graph, where the nodes are objects & the links are connections among objects, then a cluster can be defined as a connected component.

Eg: Contiguity-based clusters; where two objects are connected only if they are within a specified distance of each other.

4) Density Based: A cluster is a dense region of objects that is surrounded by a region of low density... This type is employed when the clusters are irregular or intertwined & when noise and outliers are present.

5) Shared-Property: The shared property approach includes new type ~~of~~ clusters



A triangular area is adjacent to rectangular area & the two intertwined circles. In both cases, clustering algorithm would need a very specific concept of cluster to successfully detect these clusters. This is called Conceptual clustering.

Q4) Explain K-means algorithm what are its limitations? IPE17CS044

Ans Basic K-means Algorithm:

1. Select K points
 2. repeat
 - Form K clusters by assigning each point to its closest centroid.
 - Recompute the Centroid of each cluster
 - until centroids do not change
- we choose K initial centroids, (K is user specified) namely each, no. of clusters to be used / desired.
 - Each point is then assigned to the closest centroid & collection of points assigned to the cluster is a cluster.
 - Repeat the step until no points assigned to cluster is changed:

Step: Centroid of each cluster is then updated based on the points assigned to the cluster.

Limitations:

- ~~K-means, although are more efficient, they are less susceptible to initialization problems.~~
- K-means is not suitable for all types of data
- It cannot handle non-globular clusters or cluster of different sizes & densities

- K-means also has trouble clustering data that contains outliers.
- K-means is restricted to data for which there is a notion of a center.

Q5 How are density based methods for clustering. How are parameters of the DBSCAN algo selected?

Ans) Density based clustering locates regions of high density that are separated from one another by regions of low density. DBSCAN is a simple and effective density-based algorithm that illustrates a no. of important concepts that are important for any density-based clustering approach.

Selection of DBSCAN params:

There is, of course, the issue of how to determine the parameters ϵ and Min Pts. The basic approach is to look at the behaviour of the distance from a point to its k^{th} nearest neighbor. ($k\text{-dist}$).

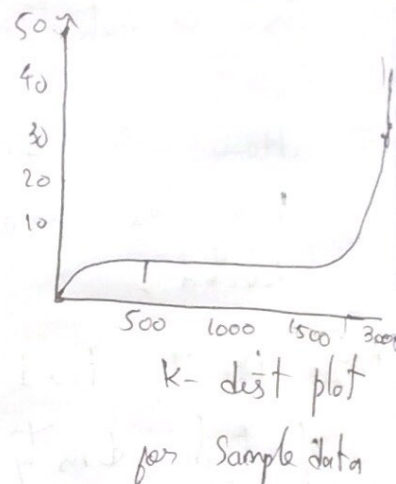
Therefore, if we compute the $k\text{-dist}$ for all the data points for some k , sort in ascending order, then plot the sorted values, we expect to see

a sharp change at the value of k -dist that corresponds to a suitable value of ϵ ps.

Example



Sample data



Q2 List and explain the desired features of cluster analysis

- Ans
- 1) Clustering for understanding
 - 2) Clustering for utility

1) Clustering for understanding: Classes or conceptually meaningful groups of objects that share common characteristics, play an important role in how people analyse & describe the world. In the context of understanding data, clusters are potential classes & cluster analysis is the study of technique for finding classes automatically.

Eg: Biology, information retrieval, Climate, Business etc

clustering for utility:

1PE17CS44

Cluster analysis provides an abstraction from individual data objects to the clusters in which these data objects ~~reside~~ reside.

∴ In context of ~~test~~ of utility, clustering analysis is the study of techniques for finding the most representative cluster prototypes.

Ex. Summarization, Compression, Efficiently finding nearest neighbors.

Q6 Explain different methods of computing distances between clusters

Ans. 1) The different methods of computing distances are

1. Single link or MIN

The proximity of 2 clusters is defined as the minimum of the distance between any two points in the two different clusters.

It is given by:

$$\min \{ d(a, b), a \in A, b \in B \}$$

a, b - points in clusters A, B

A, B - clusters A, B

2) Complete link or MAX or CLIQUE

IPETCS044

The proximity of two clusters is defined as the maximum of the distance between any two points in the two different clusters.

Given by:

$$\text{Max} \{d(a, b), a \in A, b \in B\}$$

3) Group Average:

For the group average, the proximity of two clusters is defined as the average pairwise proximity among all pairs of points in the different clusters.

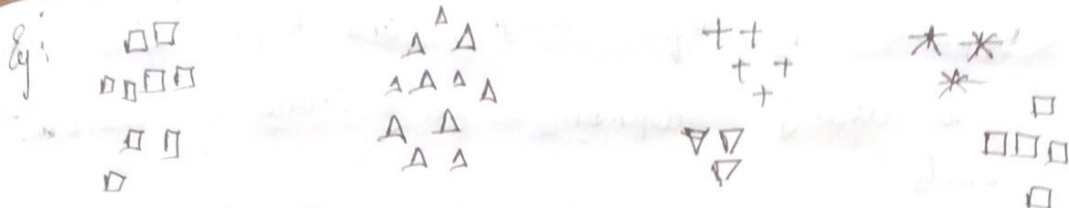
$$\text{Thus proximity } (C_i, C_j) = \frac{\sum_{\substack{x \in C_i \\ y \in C_j}} \text{proximity}(x, y)}{m_i \times m_j}$$

C_i, C_j — clusters of size m_i, m_j

Q7 Describe the following approaches to clustering with an example in each case.

Sol/ a) Partitioning Method:

A partitional clustering is simply a division of the set of the data objects into non-overlapping subsets (clusters) such that each data object is exactly one subset.



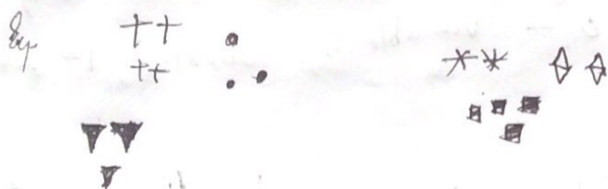
Two clusters

Four clusters

b) Hierarchical Clustering

IPETCS044

If we permit clusters to have subclusters, then obtain a hierarchical clustering, which is a set of nested clusters that are organized as a tree. Each node (cluster) in the tree is the union of its children and the root of the tree is the cluster containing all the objects.



c) Briefly outline how to compute dissimilarity between objects described by

d) Interval Scaled variable

These variables are continuous measurements of a roughly linear scale. Eg: height, latitude, meters etc.

• Standardizing measurements attempts to give all

variables an equal weight, and helps to avoid dependence on choice of measurement units

- Standardize using mean absolute deviation:

~~Standardize~~

- Distances are normally used to measure the similarity or dissimilarity between two data objects.

1) Euclidean distance, $d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots}$

2) Manhattan distance, $d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots$

3) Minkowski distance, $d(i, j) = \sqrt[q]{|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots}$

- b) Binary Variables: These variables has only two states 0/1, 0 - variable is absent, 1 - variable is present.

One approach involves computing dissimilarity matrix or if all binary variables have same weight, we have Contingency table.

Symmetric binary dissimilarity; $d(i, j) = \frac{r + s}{q + r + s + t}$

Asymmetric binary dissimilarity; $d(i, j) = \frac{r + s}{q + r + s}$

Q4 Illustrate grid-based clustering algorithm. How clusters are formed from dense-grid cells. 1PE17CS044

Ans Grid-based clustering algorithm:

1. Define a set of grid cells
2. Assign objects to the appropriate cells and compute the density of each cell.
3. Eliminate cells having a density below a specified threshold, T .
4. Form clusters from contiguous groups of dense cells.

In many cases, the data has both spatial & non-spatial attributes. In other words, some of the attributes describe the location of objects in the time or space, while other attributes describe other aspects of the objects. A

Common example is houses, that have both a location & characteristics such as price, space etc. Because of spatial autocorrelation, objects in a particular cell often have similar values for their other attributes. In such cases, it is possible to filter

the cells based on statistical properties of IPEL7CS44
one or more non-spatial attributes, and then
form clusters based on the density of the
remaining points.

Q10) Develop DENCLUE algorithm for kernel density function.

Ans.

DENCLUE algorithm:

1. Derive a density function for the space occupied by the data points.
2. Identify the points that are local maxima.
3. Associate each point with a density attractor by moving in the direction of maximum increase in density.
4. Define clusters consisting of points associated with a particular density attractor.
5. Discard clusters whose density attractor has a density less than a user-specified threshold of ϵ .
6. Combine clusters that are connected by a path of points that all have density of ϵ or higher.