

UNIT-V :- (I) Clustering

Clustering :- A collection of data objects
→ similar (one another within same group) [related]
→ dissimilar [no objects in other groups] [unrelated]

Cluster Analysis : [clustering, data segmentation]

→ Finding similarities between data according to the characteristics found in the data grouping similar data objects into clusters.

Unsupervised learning :-

→ no predefined classes (ie, learning by observations vs learning by examples : supervised)

→ Applications :-

→ A stand-alone tool to get insight into data distribution

→ As a preprocessing step for other algorithms

→ Biology :- Animal kingdom

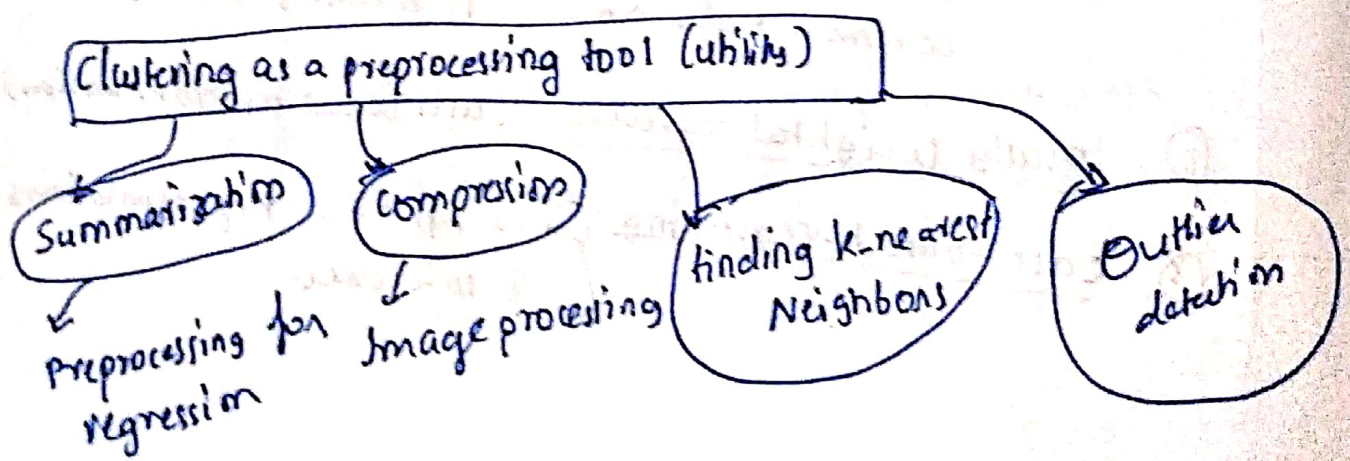
→ Information retrieval : document clustering

→ land use : similar lands

→ Marketing :- Targeted marketing programs by discovering distinct groups

→ city planning & earth quake studies :

→ Economic science :



quality (what is good clustering)

- A good clustering method will produce high quality clusters
 - high Intra-class similarity: cohesive within clusters
 - lower Inter-class similarity: distinctive between clusters.
- The quality of a clustering method depends on
 - The similarity measure used by the method
 - Its implementation
 - The ability to discover some or all hidden patterns

measuring the quality of the cluster

- dissimilarity / similarity metric
 - Similarity expressed in terms of a distance function
 - The definitions of distance functions are usually rather different for interval scaled, boolean, categorical, ordinal ratio and vector variables.
 - weights should be associated with different variables
- quality of clustering:

- This is usually separate "quality" function that measures the "goodness" of a cluster
- It is hard to define "similar enough" or "good enough" because they are subjective

Considerations of cluster Analysis

- partitioning criteria
 - single level vs hierarchical partitioning
- Separation of clusters
 - exclusive (one customer belongs to only one region) vs non-exclusive (eg.: one document belong to more than one class)
- Similarity measure
 - distance based (euclidean, vector) vs connectivity based (density, contiguity)

- Clustering space
 - Full space (often when low dimensional) vs subspaces (often in high dimensional clustering)

Requirements & challenges

- Scalability :-
 - Clustering all the data instead of samples
- Ability to deal with different kind of attributes :-
 - Numerical, binary, categorical, ordinal etc
- Constraint based clustering :-
 - user may give inputs on constraints
 - use domain knowledge to determine input parameters
- Interpretability and usability:
- Others:
 - discovery of clusters with arbitrary shape
 - ability to deal with noisy data
 - Incremental clustering and insensitivity to input order
 - High dimensionality

Major clustering approaches

- Partitioning approach:
 - Construct various partitions and then evaluate them by some criterion
eg: minimizing the sum of square errors
 - Typical methods: ~~Partitioning~~, K-Medoids, CLARANS
- Hierarchical approach:
 - creates a hierarchical decomposition of set of data by some criteria.
 - Typical methods: diana, Agnes, BIRCH, Chameleon

• Density based approach :

- Based on connectivity & density function
- Typical methods: DBSCAN, OPTICS, DENCLUE

• Grid based approach :

- Based on multiple-level granularity structure
- Typical methods: STING, waveCluster, CLIQUE

• Model Based :-

- A model is hypothesized for each of the clusters and tries to find the best fit of model to each other
- Typical method: EM, SOM, COBWEB

• Frequent pattern Based :-

- Based on analysis of frequent patterns.
- Typical methods: P-cluster

• User Guided / Constraint Based

- Clustering by considering user-specified or application specific constraints.
- Typical methods: COD (obstacles), constrained clustering

• Link based clustering :

- objects are often linked together in various ways
- massive links can be used to cluster objects: simRank_s

Link Clus

Partition Based Methods

- Partition method :- Partition a database D of n objects into a set of k clusters, such that sum of squared distances is minimized.
(where c_i is the centroid or median of cluster)

$$E = \sum_{i=1}^k \sum_{p \in c_i} (p - c_i)^2$$

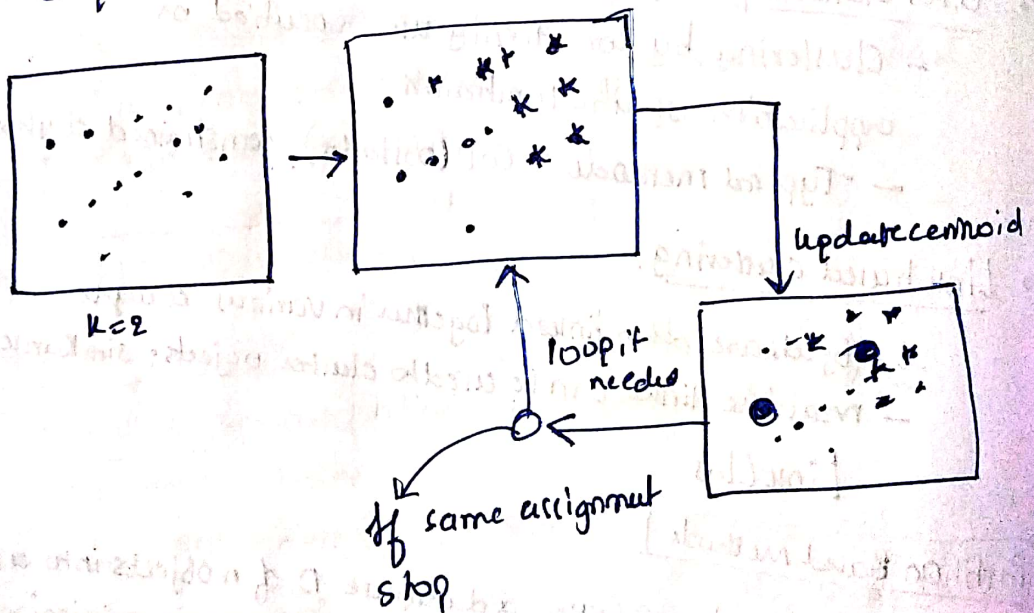
- Given k , find a partition of k clusters that optimizes the chosen partition criteria

- global optimal :- exhaustively enumerate all partitions
- Hierarchical partitioning : Each cluster is represented by the centre of cluster or one of objects (k-means, k-medoids)
- k-means : each cluster is represented by center of cluster
- k-medoids or PAM (Partition around medoids) each cluster is represented by one of objects in the cluster.

k-means clustering method

4 steps

- Partition objects into k nonempty subsets
- Compute seedpoints as the centroids of clusters of the current partitioning (centroid is the centre, mean of cluster)
- Assign each object to cluster with nearest seedpoint
- go back to step 2, when the assignment does not change stop.



Advantages:-

- efficient : $O(tkn)$, # of objects $\rightarrow n$, $k \rightarrow$ # of clusters, $t \rightarrow$ # of iterations.

mostly $k, t \ll n$

PAM ($O(k(n-k)^2)$)

CLARA : $O(k^2 + k(n-k))$

Disadvantages:

- terminates often at local optimal
- sensitive to noisy data & outliers
- NOT suitable to discover clusters with non-convex shapes

K-means
Problem?

Variations of k-means

- most of the variants of k-means differ in
 - selection of the initial k means
 - dissimilarity calculations
 - strategies to calculate cluster means
- Handling categorical data: k-modes
 - replacing means with modes
 - A mixture of categorical & numerical data: k-prototype method

Why k-medoids?

- k-means algorithm is sensitive to outliers
- But instead of taking the mean, the medoids can be used (which are centrally located object in a cluster)

PAM: A typical k-medoids problem

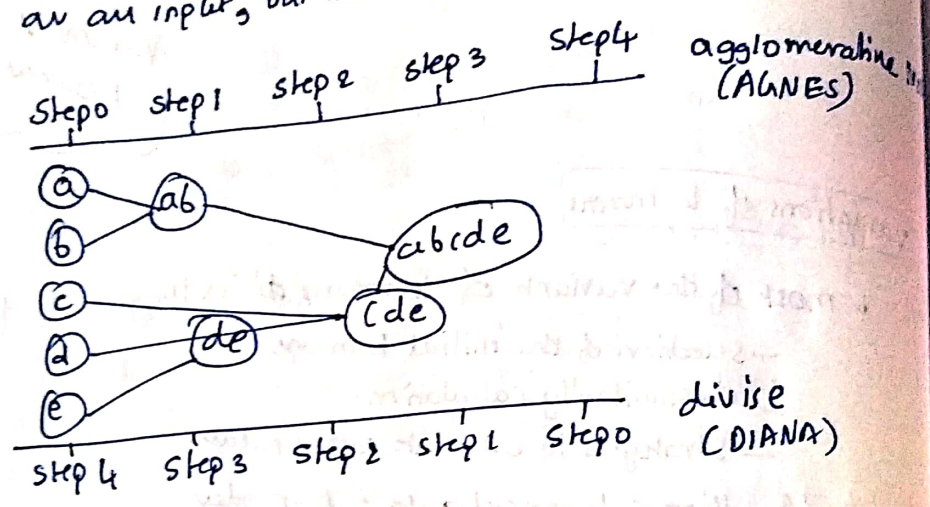
- Start from an initial set of medoids & iteratively replace one of medoids by one of non-medoids iff it improves the total distance of the resulting clustering
- PAM works effectively for small data sets; but does not scale well for large data sets (due to computational complexity)

→ Improvements on PAM

- CLARA: PAM on samples
- CLARANS: Randomized re-sampling

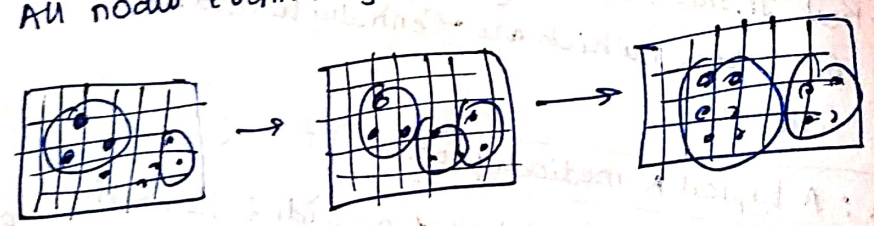
Hierarchical methods

- use distance matrix as a clustering criteria.
- This method does not require number of clusters k as an input; but needs a terminating condition



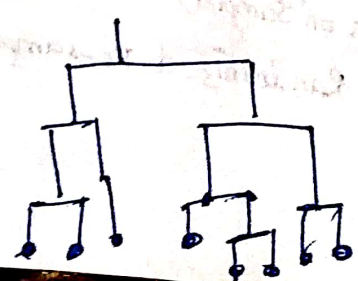
Agnes (Agglomerative nesting)

- Implemented in statistical package SPSS
- use single link method & dissimilarity matrix
- merges nodes with least dissimilarity.
- goes on a non-descending fashion
- All nodes eventually belong to same cluster



Dendrogram: How clusters are merged

- Decompose data objects into several levels of nested partitioning (tree of clusters), called a dendrogram
- A clustering of data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster



DIANA (divisive Analysis)

- inverse order of Agnes
- eventually each node forms a cluster of its own.

Distance between clusters

- single link: smallest distance between an element in one cluster and another element in another cluster.

$$\text{dist}(k_i, k_j) = \min(t_{ip}, t_{jw})$$

- complete link: longest distance between an element in one cluster & an other element in another cluster.

$$\text{dist}(k_i, k_j) = \max(t_{ip}, t_{jw})$$

- Average: Average distance between an element in one cluster and element in another cluster $\text{dist}(k_i, k_j) = \text{avg}(t_{ip}, t_{jw})$

- Centroid: distance between two centroids in two clusters

- Median: distance between two medoids in two clusters.

- centroid: $c_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$

- Radius: Square root of average distance from any point in cluster to its centroid.

$$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$

- Diameter: Square root of average mean squared distance between all pairs of points in cluster.

$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jw})^2}{N(N-1)}}$$

Extensions to Hierarchical clustering

- major weakness in agglomerative clustering:
 - can never undo what's done previously
 - Do not scale well

• Integration of hierarchical & distance based clustering.

① **Birch (1996):**

(Balanced iterative reducing & clustering using hierarchies)

• Incrementally construct a CF (clustering feature) tree, a hierarchical datastructure for multiphase clustering.

Phase 1: - Scan DB to build an initial in-memory CF tree (A multi-level compression of data that tries to preserve the inherent clustering structure of the data)

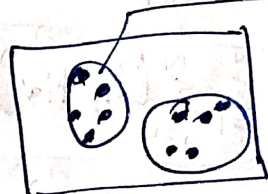
Phase 2: use an arbitrary clustering algorithm to cluster the leaf nodes of CF tree

- scales linearly
- handles only numeric data
- clustering feature vector in Birch

$CF = (N, L, SS)$

- $N \rightarrow$ number of data points
- $L \rightarrow$ Linear sum of N points $\sum_{i=1}^N X_i$
- $SS \rightarrow$ square sum of N points $\sum_{i=1}^N X_i^2$

$CF = (5, (16, 36), (54, 190))$



- (3, 4)
- (2, 6)
- (4, 5)
- (4, 7)
- (3, 8)

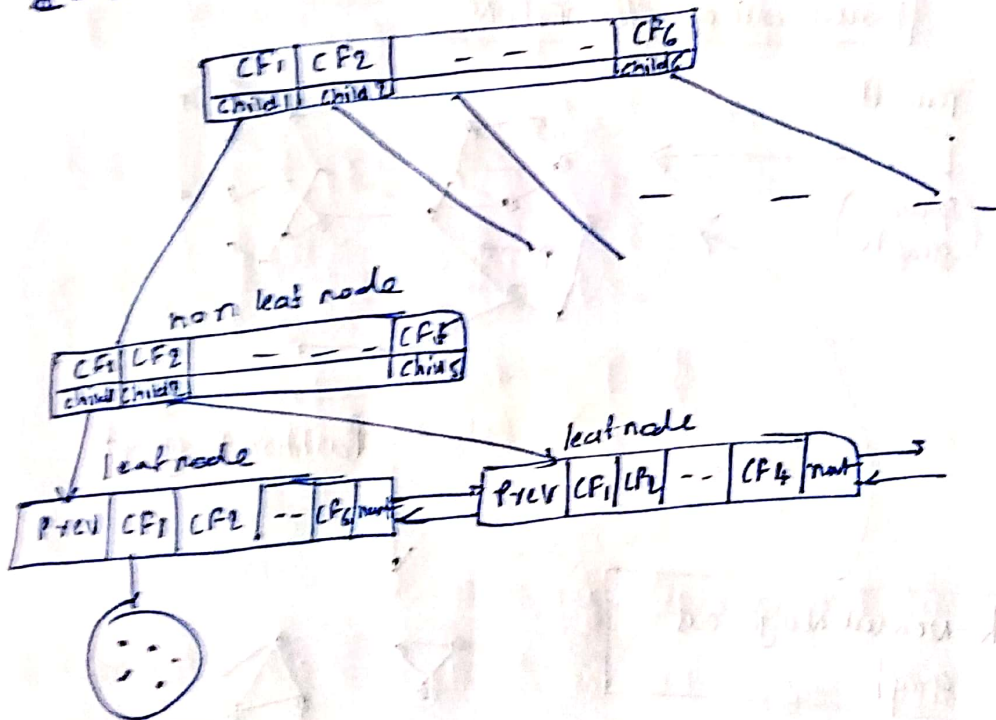
• A CF tree is a high balanced tree that stores the clustering features for a hierarchical clustering.

- A non-leaf node in a tree has children
- The non-leaf node stores the sums of CFs of their children.

• A CF tree has two parameters

- Branching factor: Max # of children
- Threshold: Max diameter.

$B = 7$
 $t = 6$



Birch Algorithm :

• cluster diameter :

$$\sqrt{\frac{1}{n(n-1)} \sum (x_i - x_j)^2}$$

- For each point in input
 - Find closest leaf entry
 - Add point to leaf entry & update CF
 - If entry diameter > max-diameter then split it.

• Algorithm: $O(n)$
 • Sensitive to insertion order of data points

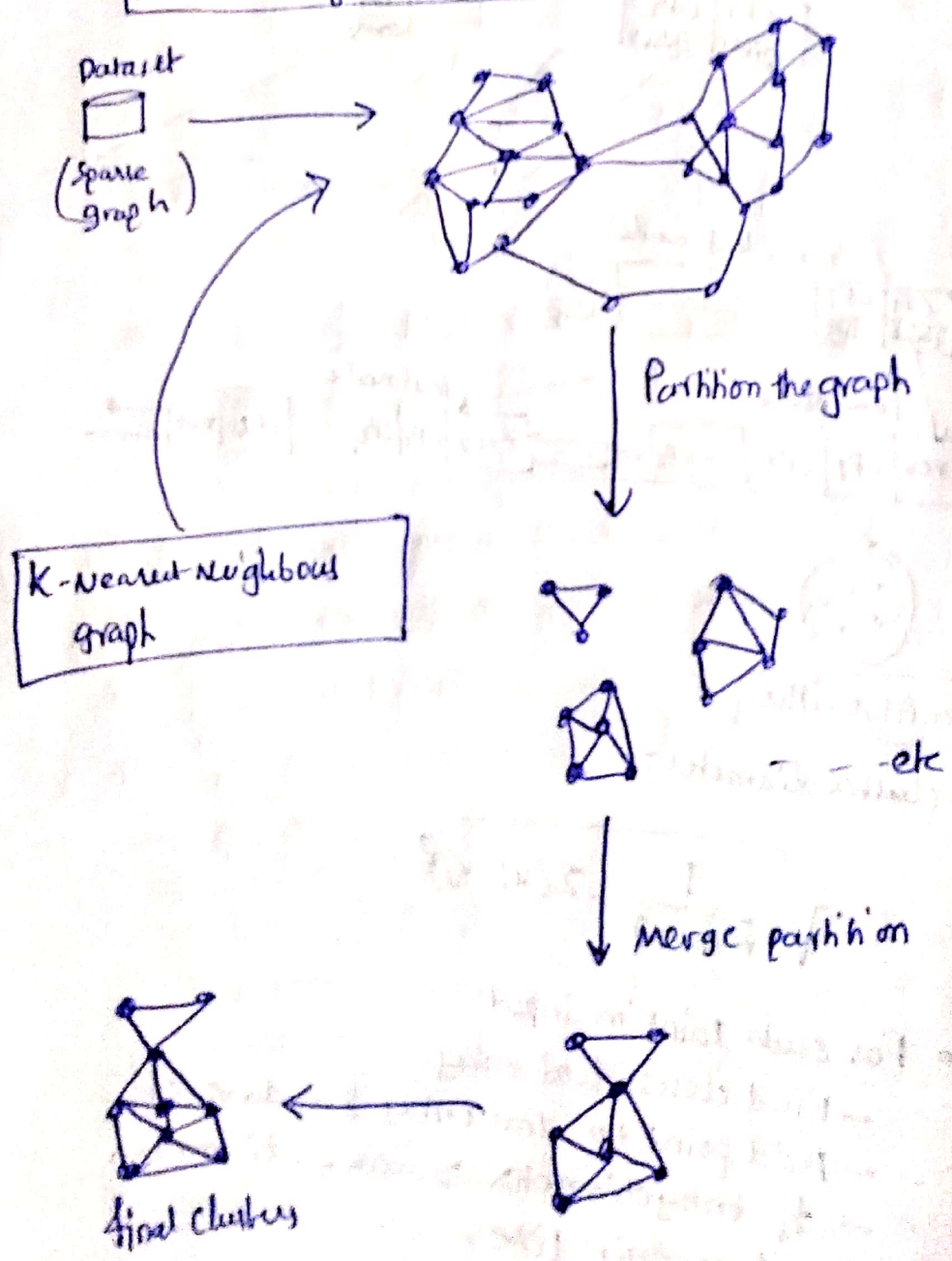
CHAMELEON (Hierarchical clustering using dynamic modelling)

- measures the similarity based on a dynamic model
 - Two clusters are merged only if H (Interconnectivity & Proximity) between two clusters are high;

Graph-Based, two phase Algorithm:

- Use a graph-partitioning algorithm: - cluster objects into a large number of relatively small sub-clusters
- Use an agglomerative hierarchical clustering algorithm: find the genuine clusters by repeatedly combining these sub-clusters.

Framework of CHAMELEON



Probabilistic Hierarchical clustering

- Algorithmic hierarchical clustering
 - Nontrivial to choose a good distance measure
 - hard to handle missing values
 - Optimization goal not clear; heuristic, local search

• Probabilistic hierarchical clustering:

- use probabilistic models to measure distance between clusters
- Generative model: set of data objects to be clustered as a sample of underlying data generation mechanism to be analyzed.
- easy to understand, same efficiency as Algorithmic hierarchical clustering, but can handle partly observed data

eg: Gaussian distribution or Bernoulli distribution

Generative model :-

- let 1-D points $x = \{x_1, \dots, x_n\}$ and assume are generated by gaussian distribution

$$N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)}$$

- Probability that point $x_i \in X$ is generated by model

$$P(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x_i - \mu)^2}{2\sigma^2}}$$

- The maximum likelihood (Task of generative model learning)

$$N(\mu_0, \sigma_0^2) = \arg \text{Max} \{ L(N(\mu, \sigma^2) : X) \}$$

A Probabilistic Hierarchical Clustering Algorithm

$$Q\{c_1, \dots, c_m\} = \prod_{i=1}^m P_i(L_i)$$

Maximum likelihood

• distance between clusters c_1 and c_2

$$\text{dist}(L_1, L_2) = -\log\left(\frac{P(L_1 \cup L_2)}{P(L_1)P(L_2)}\right)$$

• Algorithm :- progressively merge points & clusters

Input: $D = \{0, \dots, 0\}$

Output: A hierarchy of clusters

Method:

create a cluster for each object $c_i = \{0_i\}$ $1 \leq i \leq n$

for $i=1$ to n

 find pair of clusters c_i & c_j such that

$$c_i, c_j = \arg \max_{i \neq j} \left\{ \log\left(\frac{P(L_i \cup L_j)}{P(L_i)P(L_j)}\right) \right\}$$

 if $\log\left(\frac{P(L_i \cup L_j)}{P(L_i)P(L_j)}\right) > 0$, merge c_i & c_j

Density Based clustering method

• clustering based on the density (local cluster criterion) such as density connected points

• major features

— one scan

— handle noise

— discover clusters of arbitrary shape

— requires density parameters as terminating condition

Basic concepts

- Eps (Parameter) :- Maximum radius of the neighbourhood

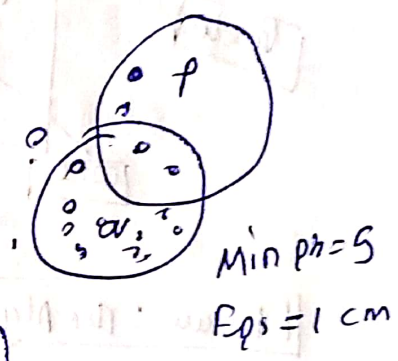
- MinPts (Parameter) :- minimum number of points in an Eps neighbourhood of that point

$N_{Eps}(P) = \{q \text{ belongs to } D \mid \text{dist}(P, q) \leq Eps\}$

Direct density reachable : A point p is directly density reachable from a point q wrt Eps, minpts if

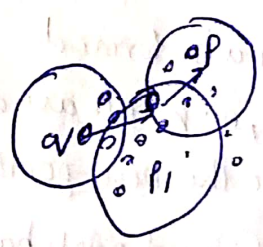
- p belongs to $N_{Eps}(q)$
- core point condition

$$|N_{Eps}(q)| \geq \text{MinPts}$$



Density - Reachable & density connected

Density reachable :-



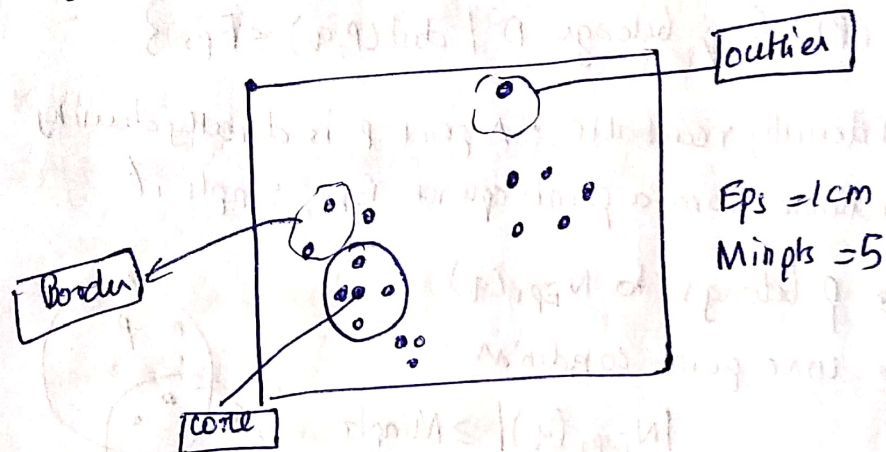
ccc A point p is density reachable from a point q wrt Eps, minpts if there is a chain of points p_1, p_2, \dots, p_n --- $p_1, p_1 = q, p_n = p$ such that p_{i+1} is directly density reachable from p_i \forall

Density-Connected

ccc A point p is density connected to a point q wrt Eps, minpts if there is a point o such that both p & q are density reachable from o wrt Eps & minpts $\forall \forall$

DBSCAN (Density-Based spatial clustering of Applications with noise)

- Relies on a density based notion of cluster. A cluster is defined as maximal set of density connected points
- Discovers clusters of arbitrary shape in spatial database with noise



DBSCAN: The Algorithm

- Arbitrary select a point P
- Retrieve all points density-reachable from P w.r.t Eps & $minpts$
- If P is a core point cluster is formed
- If P is a border point, no points are density-reachable from P & DBSCAN visit the next point of the database
- Continue the process until all of points have been processed

OPTICS: A cluster-ordering method

(Ordering points to identify clustering structure)

— good for both automatic & interactive cluster analysis, including finding intrinsic clustering structure

— can be represented graphically or using visualization techniques

- ophiu \rightarrow adder an an extension to DBSCAN

\rightarrow Index based :-

$K =$ Number of dimensions

$N = 20$

$P = 75\%$

$M = N(1-P) = 5$

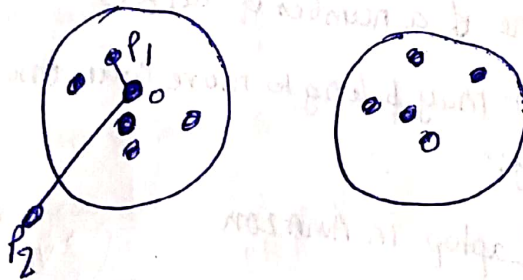
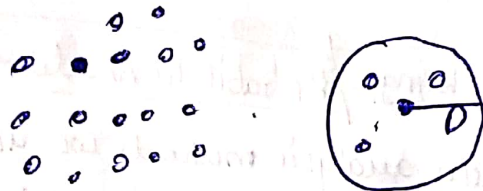
Complexity $O(N \log N)$

\rightarrow Core distance :-

min eps st point is core

\rightarrow Reachability distance :-

$\text{Max}(\text{Core-distance}(o), d(o, p))$



DENCLUE: using statistical density functions

$f_{\text{Gaussian}}(x, y) = e^{-\frac{d(x, y)^2}{2\sigma^2}}$ Increase of γ on x

$f_{\text{Gaussian}}^D(x) = \sum_{i=1}^N e^{-\frac{d(x_i, x)^2}{2\sigma^2}}$ total influence on x

$\nabla f_{\text{Gaussian}}^D(x, x_i) = \sum_{i=1}^N (x_i - x) \cdot e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$ gradient of x in direction of x_i

Major features of DENCLUB

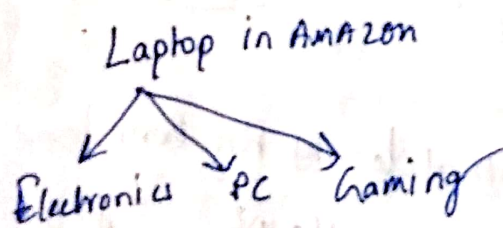
- solid mathematical foundation
- good for datasets with large amounts of noise
- needs large parameters
- faster than DBSCAN

Grid-Based Clustering method

- using multi-resolution grid data structure
- several interesting methods
 - STING (Statistical Information grid approach)
 - wavecluster (A multi-resolution clustering approach)
 - CLIQUE (Both grid-based & subspace clustering)

Model Based clustering : Probabilistic model-based clustering

- In all the cluster analysis methods we use
 - data object for each to be assigned to only one of a number of clusters
- Sometimes it may belong to more than one cluster
eg:-



III Model Based clustering Attempts to optimize the fit between the data & some mathematical model by using statistical & AI approach III

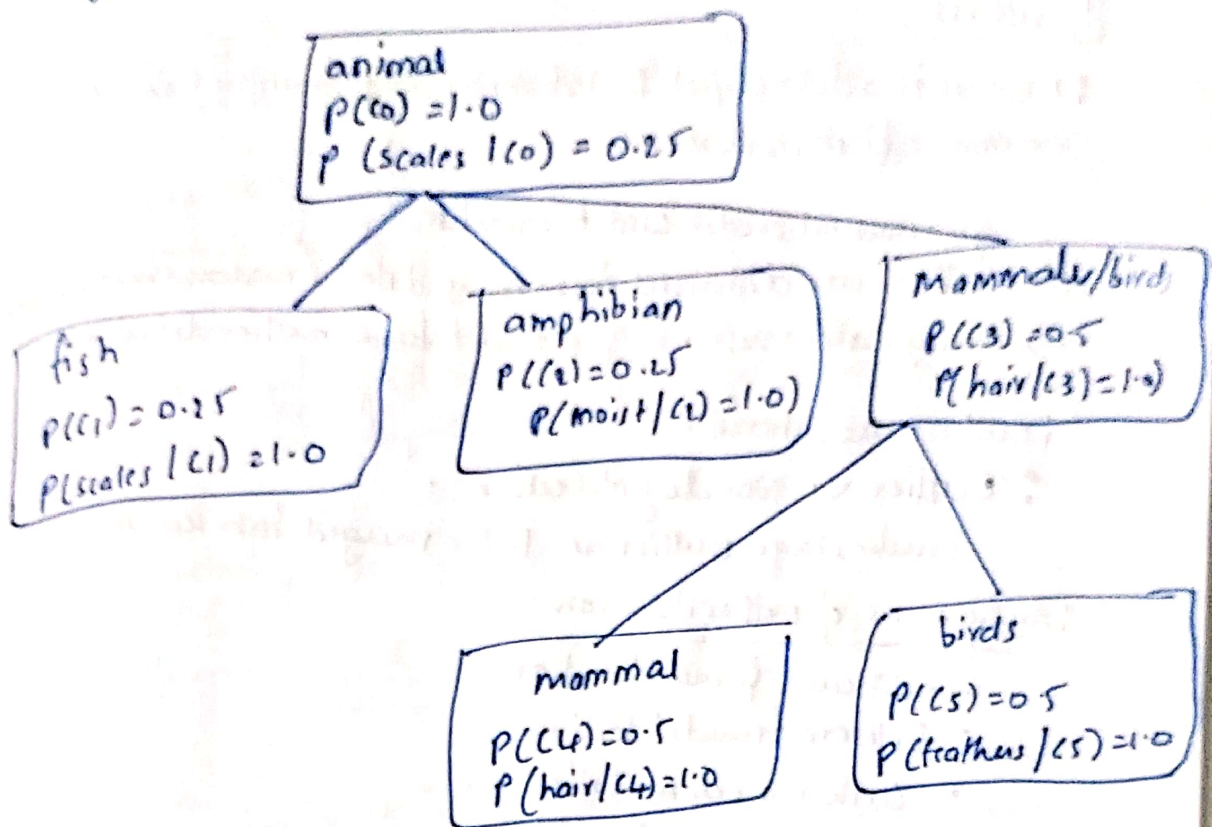
— Conceptual clustering

(finds characteristic description for each concept/class)

— COBWEB (Fisher 7)

- (Incremental conceptual learning)
- creates hierarchical clustering in form of classification tree
- each node refers to a concept &

contains probabilistic description for that object.



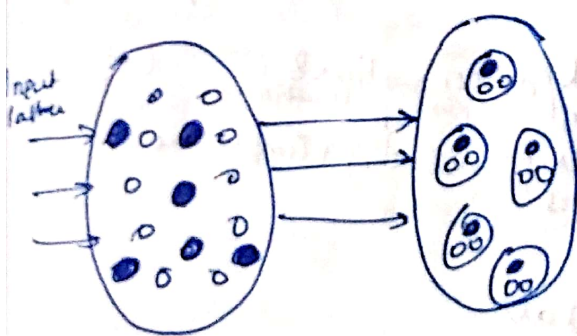
• limitations of COBWEB

- not suitable for clustering large database
- assumption that attributes are not co-related may not be true for all

— Class IT

(an extension of COBWEB) but as failure as COBWEB

— Competitive learning



* hidden layers may also be present

— Self-organizing feature maps A A A

- Clustering is performed by objects competing for a class
- Winner & its neighbours learn by adjusting their weights
- Similar to processing in human brain

Outliers

• A outlier is a data object that deviates significantly from the normal object in dataset

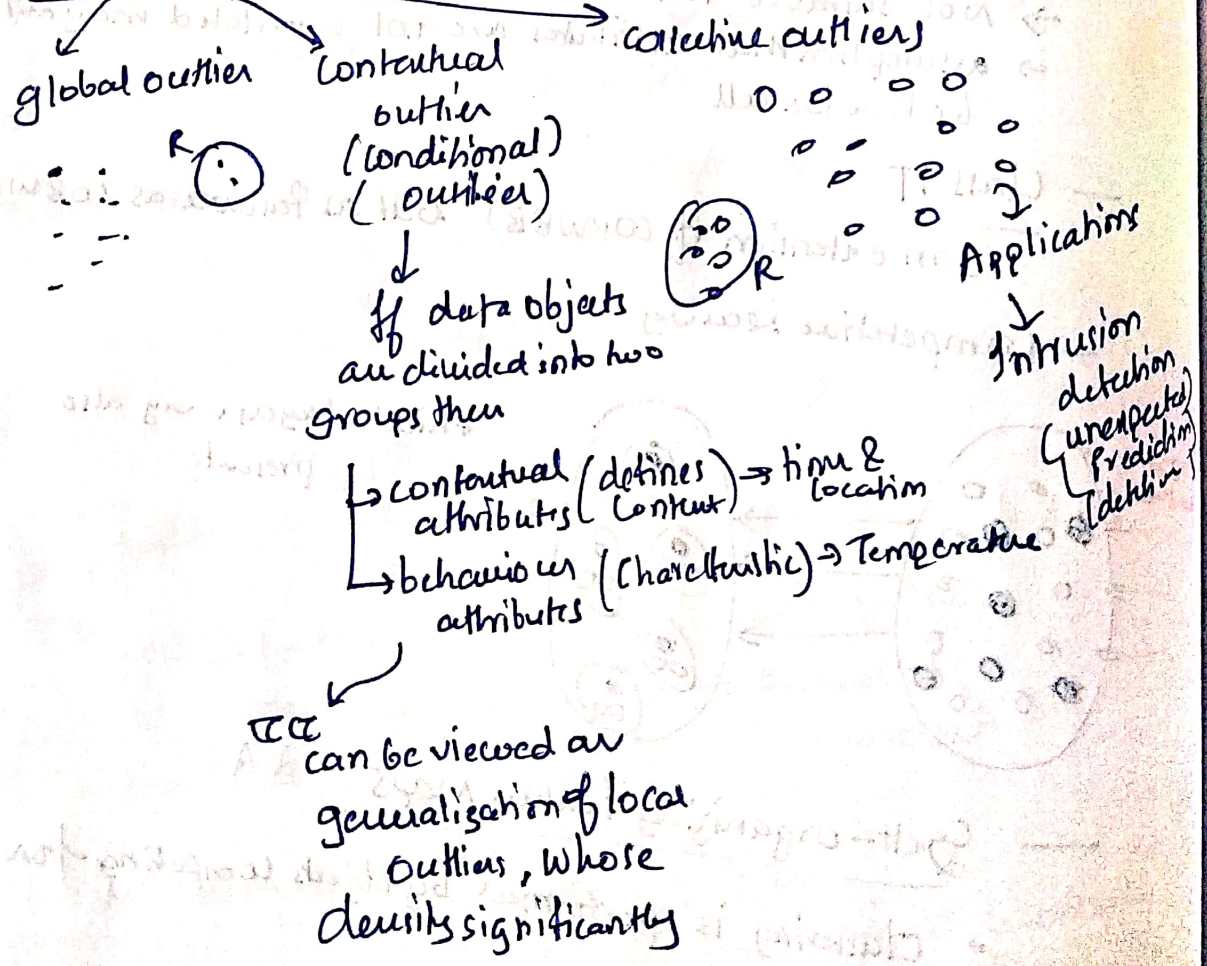
- eg: unusual credit card transaction
- Outliers are different from noisy data (random error, variance)
- noisy data must be removed before outlier detection

- outliers are interesting
- Outlier vs Novelty detection: -
early stage outlier but later merged into the model

Applications of outlier detection: -

- Credit card fraud detection
- Telecom fraud detection
- Customer segmentation
- Medical analysis

Types of outliers



Challenges of outlier detection

- modeling normal objects & outliers properly
- Application-specific outlier detection
- Handling noise in outlier detection
- Understandability

Outlier detection Methods

Based on user labelled examples of outliers

Supervised methods

- modelling outlier detection as a classification problem
- Method for learning a classifier for outlier detection effectively
- challenges: -
 - * Imbalanced classes
- Sol: make artificial outlier
- * Catch as many outliers as possible

Unsupervised Methods

- Already assuming normal objects are clustered into groups each having distinct features
- An outlier is expected to be far away from the clusters
- challenges:
 - * In virus detection it may have high false positive rate
 - * very costly

Semi supervised methods

If labelled normal objects are available

- Then train model using unlabelled data objects by using labelled examples
- Those not fitting lines are outliers

If labelled outlier objects are available

- Then take help of unsupervised methods

Based on assumptions between normal & outlier data

Statistical methods (or) Model Based methods

- find gaussian distribution to model the normal data
- alternative
- * Parametric vs non-Parametric

mean, variance, Grubbs Test

$f(x,0)$

Proximity Based methods

- An object is outlier if the nearest neighbours of object are far away

Clustering Based methods

- normal data belong to large & dense clusters
- outliers belong to small or sparse clusters or don't belong to any clusters



distance based

- challenges:
 - hard to find outliers close to each other

grid cell based algorithm
Minkowski algorithm

Local outlier factor (LOF)