

Unit

Clustering

- Intra Cluster → minimized (btw objects)
- Inter Cluster → maximized (btw clusters)

Cluster Analyze:-

summarization: reduce the size of large dataset.

Requirements of clustering

- Scalability
- Ability to deal with diff type of attributes
- Able to deal with noise and outliers
- High dimensionality
- Incorporation of user-specified constraints
- Insensitive to order of input record.

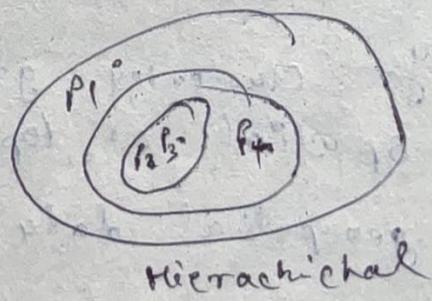
Type of Clustering

→ Partitional clustering

A division of data objects into non-overlapping subsets.
division (clusters)

→ Hierarchical clustering

A set of nested clusters organized as a hierarchical tree



- Exclusive v/s non-exclusive
 - In non-exclusive clustering point may belong to multiple clusters.
- Fuzzy clustering (non exclusive)
 - Partial v/s complete
 - ↓
 - ↳ weight & work
sum = 1
 - ↳ Could be 'border' points

Types of Clusters

- Well separated \Rightarrow K-means
- Prototype based
- Contiguity based
- Density based
- Described by an objective function

Prototype:- Cluster refers to a set of representative points that summarize a cluster of data points in clustering algorithms.

- in K-means prototype are the centroids of each cluster.
- in K-medoid prototypes are actual data points that best represent the cluster.

Contiguity:- Clustering groups data points based on their spatial or topological adjacency. (closeness or connection)

- It is used in geospatial data to identify regions that
- Clusters are formed based on proximity with adjacent data points sharing similar characteristics or closeness.

Density Based :-

ip. a dense region of points which is separated by low-density regions, from other regions of high density.

- density based clustering groups data points that are closely packed together in high-density regions and separate them from sparse region.
- here we use K-means, DBSCAN, OPTICS

Clusters defined by an objective functions

finds clusters that minimize or maximize an objective function.

Enumerate all possible ways of dividing the points into clusters and evaluate the goodness of each potential set of clusters by using the given objective function. (NP Hard)

Can have global or local objectives.

→ hierarchical clustering algorithms typically have local objectives

→ partitional algo typically have global objective

Characteristics of the Input Data and types

Good Clustering :-

- high intra-class similarity
- low inter-class similarity.

Quality → similarity
implementation

Partition clustering :-

$$E = \sum_{k=1}^K \sum_{x \in C_k} \text{Distance}(x, m_k)$$

total error
 total no. of clusters

$$\text{Distance}(x) = \sum_{d=1}^D (x_d - m_{kd})^2$$

Solve K means examples.

Strengths

- simple: easy to understand
- $O(tkn)$
- iteration
- cluster num

Weakness

- Algo is only applicable if the mean
- user need to specify k is defined
- outliers sensitive

Preventions :-

- remove outliers first.
- take some subsets of point then chance of selection of outlier is very low.

Weakness :- algo is not suitable for discovering clusters that are not hyper-ellipsoids.

(x-0) \neq 0

K medoids :-

PAM algo represent some data points.

Solve K medoids examples -

Properties :-

$$\Theta(k(n-k)^2)$$

for large val of n & k computation become very costly.

Advantage :- K medoid is more robust than K-mean in the presence of noise and outliers.

Disadvantage :-

→ K medoid is more costly than K-mean.

- K mean, Kmedoid requires their given to specify k,
- it does not work well for large dataset.

PAM (Partitioning Around Medoids) :-

- that are centrally located in clusters, finds medoid
- The goal of the algo is to minimize the avg. dissimilarity of obj. to their closest selected obj.

→ PAM works effectively for small datasets, but does not scale well for large datasets.

K-medoids properties :-

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

Advantage :- K-medoids method is more better than k-means in the presence of noise and outliers.

Disadvantage :-

- k-medoids is more costly than k-mean.
- k-medoids requires the user to specify k.
- does not scale well for large datasets

CLARA (Clustering Large Application) :-

uses a sampling based method to deal with large datasets.

→ chosen medoids will likely be similar to what would have been chosen from the whole dataset.

- draw multiple samples of the dataset
- Apply PAM (k-medoid) to each sample
- return the best clustering

Properties :-

$$O(tk(s-k)^2 + k(n-k))$$

Diagram illustrating the time complexity components:

- t : Number of samples drawn.
- k : Number of clusters.
- s : Size of each sample.
- n : Total number of objects.
- $\max_{\text{iteration}} \text{ required size}$: Maximum iteration size required for PAM.
- No. of Obj. : Total number of objects.
- No. of Clusters : Total number of clusters.

problem

- The best K medoids may not be selected during the sampling process. In this case, CLARA will never find the best clustering.
- If the sampling is biased we cannot have a good clustering.

CLARANS ("Randomized" CLARA)

(A clustering large application based on randomized search)

- The clustering process can be presented as searching a graph where every node is a potential solution that is a set of K medoids.
- Two nodes are neighbours if their sets differ by only one medoid.
- Each node can be assigned a cost that is defined to be the total dissimilarity between every object and the medoid of its cluster.
- The problem corresponds to search for a minimum on the graph.
- At each step, all neighbours of current node node are searched; the neighbour of current node node which corresponds to the deepest descent in cost is chosen as the next solution.

Graph Abstraction -

- Every node has $K(n-1)^c$ adjacent nodes
- At each step, CLARANS draws sample of neighbours to examine.
- Note that CLARA draws a sample of nodes at the beginning of search.

- therefore, CLARANS had the benefit of not confining the search to a restricted area.
- If the local optimum is found, CLARANS starts with a new randomly selected node in search for a new local optimum. The number of local optimums to search for is a parameter.
- It is more efficient and scalable than both PAM and CLARA; returns higher quality clusters.

Disadvantage:- for large value of n and k , examining $K(n-K)$ neighbours is time consuming.

Exam

Hierarchical Clustering :-

- Agglomerative vs Divisive,
- two sequential clustering strategies for constructing a tree of clusters
- Agglomerative: a bottom-up strategy
 (atomic) → Initially each data object is in its own cluster
 → Then merge these atomic clusters into larger and larger clusters.

Divisive: a top-down strategy:-

- Initially all objects are in one single cluster
- Then the cluster is subdivided into smaller and smaller clusters.

single link :-

Complete link

Avg link

$$\min \{ \dots \}$$

$$\max \{ \dots \}$$

$$\text{Avg} = \frac{\dots}{n}$$

Single-link

- Can find irregular-shaped clusters
- Sensitive to outliers, suffers the so-called chaining effects
 - In order to merge two groups, only need one pair of points to be close, irrespective of all others.

Average-link and Centroid distance

- Robust to outliers
- Tend to break large clusters

ACINES (Agglomerative Nesting) :-

→ Single link
→ least dissimilarity

UPGMA (un-weighted Pair Group Method Average) :-

→ Average link approach

DIANA :- (Divisive Analysis)

BIRCH :-

mathoh huddar solved example youtube

CURE (Clustering Using Representative) :-

Drawbacks of Traditional Clustering Algorithms :-

- Centroid-based approach (using d_{mean}) considers only one point as representative of a cluster - the cluster centroid.
- All-points approach (based on d_{min}) makes the clustering algorithm extremely sensitive to outliers.
- Both of them can't work well for non-spherical or arbitrary shaped clusters.

CURE : ~~Approach~~ Approach

- CURE is positioned between centroid based and all point extremes.
- A constant number of well scattered points is

DBSCAN Density-Based Clustering :-

Basic idea:-

- clusters are dense regions in the data space, separated by regions of lower object density
- A cluster is defined as a maximal set of densely connected points.

Each cluster has a considerably higher density of points than outside of the cluster.

DBSCAN :-

Advantages:-

- Clusters can have ~~and~~ arbitrary shape and size.
- No. of clusters is determined automatically.
- Can separate clusters from surrounding noise.

Disadvantages:

- Input parameters may be difficult to determine.
- In some situations very sensitive to input parameter setting.

Association Rule Mining

Study of "what goes with what"

- "Customers who bought X also bought Y"
- what symptoms go with what diagnosis

Transaction-based or event-based

Also called "Market Basket Analysis" and "affinity analysis".

Support % —

the supports of an itemset X is the percentage of transactions in the transaction database T that contain X .

→ support of the rule $X \Rightarrow Y$ in the transaction database T is the support of the items set $X \Rightarrow Y$ in T .

→ probability that a transaction contains $X \cup Y$.

Confidence % —

of the rule $X \Rightarrow Y$ in the transaction database T is the ratio of

the no. of transactions in T that contain $X \Rightarrow Y$ to the no. of transactions that contain X in T.

Conditional probability that a transaction having X also contains Y.

Support :-

$$\frac{\text{kitni transaction} \\ (\text{me ek saath done a}) \\ \text{rhe hain}}{\text{total transaction}} \times 100\%$$

Confidence :-

$$\frac{\text{(dono left + right)} \\ (\text{kitni bar aa rha hai})}{(\text{left side} \\ \text{wale ka count})} \times 100\%$$

The Apriori Algorithm :-

Maximum Frequent Set :-

maximum length wale combination jiska ~~support~~ support, minimum support (threshold) se jyada (\geq) ho to vo maximum frequent set hogा.

Example :-

Border Frequent Set :-

maximum length wale combination jiska support, minimum support (threshold) se kam hogा lekin uske sare proper subset ka support min support se jyada hogा to vo border freq. set hogा.

Ex.

	<u>support</u>
{A}	4 ✓
{B}	4 ✓
{C}	4 ✓
{D}	3 ✓
{A, B}	3 ✓
{B, C}	3 ✓
{A, C}	3 ✓
{A, D}	2 ✓
{B, D}	2 ✓
{C, D}	2 ✓
{A, B, C}	2 ✓
{A, B, D}	1 ✗
{A, C, D}	1 ✗
{B, C, D}	1 ✗
{A, B, C, D}	0 ✗

T_id	item
T1	{A, B, C}
T2	{A, B, D}
T3	{A, C, D}
T4	{B, C, D}
T5	{A, B, C}

Given →

min support = 2

Maximum Frequent Set :-

{A, B, C} because its all subsets are frequent itemset (support ≥ 2).

because ~~its~~ length is maximum among all frequent itemset.

Border Frequent Set :-

maximum length of itemset that is not frequent (support < 2) and all subsets are frequent set. so here we check 4 combinations

{A, B, C, D}, {B, C, D}, {A, C, D}, {A, B, D}.

so our border sets are {A, B, D}, {A, C, D}, {B, C, D}.

Challenger of Frequent Pattern Mining

challenges:-

- Multiple scans of transaction database
- Huge no. of candidates
- Tedium workload of support counting for candidates.

Improving Apriori :

- Reduce passes of transaction database scans
- Shrink no. of candidates
- Facilitate support counting of candidates

Partition Algorithm :-

The partition Algo. is based on the observation that the frequent sets are normally very few in no. compared to the set of all item sets.

As a result, if we partition the set of transactions to smaller segments such that each segment can be accommodated in the main memory, then we can compute the set of frequent sets of each of these partition.

Since each partition can fit in the main memory there will be no additional disk I/O for each partition after loading the partition into the main memory.

FP-tree Growth Algorithm :

why FP-tree and not Apriori?

→ lots of frequent patterns

- Big set of items

- low minimum support threshold

→ long patterns

FP-Tree Construction :-

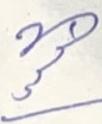
→ FP tree is constructed using 2 passes over the dataset

Pass 1:

- Scan data and find support for each item.
- Discard infrequent items.
- Sort frequent items in decreasing order based on their support.
- Use this order when building the FP-tree, so common prefixes can be shared.
- In FP-tree nodes correspond to items and have a counter.

Pass 2:

1. FP Growth reads 1 transaction at a time and maps it to a path.
2. Fixed order is used, so paths can overlap when transactions share items (when they have the same prefix).
 - In this case, counters are incremented.
3. Pointers are maintained between nodes containing the same item, creating singly linked lists (dotted lines).
 - The more paths that overlap, the higher the compression. FP-tree may fit in memory.
4. Frequent itemsets extracted from the FP-tree



Best Case:-

Same set of items

Worst Case:- every transaction has unique set of items

Measure for Association Rules :-

- Confidence ($X \rightarrow Y$) should be sufficiently high
To ensure that people who buy X will more likely buy Y than not buy Y
- Confidence ($X \rightarrow Y$) $>$ support(Y)
Otherwise rule will be misleading bcoz having item X actually reduces the chance of having item Y in the same transaction.

Statistical Relationship btw X & Y :-

$$\text{Confidence } (X \rightarrow Y) = \text{support}(Y)$$

is equivalent to :

- $P(Y/X) = P(Y)$
- $P(X, Y) = P(X) \times P(Y)$ { X and Y are independent }

If $P(X, Y) > P(X) \times P(Y)$: X & Y are positively correlated

If $P(X, Y) < P(X) \times P(Y)$: X & Y are negatively correlated

$$\text{Lift} = \frac{P(Y/X)}{P(Y)}$$

$$\text{Interest} = \frac{P(X, Y)}{P(X) P(Y)}$$

$$PS = P(X, Y) - P(X) P(Y)$$

PS: pairwise similarity
 ≥ 0 no association
 > 0 +vely
 < 0 -vely

$$\phi - \text{coefficient} = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1-P(X)]P(Y)[1-P(Y)]}}$$

Ex:-

Association Rule:

$\boxed{\text{Tea} \rightarrow \text{Coffee}}$

confidence

$$\begin{aligned} \text{Interest} &= \frac{P(X, Y)}{P(X) \cdot P(Y)} \\ &= \frac{150}{1000} \\ &= \frac{200}{1000} \times \frac{800}{1000} \end{aligned}$$

$$\left\{ \begin{array}{l} X = \text{Tea} \\ Y = \text{Coffee} \end{array} \right.$$

	Coffee	Coffee	
Tea	150	80	200
Total	650	150	800
	800	200	1000

$$= \frac{.15}{.2 \times .8} = .9375 < 1 \quad \left\{ \text{therefore negatively associated} \right.$$

$$\text{Lift} = \frac{P(Y/X)}{P(Y)} = \frac{\frac{P(Y|X)}{P(X)}}{\frac{P(Y)}{P(X)}} = \frac{P(Y|X)}{P(X) \cdot P(Y)}$$

$$PS = P(X, Y) - P(X) \cdot P(Y)$$

$$= .15 - (.2 \times .8)$$

$$= .15 - .16 = -0.01 < 0 \quad \left\{ \text{negatively associated} \right.$$

Text Mining

Data Mining

- Identify data sets
- Select features
- Prepare data
- Analyze distribution

Text Mining

- Identify documents
- Extract features
- Select features by algo
- Prepare data
- Analyze distribution

Text Mining Applications :-

- email, newspaper, patent database, company's reports.

Why dealing with Text is Tough :-

- Abstract concepts are difficult to represent
- "Countless" combinations of subtle abstract relationships among concepts.
- Many ways to represent similar concepts
Eg. space ship, flying saucer, UFO
- Concepts are difficult to visualize
- High dimensionality
- Tens or hundreds of thousands of features

Why dealing with Text is Easy :-

- High redundant data
most of the methods count on this property
- Just about any simple algo can get "good" results for simple tasks:
 - pull out "important" phrases
 - Find "meaningfully" related words
 - Create some sort of summary from documents

Word properties :-

Homonymy :— same form, but diff. meaning.

Polysemy :— same form, related meaning.

Synonymy :— different form, same meaning.

Hipponymy :— one word denotes a subclass of another.

power distribution :-

→ small no. of very frequent word.

→ big no. of low frequent word.

Stop-words :— are words that from non-linguistic view do not carry information.

ex— A, About, Above, Along, Already —

Stemming :-

different forms of the same word are usually problematic for text data analysis, because they have different spelling and similar meaning (eg. learning, learned, learning ...)

Phrase level :-

The main effect of using phrases is to more precisely identify sense.

Taxonomies / Thesaurus level :-

main function to connect different surface word forms with the same meaning into one sense (synonyms).

Ex. EuroWordNet.

Thesaurus has a

different surface word meaning into one sense

WordNet - database of lexical relations

wordNet is the most well developed and widely used lexical database for english
 → It consist from 4 databases (nouns, verbs, adjectives, and adverbs)

Category	Unique Form	Number of Senses
Noun	94474	116317
Verb	10319	22066
Adjective	20170	29881
Adverb	4546	5677

WordNet relations:-

Each wordnet entry is connected with other entries in the graph through relations.

Vector-Space Model Level :-

The most common way to deal with documents is first to transform them into sparse numeric vectors and then deal with them with linear algebra operations.

this way of forgetting about the structure doesn't harm efficiency of solving many related problems.

Typical tasks on vector-space-model are classification, clustering, visualization etc.

Word weighting

$$tfidf(w) = tf \cdot \log\left(\frac{N}{df(w)}\right)$$

The word is more important if it appears several times in a target document.

The word is more important if it appears in less documents.

Full-parsing :-

parsing provides maximum structural information per sentence.

- On the input we get a sentence, on the output we generate a parse tree.