

11-03-19

* Data Reduction

Data reduction techniques can be applied to obtain a reduced representation of the data set. Data reduction techniques can be applied to obtain ordered representation of data set i.e. much smaller in volume yet closely maintaining the integrity of the reduced data i.e. mining of K on the reduced data set should be more efficient yet produce the same analytical results.

Techniques are of 3 types -

① Dimensionality reduction - It is the process of reducing the no. of attributes (size - rows) under consideration.

② Numerosity reduction - These techniques replace the original data volume by alternative smaller forms of data representation.

③ Data compression techniques - These are applied so as to obtain a reduced or compressed representation of the original data. If the original data can be reconstructed from the compressed data without any information loss the data reduction is called lossless.

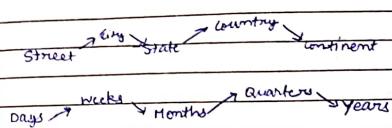
If instead we can reconstruct only an approximation of the original data, the data reduction is called lossy. Some of techniques are described below:

From category 2, a technique is used.

→ Data cube Aggregation - Rollup technique used

done previously
give example

→ Concept Hierarchy - The concept hierarchy allows the analysis of data at multiple abstraction levels



→ Clustering - Also a data reduction technique

In data reduction, the cluster representation of the data are used to replace the actual data. The effectiveness of this data technique depends upon these data's ratio.

→ Random Sampling - It can be used as a data reduction technique because it allows a large data set to be represented by a much smaller random data sample.

a) ~~**~~ SRS WOR - Simple random sample without replacement.

This is created by drawing i^{th} of N tuples from database where the probability of drawing any tuple is $\frac{1}{N}$.

b) ~~**~~ SRS WR - Simple random sample with replacement.

This is similar to the earlier method except that each time a couple is drawn from D - database, it is recorded & then replaced. After a tuple is drawn, it is placed back in database D so that it may be drawn again.

c) ~~**~~ Stratified Sample - If D is divided into mutually disjoint parts called strata, D is divided into

A stratified sample is generated by obtaining an SRS at each stratum.

(Simple random sample)

e.g. A stratified sample may be obtained from customer data where a stratum is created for each customer age group.

d) Cluster sampler - if the tuples in D are grouped into 'm' mutually disjoint clusters, then an SRS of 'k' clusters can be obtained.

e.g. Tuples in a database are usually retrieved a page at a time so that each page can be considered as a cluster. A reduced data representation can be obtained by applying SRSWOR to the pages resulting in a cluster sample of m tuples.

Attribute subset Selection :- Attribute subset selection reduces the data set size by removing irrelevant or redundant attributes(dimensions). The goal of the attribute subset selection is to find a minimum set of attributes such that the resulting probability distribution of data classes is as close as possible to the original distribution obtained using all attributes.

The best (& worst) attributes are typically determined using tests of statistical significance which assumes that attributes are independent of one another.

15-03-19

Data Transformation

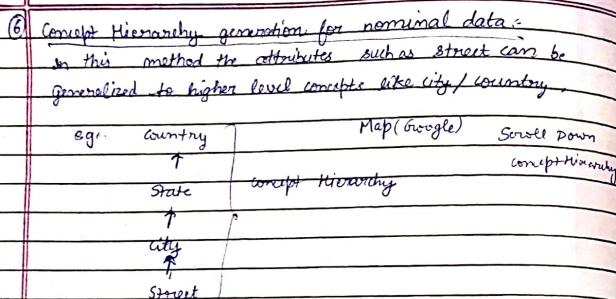
In DT, the data are transformed or consolidated into forms appropriate for mining.

Strategies for data transformation include the following :-

- ① Smoothing :- It works to remove noise from the data. The techniques include ~~mining~~, regression & clustering.
- ② Attribute Construction (Feature construction) :- In this method, new attributes are constructed & added from the given set of attributes to help the mining process.
- ③ Aggregation :- In this step, summary/aggregation operations are applied to the data. This step is typically used in constructing a data queue for data analysis at multiple abstraction levels.
- ④ Normalisation :- In this method, the attribute data are scaled so as to fall within a smaller range such as 0 to 1 or -1 to 1.

- ⑤ Discretization :- In this, the raw values of a numeric attribute are replaced by interval labels or conceptual labels.

e.g.		Make continuous value to discrete	label = Conceptual names
Table	(Name)	20, 30, 25, 23, ...	
Youth	20 - 30	frequency	
Middleage	31 - 40		

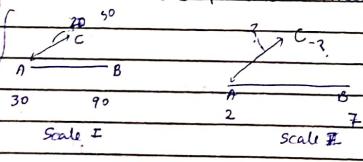


Data Transformation by Normalization

The measurement unit used can affect the data analysis e.g. changing measurement units from metres to inches or from Kg to pounds may lead to very different results. To help avoid dependence on the choice of measurement units, the data should be normalized or standardised.

This involves transforming the data to fall within a small or common range such as -1 to 1 or 0 to 1.

There are 3 prominent methods for data normalization



M-I Min Max normalization

$$v_i' = \frac{v_i - \text{Min}_A}{\text{Max}_A - \text{Min}_A} \cdot (\text{newMax}_A - \text{newMin}_A) + \text{newMin}_A$$

This method performs a linear transformation on the original data. Suppose Min_A & Max_A are the minimum & maximum values of an attribute A. Min-Max Normalisation maps a value v_i' to v_i'' in the range [NewMin_A, NewMax_A] by using the above mentioned formula.

$$\begin{aligned} \text{Min}_A &= 30 & v_i'' &= \frac{50-30}{90-30} (7-2) + 2 = 3.6 \\ \text{Max}_A &= 90 & \text{newMin}_A &= 2 \\ \text{newMax}_A &= 7 & \frac{20}{90} (5)+2 &= 5.6 \\ \Rightarrow A \text{ is } 5.6 \text{ away from } C & & \text{newMax}_A &= 7 \end{aligned}$$

M-II Z-Score Normalisation

In Z-score normalisation, the values of an attribute A are normalised based on the mean & standard deviation of A.

$$v_i'' = \frac{v_i - \bar{A}}{\sigma_A}$$

This method of transformation is used when actual min max are unknown.

M-III Normalisation by decimal scaling

This method normalises by moving the decimal pt. of values of attribute A. The no. of decimal pt. moved depends on the maximum absolute value of A.

The formula is given below:-

$$v_i'' = \frac{v_i}{10^j} \quad j \text{ is smallest integer such that } \text{Max}(|v_i'|) < 1$$

e.g. Min value = -0.985
Max. value = .975

$$\text{if } j=3 \rightarrow \frac{-0.985}{3} = -0.3283$$

$$\text{if } j=3 \rightarrow \frac{0.975}{3} = 0.325$$

18-03-19 Chapter: Mining Frequent Patterns & Associations

1. Frequent Patterns → These are patterns that appear frequently in a data set. e.g.: Frequent item sets or frequent subsequences fall in category of frequent patterns.

2. Frequent itemset → A set of items that appear frequently together in a transaction data set is called a frequent itemset. e.g. Milk & Bread

Frequent Sequential Subsequence Patterns → A subsequence such as buying first a PC, then a digital camera & then a memory card. If it occurs frequently in a shopping history database, it is a frequent sequential pattern.

Frequent itemset mining leads to the discovery of associations & correlations among items in large transactional or relational data sets. Massive amounts of data continuously being collected & stored, many industries are becoming interested in mining such patterns from their databases.

The discovery of interesting patterns can help in many business decision making processes such as catalog design, create marketing & customer shopping behavior analysis.

if $j=2$ X Not valid

if $j=4$ → Valid, Max. value

$$\frac{0.975}{4} = 0.24375$$

Chapter: Mining Frequent Patterns & Associations

• Association Rule:

$$\text{Head: } \{A, B\} \Rightarrow \{C, D\} \text{ consequent}$$

e.g. [Support = 2% /, Confidence = 60%]

It defines how strongly one itemset (or one set of items) is associated with / for other set of items.

{confidetness, antivisus} \rightarrow {handbk, pointers} positive

It is represented with an arrow symbol
left set of items: Head / Antecedent
right : Body / Consequent

Antecedent → Antecedent is associated with consequent

Transactions dataset

1.	A B C D	5.	C D B
2.	A C D	6.	B D
3.	A B E	7.	A B E
4.	A B C D E		

• Support 2% → It means that 2% of all the transactions under analysis shows that computer & A B C D activities are purchased together.

$$\text{A B C D occurs together in 2 transactions} \quad \text{ASCD will support} \\ \Rightarrow \frac{2}{10} \times 100 = 20.57 \% \quad \text{in 28370 total transactions} \\ \text{Support count = 2}$$

In 20.57% of all transactions, A B C D appear together.

Confidence 60% - It means that 60% of the customers who bought a computer & antivirus also bought harddisk & printer.

(A) (B)
(C) (D)

Those persons who bought A, B \Rightarrow 60% of those persons also bought C, D.

A, B appearing in 4 transactions. $\Rightarrow \frac{2}{4} \times 100 = 50\%$
C, D \longrightarrow 2 \Rightarrow confidence is 50%

(Ex: I/P of support & confidence is given by user.)

Association rules are considered interesting if they satisfy both a minimum support threshold & a min. confidence threshold.

$A \Rightarrow B$ A, B : itemset

$$\text{confidence}_{\text{of association rule}} = \frac{\text{Support count}(A \cup B)}{\text{Support count}(A)}$$

In general, association rule mining can be viewed as a 2 step process.

Step 1. Find all frequent itemsets

Each of these itemsets will occur atleast as frequently as predetermined minimum support count (Min-Sup).

Step 2. Generate strong association rules from the frequent itemset. These rules must satisfy minimum support & minimum confidence.

K -itemset - A set of items is referred to as an itemset. An itemset that contains K items is an K -itemset.

e.g. 2-itemset {harddisk, printer}

28-03-19 Apriori Algorithm

Apriori algo. was proposed by Aggarwal & Srikant in 1994 for mining frequent item-set. The name of the algo. is based on the fact that the algo. uses prior knowledge of frequent item set properties.

The algorithm uses an iterative approach K/a a level wise search where K -item sets are used to explore $K+1$ item sets.

Apriori starts with item sets of length 1 also K/a 1 item sets as candidates & determines their support by making a scan of the dataset 'D'. Candidates that are infrequent are discarded while the frequent ones are given as output. The frequent itemsets are then used to form candidate itemsets of length 2. This process repeats for longer itemsets until there are no more candidates count. The pseudo-code of this algo. is written below:-

(Program + English)

```

→ Apriori (D, minsup):
    1. C = { all 1-itemsets }
    2. while ( |C| > 0 )
        3. Scan D to find counts of C
        4. F = Sets of c in C with count  $\geq$  minsup
        5. Output (F)
    
```

5. $C = \text{AprioriGen}(F)$

6. $\text{AprioriGen}(F) :$
for each pair of elements X, Y in F
if X and Y , share all items except last

$$Z = X \cup Y$$

if any immediate subset of Z not in F

Prune Z (prune → return Z)

7. Transaction Id | List of Items

$T_1 : I_1, I_2, I_5$

$T_2 : I_2, I_4$

$T_3 : I_2, I_3$

$T_4 : I_1, I_2, I_4$

$T_5 : I_1, I_3$

$T_6 : I_2, I_3$

$T_7 : I_1, I_3$

$T_8 : I_1, I_2, I_3, I_5$

$T_9 : I_1, I_2, I_3$

support count
frequency of
items in D

Let minsup = 2 (IP by user)

Here items = I_1, I_2, I_3, I_4, I_5

Step 1. Make itemsets & put in C Step 2. $|C| > 0$ means there are no of items in C

C	Step 3. C	Support Count
$\{I_1\}$	$\{I_1\}$	6
$\{I_2\}$	$\{I_2\}$	7
$\{I_3\}$	$\{I_3\}$	6
$\{I_4\}$	$\{I_4\}$	2
$\{I_5\}$	$\{I_5\}$	2

$C \rightarrow$ set of candidate / item sets

Step 4.

Here minsup = 2 (given)

\Rightarrow put in F i.e. list of frequent 1 itemsets

F	Support Count
$\{I_1\}$	6
$\{I_2\}$	7
$\{I_3\}$	6
$\{I_4\}$	2
$\{I_5\}$	2

And Pruned

Step 5.

$C = \text{AprioriGen}(F)$

pass F to C as arguments

Step 6.

Make Pairs

(in this step we convert K item sets $\rightarrow K+1$ itemsets $\times C \rightarrow$ candidates)

$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_1, I_4\}$	1
$\{I_1, I_5\}$	2
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2
$\{I_2, I_5\}$	2
$\{I_3, I_4\}$	0
$\{I_3, I_5\}$	1
$\{I_4, I_5\}$	0

as minsup = 2

Here not applicable

* * Apriori Property: This property says that all known empty subsets of a frequent itemset must all be frequent.

F	Count
I_1, I_2	4
I_1, I_3	4
I_1, I_4	1

And Pruned

I_2, I_5 2

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

</

More I_2, I_3 are least
 & craft least often
 are same \Rightarrow combine

$I = X \cup Y$
 $\{I_1, \{I_2\}\} \Rightarrow I_1, I_2, I_3$
 $\{I_1, \{I_3\}\}$
 I_1, I_2, I_3
 $\{I_1, \{I_2\}\}$
 $\{I_1, \{I_3\}\}$
 I_1, I_3, I_5
 $\{I_1, I_3\}$
 $\{I_1, I_5\}$
 I_1, I_3, I_5
 $\{I_2, I_3\} \rightarrow I_1, I_3, I_4$
 $\{I_2, I_4\}$
 $\{I_2, I_3\} \rightarrow I_2, I_3, I_5$
 $\{I_2, I_5\}$
 $\{I_2, I_4\} \Rightarrow \{I_2, I_4, I_5\}$

Before counting

Check

Apriori property

\Rightarrow Check for C

Possible
subsets

I_1, I_2, I_3		
I_1, I_2	I_1, I_3	I_2, I_3

All are frequent i.e. count \geq minsup. Similarly for all

Step 1: $\begin{array}{|c|c|} \hline I_1, I_2, I_3 & \text{Count} \\ \hline I_1, I_2, I_3 & 2 \checkmark \\ \hline I_1, I_2, I_5 & 2 \checkmark \\ \hline \end{array}$

Step 2: $\begin{array}{|c|} \hline F \\ \hline \end{array}$

I_1, I_2, I_3
I_1, I_2, I_5

Step 1: combine

I_1, I_2, I_3, I_5

Step 2:

Check Apriori property

Many subsets that are not valid

\Rightarrow candidate list becomes empty.

\Rightarrow done with loop

29-03-19

Generating Association rules from frequent itemsets

Once the freq. item sets have been found, it is straightforward to generate strong association rules from them. A strong association rule satisfies both minimum support & minimum confidence.

This can be done using following equations:

$$\text{Confidence } (A \Rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

A & B are itemsets

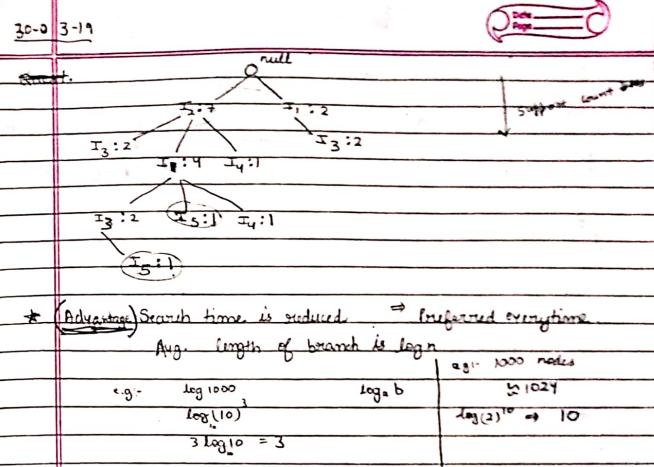
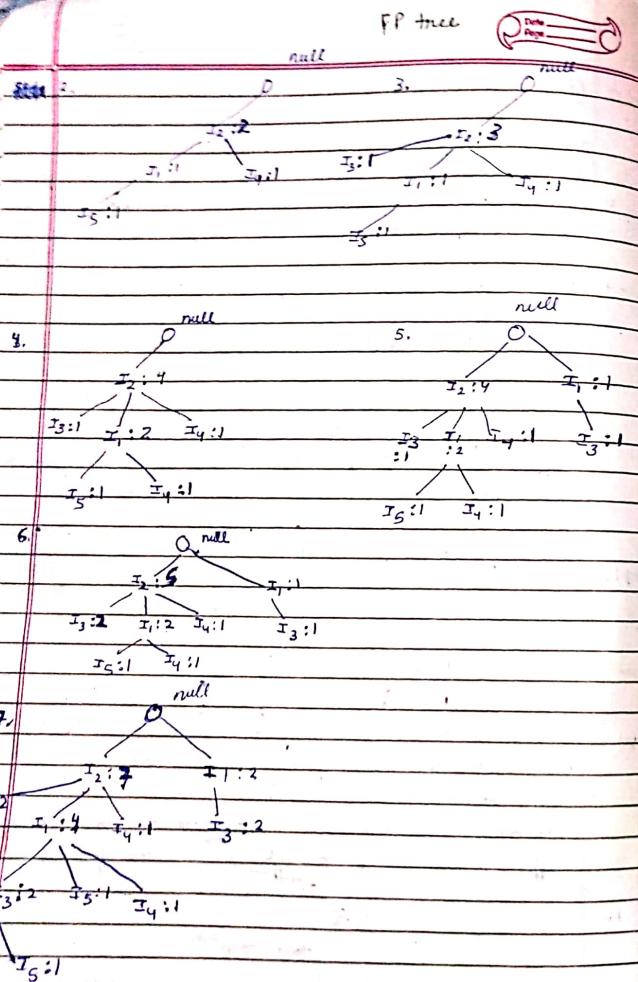
Here Support Count ($A \cup B$) is the no. of transactions containing the itemsets $A \cup B$ & Support Count (A)

is the no. of transactions containing the itemset A . Based on the regulation, the association rules can be generated as follows:

For each freq. itemset I generate all nonempty subsets of I

For all nonempty subset S of I , output the rule $I - S \Rightarrow I - S$

associated with



Generating frequent patterns from FP tree

An FP tree is mined as follows:

Start from each frequent length 1 pattern as an initial suffix pattern, construct its conditional pattern base which consists of the set of prefix paths in the FP tree co-occurring with the suffix pattern. Then construct its conditional FP tree & perform mining recursively on the tree. The pattern growth is achieved by the concatenation of a suffix pattern with the frequent patterns generated from conditional FP tree.

We first consider I_5 , which is the last item in the initial list rather than the first item. The item I_5 occurs in 2 FP tree branches. The occurrences can be found by following chain of node links.

Write item having least support count

Date Page

Item	Conditional Pattern	conditional FP tree	Frequent pattern generated
	Base		$\langle I_2 : 2, I_1 : 2 \rangle$ $\langle I_2 : 2, I_3 : 2 \rangle$ $\langle I_1 : 2 \rangle$ $\langle I_2 : 2, I_5 : 2 \rangle$ $\langle I_1 : 2, I_2 : 2, I_5 : 2 \rangle$
I_5	$\{I_2, I_1, I_3 : 1\}$ $\{I_2, I_1, I_5 : 1\}$		
I_4	$\{I_2, I_1 : 1\}$ $\{I_2 : 1\}$	$\langle I_2 : 2 \rangle$ $\langle I_2, I_4 : 2 \rangle$	I_5 is added suffix

Considering I_5 as a suffix, its corresponding tree :

prefix paths are :-
 $I_2, I_1 : 1$
 $I_2, I_1, I_3 : 1$ occurring 1 time in tree

for conditional FP tree I_2 has support count 2 in conditional pattern base

I_1 also

reject I_3 as Min Sup. Count should be 2 (considered in previous example)

I_3	$\{I_2 : 2\}$ $\{I_1 : 2\}$ $\{I_2, I_1 : 2\}$	$\langle I_2 : 4, I_1 : 1 \rangle$ $\langle I_2, I_3 : 2 \rangle$ $\langle I_1 : 2 \rangle$ $\langle I_2, I_2, I_3 : 1 \rangle$	$\langle I_2, I_3 : 2 \rangle$ $\langle I_1 : 2 \rangle$ $\langle I_2, I_1 : 1 \rangle$
I_1	$\{I_2 : 4\}$	$\langle I_2 : 4 \rangle$	$\langle I_2, I_1 : 1 \rangle$

X I_2
Don't consider
 I_2 (last item)
Table 2 from book up approach

The FP growth method transforms the problem of finding long freq. patterns into searching for shorter ones in much smaller conditional databases recursively.

A study of the FP growth method performance shows that it is efficient & scalable for mining both long & short frequent patterns and it is about an order of magnitude faster than the apriori algorithm.

01-04-19

Chapter - Classification

nominal

- Classification is a form of data analysis that extracts models describing important data classes. Such models called classifiers predict categorical class input labels (discrete & unordered). e.g. - We can build a classification model to categorise bank loan applications as either safe or risky.

Many classification methods have been proposed by this method. In Machine learning, pattern recognition & statistics.

Classification has numerous applications including fraud detection, performance prediction, manufacturing & medical diagnosis. Classification is a 2 step process consisting of:-

① Learning Step

a classification model is constructed

② Classification Step

model is used to predict class label of given data

① → In the first step, a classifier is built describing a pre determined set of data classes or concepts.

② → This is the learning step or training phase where a classification algorithm builds the classifier by analysing or learning from a training set made up of database tuples.

& their associated class labels. A tuple \vec{x} is represented by an 'n' dimensional attribute vector X

$$X = (x_1, x_2, \dots, x_n)$$

Each tuple is assumed to belong to a predefined class or determined by another database attribute called

The Class Label Attribute.
It is discrete & considered.

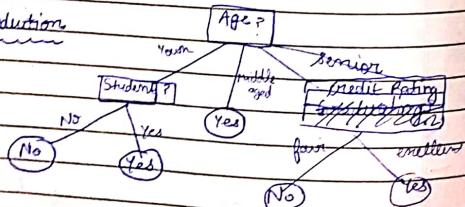
Name	Age	Income	From discussion	
ABC	25	50,000	Safe	
DEF	30	20,000	Risky	
GHI	32	25,000	Risky	
JKL	22	15,000	Risky	
MNO	28	35,000	?	

Three tuples whose class labels are unknown, fall in the category of Testing Set.

In learning step \leftarrow model (formula) is created
Algorithm of classification is generated i.e. used in next step

(The attribute to be predicted is Class Label Attribute)

* Decision Tree Induction



continuous \rightarrow discrete
Age from 60
0-10 Child
10-20 Teenager
20-60 Adult

- Decision tree induction is learning of decision trees from class label training tuples.
- Decision tree is a flow-like decision structure where each internal node denotes a test on an attribute. Each branch represents an outcome of the test & each leaf node holds a class label.

The topmost node in the tree is root node. A typical decision tree is shown above. It represents the concept of computer. It predicts whether a customer is likely to purchase a computer or not. Internal nodes are denoted by rectangles & leaf nodes are denoted by ovals. Given the tuple for which class label is unknown, the attribute values of the tuple are tested against the decision trees. A path is traced from the root to the leaf node which holds the class prediction for that tuple.

The construction of decision trees doesn't require any domain knowledge.

- The learning & classification of decision tree induction are simple & fast.
- Decision tree classifiers have good accuracy.
- Decision tree induction algos have been used for classification in many application areas such as medicine, financial analysis, astronomy, manufacturing & production.

03-09-19

K-Nearst Neighbour Algorithm

Classification \rightarrow predicting value of an attribute.

A K-nearest neighbour method is a simple technique which considers a training data set itself as the classification model. In this approach, the records of the training data set are considered as points in a d-dimensional space (d-attributes). Several possible distance functions can be used to find the similarity b/w the 2 objects.

Given a new point, 'p' to classify, the k-nearest points of p in the training data set are found & class i.e. assigned to most of these k-points is selected as the class for 'p'.

Although simple, this approach could be very slow during classification because it needs to search for the k-nearest points among the entire training data set which could be huge.

(S.ID)	Height	Weight	Class	to predict	
				62	53
1	140 cm	50 Kg	7		
2	147 "	60 "	7		
3	148 "	48 "	6		
4	150 "	60 "	8		
5	152 "	55 "	7		
6	155 "	65 "	8		
7	153 "	62 "	?		

This is numerical attribute data \Rightarrow we use Euclidean dist.

$$\text{For object 2,3} = \sqrt{(142-140)^2 + (60-48)^2}$$

Euclidean \rightarrow
dist.

It should be odd to predict majority.

Find distance b/w object 7 with all others.

$$1. (1,7) = \sqrt{(140-153)^2 + (50-62)^2} = \sqrt{169+144} = \sqrt{313} = 17.6$$

$$2. (2,7) = \sqrt{(147-153)^2 + (60-62)^2} = 6.324$$

$$3. (3,7) = \sqrt{(148-153)^2 + (48-62)^2} = 14.86$$

$$4. (4,7) = \sqrt{(150-153)^2 + (60-62)^2} = \sqrt{3^2+2^2} = \sqrt{13} = 3.605$$

$$5. (5,7) = \sqrt{(152-153)^2 + (55-62)^2} = \sqrt{1^2+7^2} = \sqrt{48} = 6.928$$

$$6. (6,7) = \sqrt{(155-153)^2 + (65-62)^2} = \sqrt{(2)^2+(3)^2} = \sqrt{4+9} = \sqrt{13} = 3.605$$

(4, 5, 6) in ascending order having least dist. value.
 ↓ ↓ ↓
 Class 8 7 8 \Rightarrow predict value = 8

Perform with k other values if not majority
 values: 3.6, 10.04, 3.6

03

partitioning based algorithm

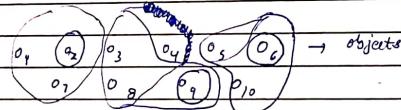
06-04-19

K-Means Algorithm

Input : K , no. of clusters, D dataset
Output : A set of K clusters

- 1) Arbitrarily choose K objects from D as initial clusters
- 2) Repeat
- 3) Assign each object to the cluster to which the object is the most similar based on the minkowski value of the objects in the cluster.
- 4) Update cluster means; that is calculate the mean value of the objects for each cluster until no change.

e.g.

Let $K = 3$

- check similarity of cluster-heads with other objects
- update cluster head $\rightarrow o_1 + o_2 + o_7 = \text{hd}_1$

a.g.	Age	Salary
	20	20000
Similarly H_2, H_3	30	30000
Mean	25	25000

- check again similarity

until stable clusters

$$E = \sum_{i=1}^k \sum_{p \in C_i} \text{dist}(p, c_i)$$

capital small

Here E is the sum of the squared error for all objects in the data set.

$p \rightarrow$ point in space representing a given object.
 $c_i \rightarrow$ It is the centroid of the cluster

This objective fn. tries to make the resulting K -clusters as compact & as separate as possible.

(This value should be min. for cluster stability)

12-04-19

Theoretical Methods

The clusters which are very similar to each other are grouped into larger clusters. These larger clusters may further be grouped into still larger clusters. In this way, a hierarchy or tree of clusters is produced

Hierarchical approach types

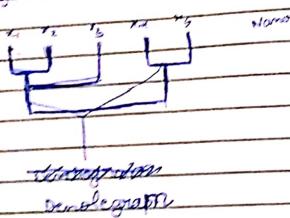
Agglomerative approach

1. $N = \sigma$
2. for each point pair $\neq D$
3. Initialize cluster $C_i = \{p\}$
4. $N = N \cup C_i$
5. repeat :
6. $(c_i, c_j) =$ closest pair of unmerged cluster
7. $C_x = \text{Merge } (c_i, c_j)$
8. $N = N \cup C_x$
9. until No unmerged cluster exist in N

63

Agglomerative approach (Cluster algorithms) -

starts with all the data objects in separate clusters & then merge nearby clusters until a single cluster is formed that contains all the objects.
The algorithm initializes each data object into its own cluster, then in each iteration, it merges the closest pair of clusters.



balanced structure Reducing & clustering using hierarchies

This is the first hierarchical method i.e. designed to be scalable for very large datasets. It treats in 2 phases:-

During the first phase it forms a hierarchical clustering. This clustering is K/a clustering feature tree (CFT).

2. In the second phase, the leaf level clustering is obtained at the end of the first one reclustered one tree clustered using any other standard clustering strategy.

For each cluster in the CFT with during the first phase the algorithm maintains a vector of summary information K/a clustering Feature.

If x_1, x_2, \dots, x_n are points in the cluster which themselves are vectors of attribute values. Its clustering feature is given as

$$\bar{x}_n, \sum x_i, \sum x_i^2$$

The clustering feature gives very useful information about clustering. It is possible to calculate size, mean & variance within a cluster.

The formula to calculate diameter of cluster is:

$$\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^{i-1} (x_i - x_j)^2}$$

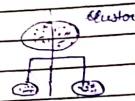
While forming the clustering feature tree, if the diameter of any leaf-level cluster exceeds a user

given threshold ' t ', the cluster is split into 2. This step is done in manner similar to splitting of a node in B+ trees.

in data structure

DS, reduce search time

In B+ tree would involve creating or updating the interior nodes that are ancestors of the leaf been split. The pseudo code of the algorithm is written as:-



1. $N = \{ \cdot \} // N$ is initial CF tree.
2. For each point p in \mathbb{R}^d
3. $m = \text{leaf node in } m, i.e., \text{closest to } p$.
4. Add p to m
5. Compute diameter M of m .
6. If $M > t$ → threshold (Exercise)
7. Split m
8. Apply another clustering algorithm to cluster the leaves of $'N'$.