

Data Preprocessing

Major steps in Data Preprocessing:-

The steps involved in the data preprocessing are

1. Data cleaning
2. Data Integration
3. Data Reduction
4. Data Transformation.

Data Cleaning:-

Data cleaning routines to clean the data by filling the missing values, smoothing noise data, removing inconsistencies, identifying or removing outliers.

Data Integration:-

Collecting the data from multiple sources like databases, datacubes, flatfiles can be called as

Data Integration.

Ex:- DataCubes

Attribute for customer identification may be referred as customer id in the database. ev-id, id other databases. which can create a confusion (or) redundancy in analyzing the data.

Having large amount of redundant data may decrease the knowledge discover process.

performance of

Hence data reduction is required

Data Reduction:-

Reducing a dataset into a smaller set can be called as Data Reduction.

Data Reduction strategies includes

i) Dimensionality Reduction

a) Numerosity Reduction.

In dimensionality Reduction data compression techniques can be used and those are.

i) wavelength Transforms.

a) Principle Component Analysis.

3) Attribute subset selection.

4) Attribute construction.

Numerosity Reduction:-

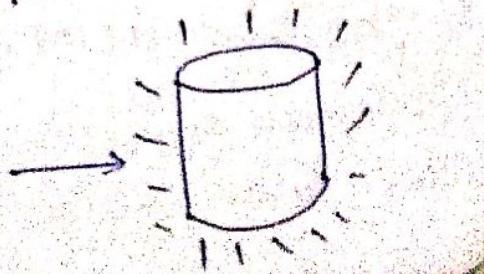
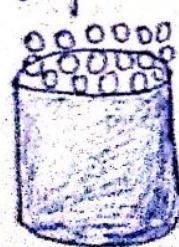
In this the data is replaced by smaller representation using either parametric models (Regression or loglinear models) or Non-parametric models (Histogram, Sampling, Clustering, data Integration).

Data Transformation:-

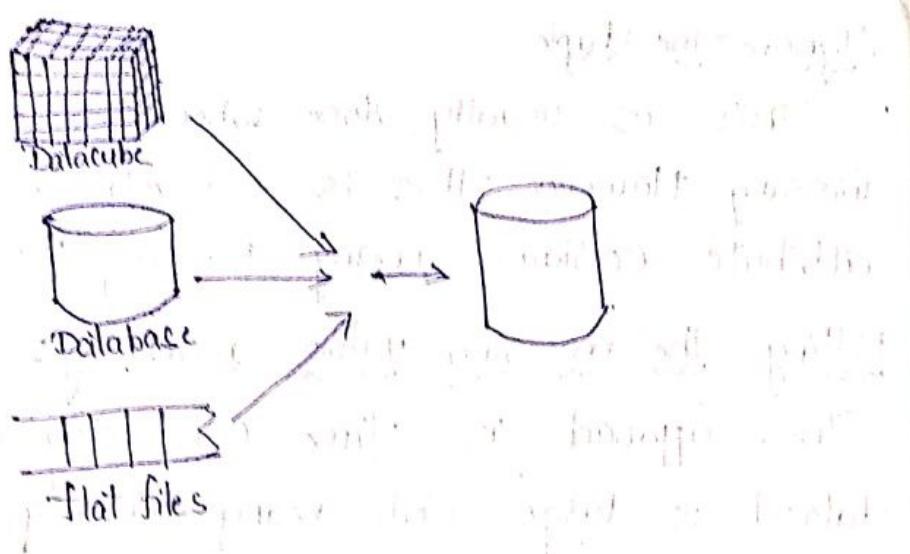
Convert the data into appropriate form for mining can be called as Data Transformation.

Forms of Data Preprocessing:-

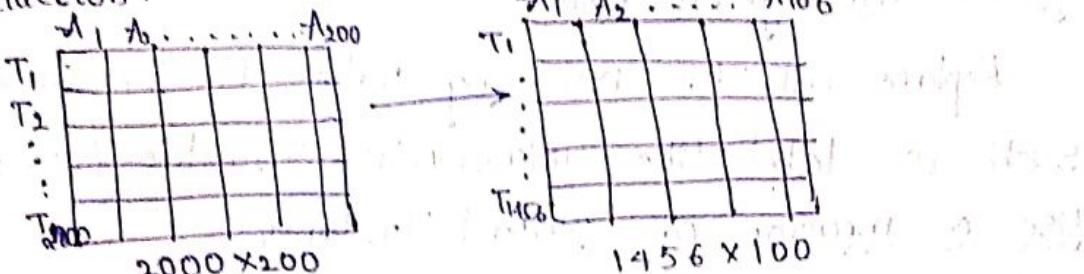
Data cleaning:-



Data Integration:-



Data Reduction:-



Data Transformation:-

$$-2, +8, 58, 100, \dots \rightarrow -0.02, 0.48, 0.58, 1.0$$

Data Cleaning:-

Data cleaning routines to clean the data while filling the missing values, smoothing noisy data, Removing (0) identifying outliers.

Missing values:- In the dataset can be filled by using different ways so, the ways are:-

1. Ignore the tuple.
2. Filling in the missing values manually
3. place the global value.
4. Use the measure of central tendency to fill the missing values.
5. Use the attribute mean (or) median for all the samples belonging to the same class.
6. Use the most probable value to fill in the missing value.

Ignore the tuple:-

This is usually done when the class label is missing. However this is not effective unless the attribute contains many missing values.

Filling the missing values manually:-

This approach is time consuming whether the given dataset is large with many missing values.

Place the global value:-

Replace all the missing values by the same constant such as label like unknown (or) infinite. (∞ , ?)

Use a measure of Central Tendency:-

In this missing values can be filled by using the mean (or) median value.

If the data is Symmetric then use the mean value to fill in the missing value.

Skewed - median.

6)

This may be determined with regression or decision tree induction.

Ex:- Use regression to fill the missing value.

x	y	$\frac{x}{x-\bar{x}}$	$\frac{y}{y-\bar{y}}$	xy	x^2	y^2	$\Sigma x = 49$	$\bar{x} = 13$	$\Sigma y = 67$	$\bar{y} = 16$
10	15	-3	-2	8	9	4				
12	16	-1	-1	8	9	4				
13	19	0	0	1	1	1				
14	17	1	8	0	0	4				
15	?	2	0	0	1	0				
				4						

$$y - \bar{y} = \frac{\Sigma xy}{\Sigma x^2} (x - \bar{x})$$

$$y - 17 = \frac{7}{11}(x - 13)$$

$$11(y - 17) = 7(x - 13)$$

$$11y - 187 = 7x$$

$$11y = 187 + 14$$

$$y = \frac{201}{11} = 18.27 \approx 19.$$

29/12/18.

Noisy Data:-

Noise is a random error (or) variance in a measured variable.

Noisy data can be smoothed by using different processors by binning, Regression, Outlier Analysis.

Binning:-

Binning methods smooth a soiled data values.

Different Binning methods are:-

1. Bin by means

2. Bin by median.

3. Bin by boundaries.

Apply the binning methods for the following data.

4, 8, 15, 21, 21, 24, 25, 28, 34

(a) Partitions into equal frequencies

Bin1: 4, 8, 15

Bin2: 21, 21, 24

Bin3: 25, 28, 34

3. Bin by boundaries

Bin1: 4, 4, 15

Bin2: 21, 21, 24

Bin3: 24, 25, 34

1. Bin by means

Bin1: 9, 9, 9

Bin2: 22, 22, 22

Bin3: 29, 29, 29

2. Bin by median:-

Bin1: 8, 8, 8

Bin2: 21, 21, 21

Bin3: 28, 28, 28

Bin by mean:-

Replace each bin value with its corresponding value.

Bin by median:-

Replace each bin value with its median value.

Bin by boundaries:-

1. Write the boundaries of a bin.

2. Fill the middle values of the bin with the desired boundary.

* Apply the binning procedure for the following elements.

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 50, 70. Use the bin depth in 3.

Partition the data into equal frequencies.

Bin 1: 13, 15, 16, 16, 18, 20, 20, 21, 22

Bin 2: 22, 25, 25, 25, 25, 25, 30, 33, 33, 35

Bin 3: 35, 35, 35, 36, 40, 45, 46, 50, 70.

Bin by means:-

Bin 1: 18, 18, 18, 18, 18, 18, 18, 18, 18

Bin 2: 28, 28, 28, 28, 28, 28, 28, 28, 28

Bin 3: 44, 44, 44, 44, 44, 44, 44, 44, 44

Bin by median:-

Bin 1: 19, 19, 19, 19, 19, 19, 19, 19, 19

Bin 2: 25, 25, 25, 25, 25, 25, 25, 25, 25

Bin 3: 40, 40, 40, 40, 40, 40, 40, 40, 40

Bin by boundaries:-

Bin 1: 13, 13, 13, 13, 22, 22, 22, 22, 22

Bin 2: 22, 22, 22, 22, 22, 35, 35, 35, 35

Bin 3: 35, 35, 35, 35, 35, 35, 35, 35, 70.

Regression:-

2/01/2019

The regression can also be used for smoothing the data.

Regression can be performed in two ways

1. Linear Regression.

2. Multiple Regression.

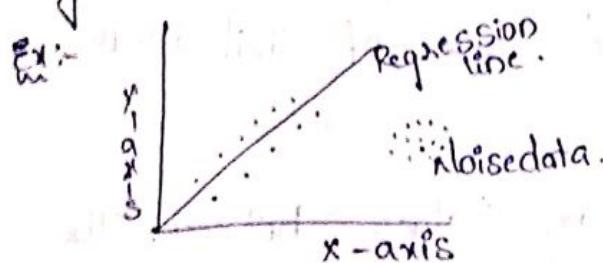
Linear Regression:-

Linear Regression involves finding the best line to fit two attributes so that one attribute can be used to predict the other.

Multiple Regression:-

Multiple Regression can be used when there are more than two attributes.

(Regression is completed.)

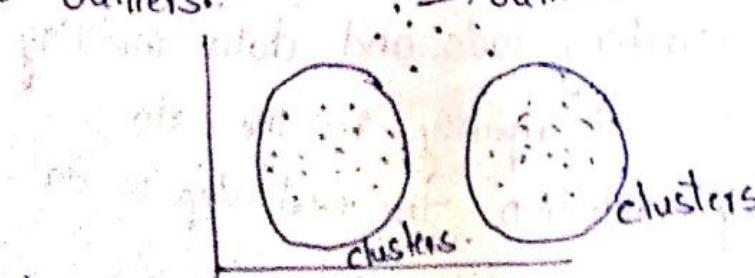


Outlier:-

Outlier can be detected by using clustering.

Objects that are similar to each other can form a group.

Objects that fall outside the groups can be called as outliers.



Data Cleaning as a Process:-

The first step in the data cleaning as a process is discrepancy detection.

Discrepancy can be caused by several factors.

- 1) Poorly designed forms with many fields as optional.
- 2) Human error in data entry.
- 3) data decay (Outdated Address).

The data should also be examined by regarding the following rules

1. Unique Rule

2. Consecutive Rule.

3. Null Rule.

Unique Rule:-

Each value of the given attribute must be different from all the other values of that attribute.

Consecutive Rule:-

There can be no missing values between the lowest and highest values of that attribute.

Null Rule:-

Null Rule specifies the use of blanks, question marks, special characters or other strings that may indicate null condition and how such values should be handled.

There are data scrubbing tools and data auditing tools for finding the discrepancy in the data.

After the discrepancy detection the next step is data transformation.

1. Data migration tools.

2. ETL tools (Extraction Transformation and Loading)

31/01/19
Data Integration:-

Merging or Combining - the data from multiple sources.

There are some important issues in the data integration.

1) Redundancy & co-relation Analysis.

2) Entity identification problem.

3) Tuple duplication.

4) Data value conflict detection & Resolution.

Redundancy and co-relation Analysis:-

Redundancy can occur if one attribute is derived from another attribute or a set of attributes.

Inconsistency can also occurs due to the redundancy in the data.

(Redundacy may be detected by some of the redundancy may be detected by co-relation analysis.

Co-relation Analysis:-

Given two attributes, the analysis of how strongly one attribute implies the other attribute based on the attribute data.

For the nominal data - Chi-square test is applied.

Numerical data - Correlation Coefficient and Co-variation can be used.

Chi-Square Test:-

Given two attributes - the gender and preferred veg readings are given in the table.

Gender:

	Male	Female	Total
Fiction	250 (90)	200 (360)	450
Non-fiction	50 (210)	1000 (840)	1050
Total	300	1200	1500

Test that is the reading is same for both the genders.

Sol- Null hypothesis: There is a preference between the gender
tabulated value for 1 of χ^2 is $= 10.828$.

$$\text{cal } \chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Expected Frequencies:-

$$E_{11} = \frac{\text{count(male)} \times \text{count(fiction)}}{\text{total}}$$

$$E_{11} = \frac{300 \times 450}{1500} = 90$$

$$E_{12} = \frac{C(f) \times \text{count(fiction)}}{\text{total}} = \frac{1200 \times 450}{1500} = 360$$

$$E_{21} = \frac{C(m) \times c(m)}{\text{total}} = \frac{1050 \times 300}{1500} = 210$$

$$E_{22} = \frac{c(f) \times c(m)}{\text{total}} = \frac{1200 \times 1050}{1500} = 840$$

$$\begin{aligned} \text{cal } \chi^2 &= \frac{(250-90)^2}{90} + \frac{(200-360)^2}{360} + \frac{(50-210)^2}{210} + \frac{(1000-840)^2}{840} \\ &= \frac{25600}{90} + \frac{25600}{360} + \frac{25600}{210} + \frac{25600}{840} \end{aligned}$$

$$\approx 507.9.$$

∴ degree of freedom is $(c-1)(r-1) = (2-1)(2-1) = 1$

If the calculated chisquare is greater than the tabulated chisquare then reject null hypothesis.
Otherwise,

If calculated chisquare is less than the tabulated chisquare then accept the null hypothesis

$507.9 > 10.828$ is true, hence reject the null hypothesis

Hence there is no preferences between the genders.

In an investigation, all the machine performances the following results are obtained.

	No. of units inspected	No. of defects
M ₁	375 (378.6)	17 (17.6) 392
M ₂	450 (455.9)	22 (21.5) 472
	825	39
		864

Test whether there is any significant performance of two machines.

Sol:- There is a significant performance of two machines.
The tabulated value for 1 of χ^2 is = 10.828.

Expected frequencies -

$$e_{11} = \frac{825 \times 392}{864} = 378.6$$

$$e_{12} = \frac{39 \times 392}{864} = 17.6$$

$$e_{21} = \frac{172 \times 825}{864} = 455.9$$

$$e_{22} = \frac{39 \times 472}{864} = 21.5$$

$$\text{cal } \chi^2 = \frac{(375 - 378.6)^2}{378.6} + \frac{(450 - 455.9)^2}{455.9} + \frac{(17 - 17.6)^2}{17.6} + \frac{(22 - 21.5)^2}{21.5}$$

$$\text{cal } \chi^2 = 0.015$$

$\text{cal } \chi^2 < \text{tab } \chi^2$
 $0.015 < 10.828$ hence accept the null hypothesis.

* On the basis of information given below about the treatment of 200 patients suffering from a disease, state whether the new treatment is comparatively superior to the conventional treatment.

	Favourable	Non-favourable	Total
New	60 (45)	30 (45)	= 90
Conventional	40 (55)	70 (55)	= 110
	100	100	200

New Treatment is the Superior to the Conventional Treatment (or) New and Conventional treatment are independent.

The χ^2 value for the degree of freedom of 1 is

10.828.

Expected values:-

$$\text{Cal } \chi^2 = \sum_{i=1}^C \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{11} = \frac{C(n) \times C(f)}{n} = \frac{90 \times 100}{200} = 45$$

$$E_{12} = \frac{C(n) \times C(nf)}{n} = \frac{90 \times 100}{200} = 45$$

$$E_{21} = \frac{C(c) \times C(f)}{n} = \frac{110 \times 100}{200} = 55$$

$$E_{22} = \frac{C(c) \times C(nf)}{n} = \frac{110 \times 100}{200} = 55$$

$$\text{Cal } \chi^2 = \frac{(60 - 45)^2}{45} + \frac{(30 - 45)^2}{45} + \frac{(40 - 55)^2}{55} + \frac{(70 - 55)^2}{55}$$

$$\text{Cal } \chi^2 = 18.18$$

$\text{cal}(x^2) > \text{tab}(x^2)$ so reject the hypothesis.

∴ new and conventional treatment are dependent.

Calculating the Correlation for the numerical data:-

For the Numerical data Correlation analysis can be done by co-relation coefficient.

Consider two numeric attributes A & B and a set

of n observations $\{(a_1, b_1), \dots, (a_n, b_n)\}$

The mean values of A and B respectively are also known as expected values of A & B.

$$\bar{A} = \frac{\sum_{i=1}^n a_i}{n}, \quad \bar{B} = \frac{\sum_{i=1}^n b_i}{n}$$

The co-variance between A and B defined as

$$\text{cov}(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

The co-relation coefficient of A, B can be defined as

$$r_{A,B} = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B}$$

Note:- If the value of r is limited to the $-1 \leq r \leq 1$

If $r_{A,B} = 0$ then there is no co-relation between A, B.

If $r_{A,B} > 0$ then there is +ve correlation between A, B.

If $r_{A,B} < 0$ then there is -ve correlation between A, B.

* Let the two attributes temperature and ice cream values

are given below.

$(b_i - \bar{B})$	Temperature (A)	Icecream (B)	$a_i - \bar{A}$	$b_i - \bar{B}$	$(a_i - \bar{A})(b_i - \bar{B})(a_i - \bar{A})$
84.64	6	20	1	9.2	18.4
0.64	5	10	-0.8	-0.8	1
10.24	4	14	0	+3.2	0
33.64	3	5	-1	-5.8	5.8
33.64	2	5	-2	-5.8	11.6
162.8	$\bar{A} = 20$	$\bar{B} = 54$	0		$\frac{4}{10}$

$$\bar{A} = \frac{20}{5} = 4, \bar{B} = \frac{54}{5} = 10.8$$

Let us assume that Temperature is A and ice cream is B

$$\text{COV}(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{a})(b_i - \bar{b})}{n} = 35/5 = 7$$

$$\sigma_A^2 = \frac{\sum_{i=1}^n (a_i - \bar{a})^2}{n} = \frac{10}{5} = 2.$$

$$\sigma_B^2 = \frac{168.8}{5} = 32.56.$$

$$\sigma_A = \sqrt{2} = 1.42.$$

$$\sigma_B = 5.70.$$

$$r_{A,B} = \frac{\text{COV}(A, B)}{\sigma_A \sigma_B} = \frac{7}{(1.42)(5.70)} = 0.86.$$

Hence the two attributes temperature and icecream are in positive.

- * In one database height in inches and weight in kgs are given. Find if there is any significant correlation between the heights and weights for the given data.

Height in inches	Weight in kgs
[A]	[B]
57	13
59	117
62	126
63	126
64	130
65	129
55	111
58	116
57	112
$\Sigma A = 540$	
$\Sigma B = 1080$	

Heights in ft (A)	Weights in kg (B)	$a_i - \bar{A}$	$b_i - \bar{B}$	$(a_i - \bar{A})(b_i - \bar{B})$	$(a_i - \bar{A})^2$	$(b_i - \bar{B})^2$
57	113	-3	-7	21	9	49
59	117	-1	-3	3	1	9
62	126	2	6	12	4	36
63	126	3	6	18	9	36
64	130	4	10	40	16	100
	129	5	9	45	25	81
65	111	-5	-9	45	25	81
55	116	-2	-4	8	4	16
58	112	-3	-8	24	9	64
				216	102	472
$\Sigma A = 540$						
$\Sigma B = 1080$						

$$\bar{A} = \frac{540}{9} = 60, \quad \bar{B} = \frac{1080}{9} = 120$$

$$\rho = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B}, \quad \text{cov}(A, B) = \sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})$$

$$\sigma_A = \sqrt{\frac{\sum_{i=1}^n (a_i - \bar{A})^2}{n}}$$

$$\sigma_A = \sqrt{\frac{102}{9}} = \sqrt{11.3} = 3.36.$$

$$\sigma_B = \sqrt{\frac{\sum_{i=1}^n (b_i - \bar{B})^2}{n}} = \sqrt{\frac{472}{9}} = 7.24.$$

$$\rho = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B} = \frac{24}{3.36 \times 7.24} = 0.986.$$

$\rho = 0.986$, it means the two attributes are in positive correlation.

∴ The two attributes are strongly related. Hence one attribute can be used in the database instead of 2 attributes.

Entity Identification Problem:-
How can equivalent realworld entities from the multiple sources can be matched up.
This is called as entity identification problem.
Ex:- How the data analyst (or) a Computer be sure that the customer id in one database and customer number in other refer to same attribute.
When matching attributes from one database to another database during integration, a special attention must be paid to the structure of the data.

Tuple Duplication:-

In addition to detecting redundancies between the attributes duplication should also be detected at Tuple level.

The use of these normalized tables is another source of the redundancy.

Data Value Conflict detection and Resolution:-

Data Integration also involves detection and resolution of data values.

For the same realworld entities attributes values from different source may differ. This may be due to difference in scaling, representation, encoding.

Example:-

A weight attribute may be stored in the metric units in one system and british imperial in another.

Variance is a special case of co-variance where the two attributes are identical.

Variance:-

The Variance of n observations $(a_1, a_2, a_3, \dots, a_n)$ formulae for Variance in a numeric attribute A is

$$\sigma_A^2 = \frac{1}{n} \sum_{i=1}^n (a_i - \bar{A})^2$$

Ex:- Find the Variance for the values 30, 36, 44, 50, 52, 60, 65, 70, 70, 110.

11/01/19.

Data Reduction:-

Data Reduction Techniques can be applied to obtain a reduced representation of the data set i.e much smaller in volume, yet closely maintains the integrity of the original data.

Data Reduction Techniques involves Numerosity Reduction

Dimensionality Reduction, data Comprehension

* Under Dimensionality Reduction we have 13 Techniques

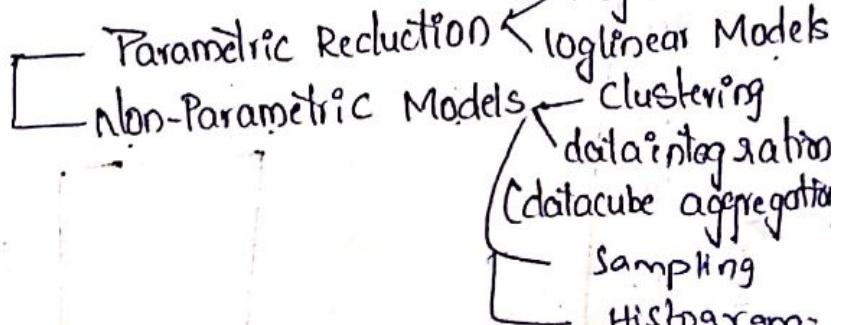
They are:-

1. Wavelet transforms

2. PCA

3. Attribute subset Selection

* Numerosity Reduction



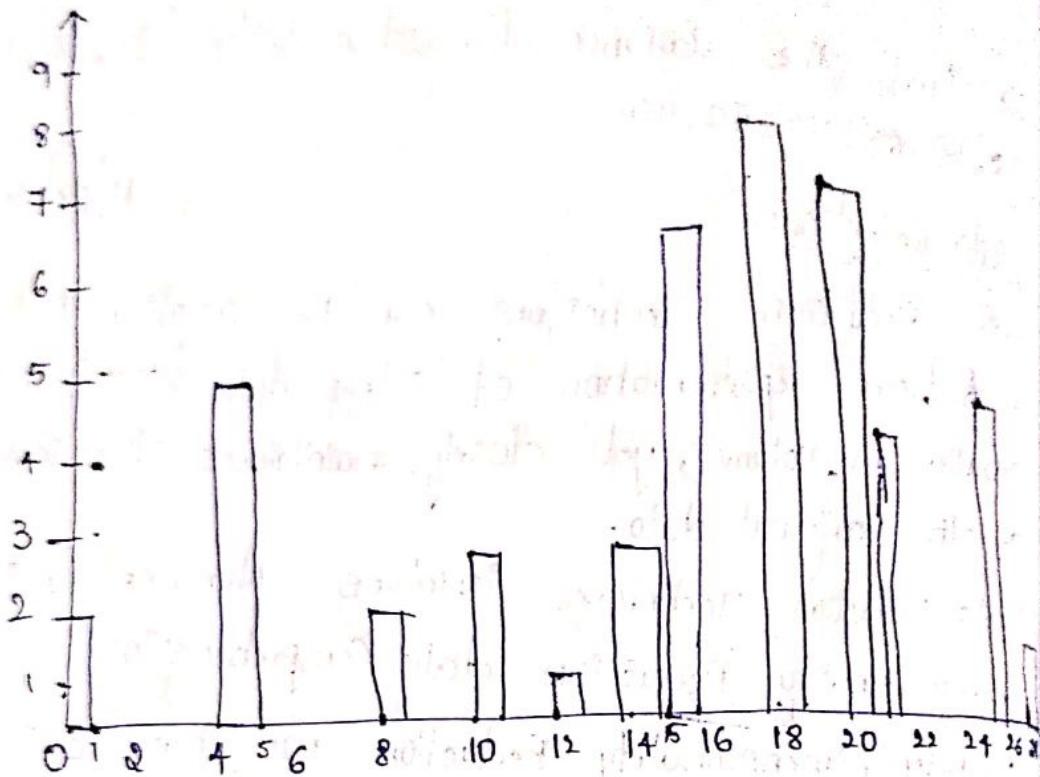
Histograms:-

Histograms uses binning to approximate data distribution.

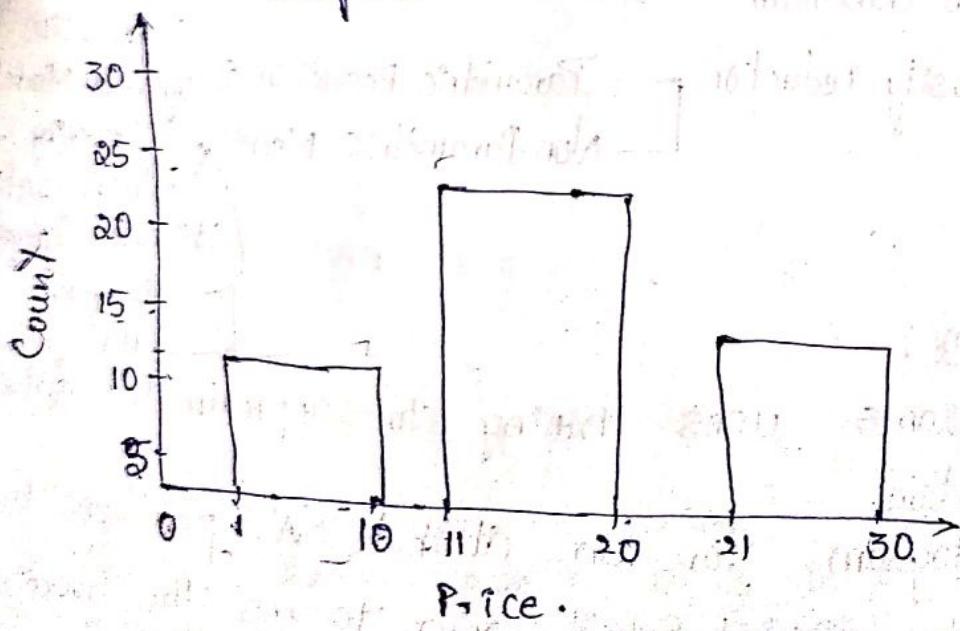
A histogram for an attribute A partitions the data into disjoint subset refer to as the buckets or wings bins.

* Construct the histogram by using Singleton buckets for the following data which contains

prices for commonly sold items.



Equal width Histogram

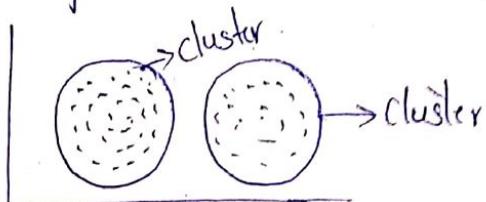


Clustering:-

29/01/19

Forming the given data into groups which are similar,
clustering is the method of identifying similar group of
data in dataset.

Objects within the cluster are similar to each other and
dissimilar to objects of the other clusters.



* Diameter can be used to measure the quality of a cluster.

* Centroid distance can also be used to measure the quality of a cluster.

Sampling Techniques that can be used for data reduction
are as follows:-

1. Simple Random Sampling without replacement (SRSWOR)

2. Simple Random Sampling with replacement (SRSWR)

3. Cluster Sampling

4. Stratified Sampling.

SRSWOR:-

Suppose that a large dataset 'D' contains 'n' tuples.

SRSWOR can be created by drawing 's' of n tuples

from 'D'.

The probability of drawing any tuple in 'D' is $1/n$ i.e.
all the tuples are equally likely to be sampled.

SRSWR:-

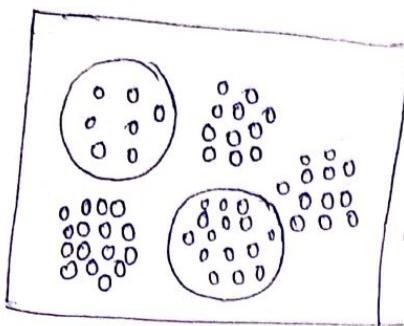
This is similar to SRSWOR except that each type a
tuple is drawn from 'D', it is recorded and then
replaced.

Cluster Sampling:-

Divide the data into groups called clusters.

Randomly select some of the groups for sampling.

Ex:-



Stratified Sampling:-

Divide the data into groups called strata.

Select randomly from individual element of each strata.

Ex:- There are 50 states (dividing the whole population of the country into states)

Randomly select 3 people from each state.

$$3 \times 50 = 150 \text{ members.}$$

i.e. Stratified Sampling count is 150 members.

Example for Cluster Sampling:-

Selecting atleast two states (All the members of cluster)

Example for SRSWOR:-

Assume that there are 6 red, 5 blue and 4 green balls in a dataset.

What is the probability of getting red ball, green ball and blue ball in one graph.

$$P(R \cap B \cap G) = \frac{8}{15} \times \frac{5}{14} \times \frac{4}{13} = \frac{4}{91} \quad \text{With Replacement:-}$$

$$P(R) = \frac{6}{15}$$

$$P(B) = \frac{5}{14}$$

$$P(G) = \frac{4}{15}$$

$$P(R) = \frac{6}{15}$$

$$P(B) = \frac{5}{15}$$

$$P(G) = \frac{4}{15}$$

$$P(R \cap B \cap G) = \frac{8}{15} \times \frac{5}{14} \times \frac{4}{13} = \frac{8}{225}$$

Performing aggregation operation on attribute values to analyse the data.

Ex:-

Year	2008
Q ₁	\$ 15000
Q ₂	\$ 16000
Q ₃	\$ 18000
Q ₄	\$ 24000

To know the annual sales of year 2008 then perform aggregation operation of all quarters.

Parametric Models:-

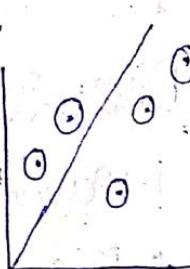
1. Regression:-

Regression and log linear models can be used to approximate the given data.

In the linear Regression, the data are modelled to fit a straight line (Regression line) (Response Variable). For example, a random variable y , can be modelled as a linear function of another random variable x (Predictor Variable) (x Independent Variable) $y = mx + c$.

Find the regression line for the following data.

x	y	(x - \bar{x})	(y - \bar{y})	(x - \bar{x})(y - \bar{y})	(x - \bar{x}) ²
1	3	-2	-0.6	1.2	4
2	4	-1	+0.4	-0.4	1
3	2	0	-1.6	0	0
4	4	1	0.4	0.4	1
5	5	2	1.4	2.8	4
$\bar{x} = 3$		$\bar{y} = 3.6$		$\frac{10}{4}$	$\frac{10}{4}$



$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{10}{10} = 1$$

$$y = mx + c$$

$$3.6 = (0.4)(3) + c$$

$$3.6 = 1.2 + c$$

$$c = 3.6 - 1.2$$

$$c = 2.4$$

$y = 0.4x + 2.4$ is the required equation.

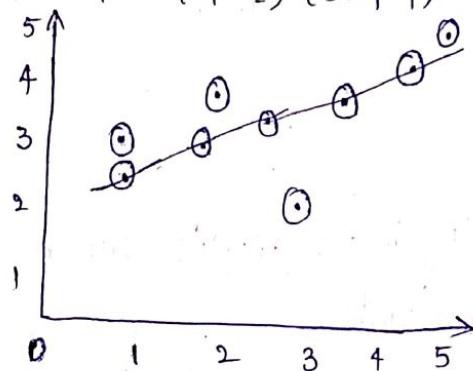
$$x=1 \Rightarrow y = 0.4 + 2.4 = 2.8 \Rightarrow (1, 2.8)$$

$$x=2 \Rightarrow y = 0.8 + 2.4 = 3.2 \Rightarrow (2, 3.2)$$

$$x=3 \Rightarrow y = 1.2 + 2.4 = 3.6 \Rightarrow (3, 3.6)$$

$$x=4 \Rightarrow y = 1.6 + 2.4 = 4.0 \Rightarrow (4, 4.0)$$

$$x=5 \Rightarrow y = 2.0 + 2.4 = 4.4 \Rightarrow (5, 4.4)$$



Loglinear Model:-

Loglinear Models allows a higher dimensional data space to be constructed from lower dimensional spaces.

Regression and loglinear models can be applied on a sparse data and skewed data.

Dimensionality Reduction:-

31/01/19

1. Wavelets transform.

2. PCA.

3. Attribute Subset Selection.

Attribute Subset Selection:-

Extracting the subset of the attributes from the given 'n' no. of attributes can be called as

Attribute Subset Selection

There are 4 procedures for attribute subset selection

1. Stepwise forward selection.

2. Stepwise Backward Elimination.

3. Combination of Stepwise forward and Backward Selection.

4. Decision tree induction.

Stepwise forward Selection:-

1. The procedure starts with empty set.

2. Identify the important attribute and add it to the previous set.

Ex:- $\{A_1, A_2, A_3, A_4, A_5, A_6, A_7\}$. Assume $\{A_1, A_4, A_6\}$ are required.

$$S_1: \emptyset$$

$$S_2: \{A_1\}$$

$$S_3: \{A_1, A_4\}$$

$$S_4: \{A_1, A_4, A_6\}$$

Stepwise Backward Selection:-

1. Procedure starts with full set.

2. Remove the unnecessary attributes from the remaining set at each and every iteration.

Ex:- $\{A_1, A_2, A_3, A_4, A_5, A_6, A_7\}$ Assume $\{A_1, A_4, A_6\}$ are required.

$$S_1: \{A_1, A_2, A_3, A_4, A_5, A_6\}$$

$$S_2: \{A_1, A_2, A_3, A_4, A_6\}$$

$$S_3: \{A_1, A_2, A_4, A_6\}$$

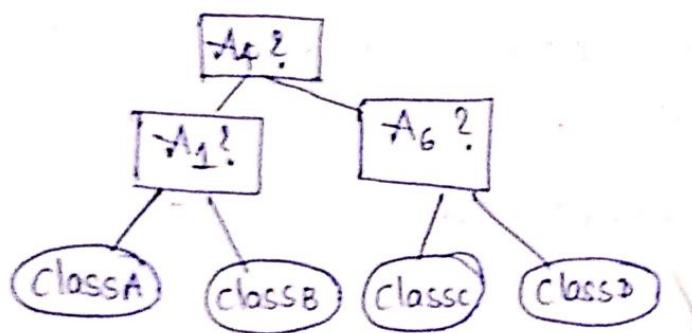
$$S_4: \{A_1, A_4, A_6\}$$

Combination of stepwise forward and backward selection.
At each step the procedure selects the best attribute and remove the worst from among the remaining attributes.

Decision Tree Induction:-

Decision Tree Induction constructs a flow chart like a tree structure.

Ex:-



Wavelet Transform:-

Wavelet Transform is a non-linear Signal processing Technique when it is applied to the data vector it transform into numerically a different data vector called ' \hat{x} '.

The elements of \hat{x} can be called as wavelet coefficients. Select the strongest wavelet coefficients from the \hat{x} .

Principle Components:-

Principle Component analysis can be used for data reduction.

Processor for identifying principle Components.

Procedure for PCA:-

1. Find the covariance matrix for the given two data points.

2. Find the covariance matrix for the given data i.e. covariance matrix for the attributes x, y can be

$$\begin{matrix} x \\ y \end{matrix} \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} \text{ where } \text{cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

3. Find the Eigen values $|S - \lambda I| = 0$. λ_1, λ_2 can be called as eigen values.

3. Find the eigen vector namely $\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$

4. Select the eigen vector values which is corresponding to the highest eigen vector.

$$3 \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

5. The selected eigen vector values can be called as principle components.

* Find the principle components for the given data.

x	2	1	0	-1
y	4	3	1	0.5

x	y	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$
2	4	1.5	1.875	2.812
1	3	0.5	0.87	0.435
0	1	-0.5	-1.125	0.562
-1	0.5	-1.5	-1.62	2.43
$\bar{x} = 0.5$		$\bar{y} = 2.125$		$E = 6.239$

$$\text{cov}(x,y) = \frac{6.239}{3} = 2.079$$

$$\text{Cov}(x, x) = \sum_{i=1}^n \frac{(x_i - \bar{x})(x_i - \bar{x})}{n-1} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} = 1.60$$

$$\text{Cov}(x, x) = 5/3 = 1.66 \quad \text{Cov}(y, y) = \sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n-1} = \frac{8.16}{3} = 2.72$$

Now the covariance matrix is $S = \begin{bmatrix} 1.60 & 2.079 \\ 2.079 & 2.72 \end{bmatrix}$

Step 2: Find eigen values

$$|S - \lambda I| = 0$$

$$\begin{vmatrix} 1.66 - \lambda & 2.079 \\ 2.079 & 2.72 - \lambda \end{vmatrix} = 0$$

$$\Rightarrow (1.66 - \lambda)(2.72 - \lambda) - (2.079)^2 = 0$$

$$\Rightarrow \lambda^2 - 4.38\lambda + 0.20 = 0$$

$$\lambda_1 = 4.34$$

$$\lambda_2 = 0.056.$$

Step 3: Find eigen vectors

Substitute λ_1 in $(S - \lambda_1 I)$

$$\begin{bmatrix} 1.66 - 4.34 & 2.079 \\ 2.079 & 2.72 - 4.34 \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = 0$$

$$\begin{bmatrix} -2.68 & 2.079 \\ 2.079 & -1.62 \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = 0$$

$$-2.68 a_{11} + 2.079 a_{12} = 0 \rightarrow ①$$

$$2.079 a_{11} - 1.62 a_{12} = 0 \rightarrow ②$$

Orthogonal projects on $a_{11}^2 + a_{12}^2 = 1$

$$a_{12} = 1.029 a_{11} \rightarrow ③$$

$$a_{11}^2 + (1.029 a_{11})^2 = 1$$

$$a_{11}^2 + 1.064 a_{11}^2 = 1 \Rightarrow a_{11}^2 (1 + 1.064) = 1$$

$$\hat{a}_{11} = \frac{1}{\sqrt{2.664}} \Rightarrow a_{11} = \sqrt{\frac{1}{2.664}} = 0.64.$$

$$a_{12} = 1.29 \times 0.64 = 0.825$$

Step 4:- Substitute λ_2 in $|S - \lambda I|$

$$\begin{bmatrix} 1.64 - 0.05 & 0.079 \\ 0.079 & 2.72 - 0.03 \end{bmatrix} \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix} = 0$$

$$a_{21} = 0.743$$

$$a_{22} = 0.829$$

$$\begin{bmatrix} 4.34 \\ 0.05 \end{bmatrix} = \begin{bmatrix} 0.64 & 0.729 \\ 0.743 & 0.829 \end{bmatrix}$$

- From this a_{11}, a_{12} are principal Components.

Dissimilarity measures for finding numeric attributes:-

There are '4' methods for finding dissimilarity b/w numeric attributes.

1) Euclidean distance.

2) Manhattan distance.

3) Minkowski distance.

4) Supremum distance.

Let $i = (x_{i1}, x_{i2}, x_{i3}, \dots, x_{in})$

$j = (x_{j1}, x_{j2}, x_{j3}, \dots, x_{jn})$

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

2) Manhattan distance,

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

3) Minkowski distance,

$d(i, j) = \sqrt[b]{(x_{i1} - x_{j1})^b + (x_{i2} - x_{j2})^b + \dots + (x_{in} - x_{jn})^b}$

4) Supremum distance,
 $d(i,j) = \lim_{h \rightarrow 0} \left(\sum_{j=1}^h |x_{ij} - x_{jf}|^h \right)^{1/h}$ i.e., the attribute that gives the max difference in values of the two objects.

Eg:- Let $x_1 = 1, 2$

$x_2 = 3, 5$ are the given two points.

1) Euclidean distance:-

$$\text{Euclidean distance } d(x_1, x_2) = \sqrt{(1-3)^2 + (2-5)^2}$$

$$= \sqrt{(-2)^2 + (-3)^2}$$

$$= \sqrt{4+9} = \sqrt{13} = 4$$

2) Manhattan distance:-

$$d(x_1, x_2) = |1-3| + |2-5| = |-2| + |-3| = 2+3 = 5$$

3) Minkowski distance:-

$$d(x_1, x_2) = \sqrt[2]{(1-3)^2 + (2-5)^2}$$

$$= \sqrt{4+9} = \sqrt{13}$$

4) Supremum distance-

$$\max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|)$$

$$= \max(|1-3|, |2-5|)$$

$$= \max(|-2|, |-3|)$$

$$= \max(2, 3)$$

=

In data compression transformations are applied to obtain a reduced and compressed representation of the original data.

If the original data can be reconstructed from the compressed data without any information loss then the data reduction is called lossless.

If we can reconstruct only an approximation of the original data then the data reduction is called lossy.

Data Transformation:-

Converting the data into its appropriate form for mining can be called as data transformation.

Strategies include in data transformation. They are several strategies:-

- (i) Smoothing
- (ii) Attribute Construction
- (iii) Aggregation
- (iv) Normalisation.
- (v) Discretization.
- (vi) Concept hierarchy generation.

(i) Smoothing:- To remove the noisy in the data. Techniques include binning, regression & clustering.

(ii) Attribute Construction:- When the new attributes are constructed & added from the given set of attributes to help the mining process.

(iii) Aggregation:- Where summary or aggregation operations are applied to the data.

Eg:- The daily sales data may be aggregated to conclude monthly & annual total amount.

(iv) Normalisation:- where the attribute data (or) scale so, as to form within a smaller range such as (0.0 to 1.0) (-1.0 to 1.0)

(v) Discretization:- where the raw values of numeric attribute are replaced by interval labels (0-10, 11-20, etc.)

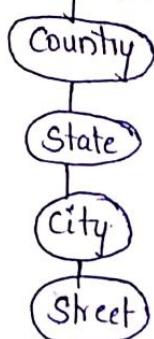
(vi)

Conceptual labels (Eg:- youth, adult, senior).

(vi) Concept hierarchy generation:- where the attributes such as street can be generalised to higher level concepts like city (or) country.

Street is nothing but lower level concept.

Eg:-



Data transformation by normalisation:-

Techniques include for Normalisation are:-

- 1) Min Max Normalisation.
- 2) \bar{z} -Score Normalisation.
- 3) Decimal Scaling.

1) Min Max Normalisation:-

It performs a linear transformation with original data.

Suppose that, min_A & max_A are the min & max values of an attribute A.

→ map a value v_i of A to v_i^1 in the range [new minA, new maxA] by computing $v_i^1 = v_i - \frac{\text{min}_A}{\text{max}_A - \text{min}_A} \times (\text{new maxA} - \text{new minA})$.

Eg:- let A be the numeric attribute with observed values $v_1, v_2, v_3, \dots, v_n$.

→ The min & max values for the attribute "income" are

\$12,000 & \$98,000 respectively.

→ Map income to the range [0.0, 1.0]

→ A value v_i for income attribute is \$73,600.

$$v_i = \$73,600$$

$$\text{min}_A = \$12,000$$

$$\text{max}_A = \$98,000$$

$$n - \text{min}_A = 0.0$$

$$n - \text{max}_A = 1.0$$

$$v_i^1 = \frac{73,600 - 12,000}{98,000 - 12,000} (1.0 + 0.0) + 0.0$$

$$v_i^1 = 0.716$$

2) τ -Score Normalization:-

Zero mean Normalisation- The values for an attribute A are normalised based on the mean & standard deviation.

τ -score normalisation transforms a value v_i of A into v_i^1 by computing

$$v_i^1 = \frac{v_i - \bar{A}}{\sigma_A}$$

Ex:- The mean & standard deviation of the values attribute income are \$ 54,000 & \$ 10,000 respectively.

$$v_i^1 = \text{when } v_i = \$73,600.$$

$$v_i^1 = \frac{73,600 - 54,000}{10,000}$$

$$v_i^1 = 1.225.$$

Normalization by Decimal Scaling:-

6/2/19

Decimal Scaling normalises by moving the decimal point of values of attribute A.

The no. of decimal points moved based on the maximum absolute value of A.

A value v_i of A is normalized to v_i^1 by computing v_i^1 , where 'j' is the smallest integer such that $\max(v_i^1) < 1^{10}$.

Ex:- Suppose the recorded values of A range from -986 to 917. Then find the normalized values for -986 to 917.

$$v_i^1 = \frac{-986}{1000} = 0.986, \quad v_i^1 = \frac{917}{1000} = 0.917.$$

Normalized values of A is 0.986 to 0.917.

Let the attribute income contains the values

T-Score Normalisation:-

The variation of T-score normalization replaces the standard deviation by mean absolute deviation of A.

The mean absolute deviation of A is denoted by

$$S_A = \frac{|v_1 - \bar{A}| + |v_2 - \bar{A}| + |v_3 - \bar{A}| + \dots + |v_n - \bar{A}|}{n}$$

$$S_{V_i}^1 = \frac{|v_i - \bar{A}|}{S_A}$$

Similarity measures Cosine Similarity:-

It is a measure of similarity that can be used to compare documents.

Let x, y be the two vectors for comparison.

Then Similarity of $x, y = \frac{x \cdot y}{\|x\| \cdot \|y\|}$

Normalization $\|x\|$

where the norm $\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$

& norm $\|y\| = \sqrt{y_1^2 + y_2^2 + \dots + y_n^2}$

Let there are '5' documents & the words that are commonly appear in all the documents are appeared below.

so now calculate as to which words are common in all the documents.

so now calculate as to which words are common in all the documents.

so now calculate as to which words are common in all the documents.

so now calculate as to which words are common in all the documents.

so now calculate as to which words are common in all the documents.

so now calculate as to which words are common in all the documents.

team	coach	bacley	baseball	soccer	penalty	Score, win
Doc1	5	0	3	0	2	0 2
Doc2	3	0	2	0	1	0 1
Doc3	0	4	0	2	1	0 3
Doc4	0	1	0	0	1	2 0

Q8 Assume that Doc1 is x and Doc2 is y then

$$x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

$$x^*y = 5 \cdot 3 + 3 \cdot 2 + 2 \cdot 1 + 2 \cdot 1 = 25$$

$$\|x\| = \sqrt{5^2 + 3^2 + 2^2 + 2^2} = \sqrt{25 + 9 + 4 + 4} = \sqrt{42} \approx 6.48$$

$$\|y\| = \sqrt{3^2 + 2^2 + 1^2 + 1^2 + 1^2 + 1^2} = \sqrt{12} = 4.12$$

$$\text{Similarity}(Doc1, Doc2) = \frac{25}{6.48 \times 4.12} = 0.986$$

7/02/19

Dissimilarity measure for nominal attribute:-

If an attribute can take two or more states then that attribute can be called as nominal attribute.

The dissimilarity of an attribute can be represented by dissimilarity matrix.

Ex:-

The dissimilarity matrix for the given '4' objects and one attribute

$$\begin{bmatrix} d(1,1) & & & \\ d(2,1) & d(2,2) & & \\ d(3,1) & d(3,2) & d(3,3) & \\ d(4,1) & d(4,2) & d(4,3) & d(4,4) \end{bmatrix}$$

- where $d(i,j) = 1$, means two objects are dissimilar
 $d(i,j) = 0$ means the '2' objects are not dissimilar
- The Similarity of $\text{sim}(i,j) = 1 - d(i,j)$
- The $\text{sim}(i,j) = 1$, means the '2' objects are similar.
 $\text{sim}(i,j) = 0$ means the '2' objects are not similar.
- The dissimilarity of '2' objects for a nominal attribute.

$$d(i,j) = \frac{P-m}{P}$$

where 'P' is no. of attributes & 'm' is no. of matches.

Problem:- Find the dissimilarity matrix for the given nominal data.

Id	Type1	Type2
1	10	A
2	30	B
3	10	A
4	20	D

Sol:

$$d(1,1) = 0$$

$$d(2,1) = \frac{2-0}{2} = \frac{2}{2} = 1 \quad [P=2, m=0]$$

$$d(2,2) = 0$$

$$d(3,1) = \frac{2-2}{2} = \frac{0}{2} = 0. \quad [P=2, m=2]$$

$$d(3,2) = \frac{2-0}{2} = \frac{2}{2} = 1 \quad [m=0, P=2]$$

$$d(3,3) = 0$$

$$d(4,1) = \frac{2-0}{2} = \frac{2}{2} = 1 \quad [P=2, m=0]$$

$$d(4,2) = \frac{2-0}{2} = \frac{2}{2} = 1 \quad [P=2, m=0]$$

$$d(4,3) = \frac{2-0}{2} = \frac{2}{2} = 1 \quad [P=2, m=0]$$

$$d(4,4) = 0$$

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0 & 1 & 0 & \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

From the above data matrix all the objects are dissimilar except the objects 3, 4.

Ex:- Find the dissimilarity matrix from a given data.

Obj	Test 1.
1	codeA
2	codeB
3	codec
4	codeA

$$d(1,1) = \frac{1-1}{1} = 0$$

$$d(2,1) = \frac{1-0}{1} = 1 \quad (P=1, m=0)$$

$$d(2,2) = 0$$

$$d(3,1) = \frac{1-0}{1} = 1 \quad (P=1, m=0)$$

$$d(3,2) = \frac{1-0}{1} = 1 \quad (P=1, m=0)$$

$$d(3,3) = 0$$

$$d(4,1) = \frac{1-1}{1} = 0 \quad (P=1, m=1)$$

$$d(4,2) = \frac{1-0}{1} = 1 \quad (P=1, m=0)$$

$$d(4,3) = \frac{1-0}{1} = 1 \quad (P=1, m=0)$$

$$d(4,4) = 0$$

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

From the above matrix all the attributes are dissimilar except object 4, 1.

Dissimilarity measure for ordinal attribute:-

The values of an ordinal attribute have a meaningful order or ranking above the attribute.

Ex:- 1. The sequence of small, medium, large.

2. Grade A, Grade B, Grade C

3. Excellent, good, average.

Dissimilarity of an ordinal attribute can be computed by the following procedure.

1. Give the rank for the given states.

2. Normalize the given rankings to $[0.0, 1.0]$ by computing

$$\bar{r}_{if} = \frac{r_{if} - 1}{m_f - 1} \quad m_f = \text{no. of states}, r_{if} = \text{ranking for a state.}$$

3. Compute $d(i,j)$ by $d(i,j) = |\text{Rank of } i - \text{Normalized Rank of } j|$.

Find the dissimilarity matrix for the given data.

Object	test
1	Excellent
2	fair
3	good
4.	Excellent

i) Excellent - 1, good - 2, fair - 3 (ranking)

$$2) \bar{r}_{if} \text{ for excellent} = \frac{1-1}{3-1} = 0. \quad m_f = 3, r_{if} = 1$$

$$\bar{r}_{if} \text{ for good} = \frac{2-1}{3-1} = \frac{1}{2} = 0.5 \quad m_f = 3, r_{if} = 2$$

$$\bar{r}_{if} \text{ for fair} = \frac{3-1}{3-1} = \frac{2}{2} = 1 \quad m_f = 3, r_{if} = 3$$

$$\therefore [1-0, 2-0.5, 3-1]$$

$$3) d(1,1) = 0. \quad d(3,1) = |(1RG - NRE)|$$

$$d(2,1) = |1-0| = 1. \quad = |0.5-0| = 0.5$$

$$d(1,2) = 0$$

$$d(3,1) = |NRG - NRE| = |0.5 - 1| = 0.5$$

$$d(3,3) = 0.$$

$$d(4,1) = |NRE - NRG| = 0 - 0 = 0.$$

$$d(4,2) = |NRE - NRF| = |0 - 1| = 1$$

$$d(4,3) = |NRE - NRG| = |0 - 0.5| = 0.5 \quad d(4,4) = 0.$$

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1 & 0.5 & 0 \end{bmatrix}$$

All the objects are dissimilar except the objects 1,4.

From the above matrix, the objects 1 and 2 are most dissimilar and also the objects 4 and 2 are most dissimilar.

Dissimilarity measure for numeric data:-

9/02/19

$$\text{The dissimilarity } d_{ij}(t) = \frac{|x_{if} - x_{jf}|}{\max - \min}$$

Find the dissimilarity matrix for the given data.

Objects attribute(test-3)

1	45
2	22
3	64
4	28

$$d_{11} = 0$$

$$d_{21} = \frac{|22 - 45|}{64 - 22} = \frac{23}{42} = 0.547$$

$$d_{22} = 0$$

$$d_{31} = \frac{|64 - 45|}{64 - 22} = \frac{19}{42} = 0.452$$

$$d_{32} = \frac{|64 - 22|}{64 - 22} = \frac{42}{42} = 1$$

$$d_{33} = 0$$

$$d_{41} = \frac{|28 - 45|}{64 - 22} = \frac{17}{42} = 0.40$$

$$d_{42} = \frac{|28 - 22|}{64 - 22} = \frac{6}{42} = 0.142$$

$$d_{43} = \frac{|28 - 64|}{64 - 22} = \frac{36}{42} = 0.857$$

$$d_{44} = 0$$

From the above data matrix, all the objects are dissimilar except the objects 3 and 2 are similar.

* Calculating the dissimilarity for mixed data:

$$d_{ij} = \frac{\sum_{f=1}^P \delta_{ij}^{(f)} \cdot d_{ij}^{(f)}}{\sum_{f=1}^P \delta_{ij}^{(f)}}$$

where $\delta_{ij}^{(f)} = 0$, if x_{if} or x_{jf} are missing

(or)
 x_{if} and x_{jf} are zero otherwise it's
 corresponding value is 1.

Find the dissimilarity matrix for mixed data.

Object	(nominal)	(ordinal)	(numeric)
	test1	test2	test3
1	codeA	Excellent	45
2	codeB	Fair	44
3	codec	Good	64
4	codex	Excellent	28

$$\begin{bmatrix} 0 \\ 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 & 0 \\ 0.5 & 0.5 & 0 \\ 0 & 1 & 0.5 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0.55 & 0 \\ 0.45 & 1 & 0 \\ 0.40 & 0.14 & 0.86 & 0 \end{bmatrix}$$

$$d_{21} = \frac{1 \times 1 + 1 \times 1 + 1 \times 0.55}{1+1+1} = \frac{1+1+0.55}{3} = \frac{2.55}{3} = 0.8$$

$$d_{31} = \frac{1 \times 1 + 0.5 + 1 \times 0.45}{3} = \frac{1+0.5+0.45}{3} = \frac{1.95}{3} = 0.65$$

$$d_{32} = \frac{1 \times 1 + 1 \times 0.5 + 1 \times 1}{3} = \frac{1+0.5+1}{3} = \frac{2.5}{3} = 0.83$$

$$d_{42} = \frac{1 \times 1 + 1 \times 1 + 1 \times 0.14}{3} = \frac{2.14}{3} = 0.71, \quad d_{41} = \frac{0+0+0.4}{3} = 0.13$$

$$d_{43} = \frac{1 \times 1 + 1 \times 0.5 + 1 \times 0.86}{3} = \frac{1+0.5+0.86}{3} = \frac{2.36}{3} = 0.78$$

$$= \begin{bmatrix} 0 \\ 0.85 & 0 \\ 0.65 & 0.83 & 0 \\ 0.13 & 0.71 & 0.78 & 0 \end{bmatrix}$$

Cosine Similarity:-

* Find the dissimilarity between the documents

	team	coach	hockey	baseball	soccer	penalty	score	win	loss	sec
doc1	3	0	3	0	2	0	0	2	0	0
doc2	4	0	2	0	1	1	0	1	0	1
doc3.	0	8	0	3	1	0	0	3	0	1
doc4	1	2	0	0	1	3	2	0	4	0

$$\text{Cos}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$$\text{Where } \|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

$$\|y\| = \sqrt{y_1^2 + y_2^2 + \dots + y_n^2}$$

Sol:- Assume that doc1 & doc2 are x & y respectively

$$\text{doc1} = x = (3, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$\text{doc2} = y = (4, 0, 2, 0, 1, 0, 1, 0, 1)$$

$$x \cdot y = (3 \times 4 + 3 \times 2 + 2 \times 1 + 2 \times 1 + 0 \times 1) = 12 + 6 + 2 + 2 = \underline{\underline{22}}$$

$$\text{Norm} \|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

$$= \sqrt{3^2 + 3^2 + 2^2 + 2^2} = \sqrt{9+9+4+4} = \sqrt{26} = 5.09.$$

$$\text{Norm} \|y\| = \sqrt{16+4+1+1+1+1} = \sqrt{24} = 4.89.$$

$$\text{Sim}(\text{doc}_1, \text{doc}_2) = \frac{22}{5.09 \times 4.89} = 0.88$$

\therefore The '2' doc's are more likely to be similar.

Assume that doc_1 and doc_3 are P, q respectively.

$$\text{doc}_3 = P = (0, 1, 8, 0, 1, 3, 1, 0, 0, 3, 0, 0)$$

$$\text{doc}_1 = P = (3, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$P = (0+0+0+0+2+0+0+6+0+0) = 8$$

$$\text{Norm} \|P\| = \sqrt{3^2 + 3^2 + 2^2 + 2^2} = \sqrt{9+9+4+4} = \sqrt{26} = 5.09.$$

$$\text{Norm} \|q\| = \sqrt{8^2 + 3^2 + 1^2 + 3^2} = \sqrt{64+9+1+9} = \sqrt{83} = 9.11$$

$$\text{Sim}(\text{doc}_1, \text{doc}_3) = \frac{8}{5.09 \times 9.11} = \frac{8}{46.36} = 0.17$$

Similarity between binary attributes:-

Find the patients which are having similar disease.

	gender	fever	cough	test1	test2	test3	test4
jack	M	N	Y	P	N	N	N
jim	M	Y	N	N	N	P	N
mary	F	Y	Y	P	N	N	N

Sol:- Convert the given data set into binary format.

	gender	fever	cough	test1	test2	test3	test4
jack	M	0	1	1	0	0	0
jim	M	1	0	0	0	1	0
mary	F	1	1	1	0	0	0

$$disSim(i,j) = \frac{r+s}{r+s+t}$$

		object j	object i
		1	0
object i	1	r	s
	0	t	u

$$dis(jack, jim) = \frac{2+2}{0+2+2} = 1$$

$$r=0$$

$$s=2$$

$$t=2$$

$$sim(i,j) = 1 - disSim(jack, jim) = 1 - 1 = 0$$

$$dis(jack, mary) = \frac{0+1}{2+0+1} = 1/3 = 0.3$$

$$r=2$$

$$s=0$$

$$t=1$$

$$u=3$$

$$sim(i,j) = 1 - dis(jack, mary) = 1 - 0.3 = 0.7$$

There is similarity between jack and mary, they have similar disease.

Association Rule:-

Find the frequent item sets that are frequently purchased together for the given dataset.

TransID List of Items

Tid	IDe
T ₁	I ₁ , I ₂ , I ₃
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

Apriori Algorithm:-

- Scan the dataset D for count of each candidate.

C ₁	L ₁
Itemset	Sup.Count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Scan D for count of each candidate

Compare Sup.Count with min. Sup Count

Itemsets	Support count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Sup = Support

Let min. supcount = 2

- Generate C₂ from L₁

$$C_2 = L_1 \bowtie L_1$$

C ₂	C ₂	Generate G ₂ from L ₁
I ₁ , I ₂		
I ₁ , I ₃		
I ₁ , I ₄		
I ₁ , I ₅		
I ₂ , I ₃		
I ₂ , I ₄		
I ₂ , I ₅		
I ₃ , I ₄		
I ₃ , I ₅		
I ₄ , I ₅		

Scan D for count of each candidate.

Itemsets	Supcount
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

Itemset	Supp.count
I ₁ , I ₂	4
I ₁ , I ₃	4
I ₁ , I ₅	2
I ₂ , I ₃	4
I ₂ , I ₄	2
I ₂ , I ₅	2

Compare the Sup.count with min Sup.count.

Step 3: Generate C₃ from L₂. C₃

C ₃	Itemset	Supp.count	L ₃
I ₁ , I ₂ , I ₃	I ₁ , I ₂ , I ₃	2	
I ₁ , I ₂ , I ₅	I ₁ , I ₂ , I ₅	2	
I ₁ , I ₂ , I ₄	I ₁ , I ₂ , I ₄	1	
I ₁ , I ₃ , I ₅	I ₁ , I ₃ , I ₅	1	
I ₂ , I ₃ , I ₄	I ₂ , I ₃ , I ₄	0	
I ₂ , I ₃ , I ₅	I ₂ , I ₃ , I ₅	1	
I ₂ , I ₄ , I ₅	I ₂ , I ₄ , I ₅	0	

Step 4: Generate C₄ from L₃.

C ₄	Itemset	Supp.count
I ₁ , I ₂ , I ₃ , I ₅	I ₁ , I ₂ , I ₃ , I ₅	1

From the given dataset, the frequent purchased items are I₁, I₂, I₃ and I₁, I₂, I₅.

* TID	List of Items
T ₁	L ₁ , L ₂ , L ₃ , L ₄ , L ₅
T ₂	L ₄ , L ₅
T ₃	L ₂ , L ₄
T ₄	L ₁ , L ₃
T ₅	L ₁ , L ₃ , L ₅
T ₆	L ₁ , L ₄ , L ₅

min count = 2.

1. Scan the dataset 'D' for count of each candidate.

Itemset	Support Count
L ₁	5
L ₂	2
L ₃	3
L ₄	3
L ₅	4

Compare the support count with minSupCount

Itemset	Support Count
L ₁	5
L ₂	2
L ₃	3
L ₄	3
L ₅	4

2. Generate C₂ from L₁, c₂

Itemset	Support Count
L ₁ , L ₂	1
L ₁ , L ₃	3
L ₁ , L ₄	2
L ₁ , L ₅	4
L ₂ , L ₃	1
L ₂ , L ₄	2
L ₂ , L ₅	1
L ₃ , L ₄	1
L ₃ , L ₅	2
L ₄ , L ₅	2

Scan 'D' for count of each candidates

Generate C₂ from L₁

Itemset	Support Count
L ₁ , L ₃	3
L ₁ , L ₄	2
L ₁ , L ₅	4
L ₂ , L ₄	2
L ₃ , L ₅	2
L ₄ , L ₅	2

3.

Itemset	Support Count
L ₁ , L ₃ , L ₄	1
L ₁ , L ₃ , L ₅	2
L ₁ , L ₄ , L ₅	2
L ₂ , L ₃ , L ₄	1
L ₁ , L ₂ , L ₄	1
L ₃ , L ₄ , L ₅	1

Scan 'D' for count of each candidate

Generate C₃ from L₂

Itemset	Support Count
L ₁ , L ₃ , L ₅	2
L ₁ , L ₄ , L ₅	2

4.

Itemset	Support Count
L ₁ , L ₃ , L ₄ , L ₅	1

Scan 'D' for count of each candidate

From the given dataset, the frequent purchased items are L₁, L₃, L₅ and L₁, L₄, L₅.

Generating Association Rules from frequent statements are:

$\{I_1, I_2, I_5\}$

$\{I_1, I_2, I_3\}$

Association Rule for	I_1, I_2, I_5	Support count	Confidence	Confidence percentage
$\{I_1, I_2\} \rightarrow I_5$	2	$2/4 = 0.5$	50%	
$I_1, I_5 \rightarrow I_2$	2	$2/2 = 1$	100%	
$I_2, I_5 \rightarrow I_1$	2	$2/2 = 1$	100%	
$I_1 \rightarrow I_2, I_5$	2	$2/6 = 0.33$	33%	
$I_2 \rightarrow I_1, I_5$	2	$2/7 = 0.28$	28%	
$I_5 \rightarrow I_1, I_2$	2	$2/4 = 0.5$	50%	100%

$$\text{Confidence}(A \rightarrow B) = P(B|A) = \frac{\text{Support count}(A \cup B)}{\text{Support count}(A)}$$

$$\text{Confidence}(I_1, I_2 \rightarrow I_5) = \frac{\text{Support count}(I_1, I_2, I_5)}{\text{Support count}(I_1, I_2)} = 2/4$$

$$\text{Confidence}(I_1, I_5 \rightarrow I_2) = \frac{\text{Support count}(I_1, I_5, I_2)}{\text{Support count}(I_1, I_5)} = 2/2 = 1$$

$$\text{Confidence}(I_2, I_5 \rightarrow I_1) = \frac{\text{Support count}(I_2, I_5, I_1)}{\text{Support count}(I_2, I_5)} = 2/2 = 1$$

$$\text{Confidence}(I_1 \rightarrow I_2, I_5) = \frac{\text{Support count}(I_1, I_2, I_5)}{\text{Support count}(I_1)} = 2/6 = 0.33$$

$$\text{Confidence}(I_2 \rightarrow I_1, I_5) = \frac{\text{Support count}(I_1, I_2, I_5)}{\text{Support count}(I_2)} = 2/7 = 0.28$$

$$\text{Confidence}(I_5 \rightarrow I_1, I_2) = \frac{\text{Support count}(I_1, I_2, I_5)}{\text{Support count}(I_5)} = 2/2 = 1$$

2) Given the itemsets milk(m), coke(c), jam(j), pepsi(p), Bread(b)

Find the strong association rules when the minimum support count = 3 and confidence = 70%.

TID	Itemsets
T ₁	M, C, j
T ₂	M, P, j
T ₃	M, j
T ₄	C, B
T ₅	B, j, C
T ₆	M, B, j
T ₇	C, B, j
T ₈	B, C
T ₉	M, C, B, j
T ₁₀	M, P, C, j

Sol:- Step 1:- Scan the 'D' dataset for count of each candidate.

Itemset	Support count	Itemset	Support count
M	6	M	6
C	7	C	7
j	8	j	8
P	2	B	6
B	6		

Scan 'D' for count of each candidate.

Compare

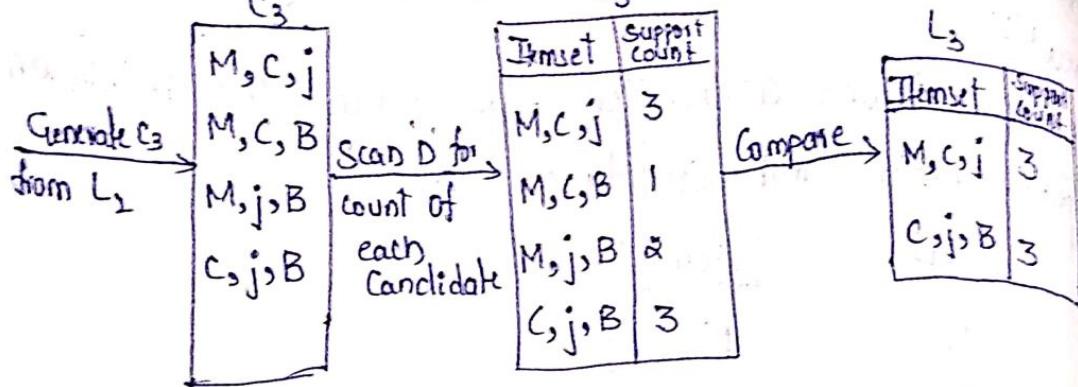
Q. Generate C_2 from L_1 .

Itemset	Itemset	Itemset
Generate C_2 from L_1	M, C	M, C
M, j	Scan D for count of each Candidate.	M, j
M, B		M, B
C, j		C, j
C, B		C, B
j, B		j, B

Support count

Itemset	Support count
M, C	3
M, j	6
M, B	2
C, j	5
C, B	5
j, B	4

3. Generate C_3 from L_2 .



4. Generate C_4 from L_4 .

C_4	
M, C, j, B	

Generate C_4 from L_4

Itemset	Support Count
M, C, j, B	1

Scan 'D' for count of each candidate

From the given dataset, the frequent purchased items are M, C, j and C, j, B .

Partitioning Algorithm for frequent itemsets:- 19/02/19.

Transaction	Items
T_1	Bread, Jelly, peanut, Butter (PNB)
T_2	Bread, peanut, Butter (PNB)
T_3	Bread, Milk, peanut, Butter (PNB)
T_4	Beer, Bread
T_5	Beer, Milk

- 1. Divide the given transactions into '2' equal parts.
- 2. Solve the equal parts by applying apriori algorithm.
- 3. Combine all the transactions.

Sol:- 1. Divide the transactions T_1 , & T_2 in one group and T_3 , T_4 , T_5 in another group.

Item	Support Count	Item	Support Count
Bread	2	Bread	2
Jelly	1	Jelly	1
Peanut Butter	2	Peanut Butter	2

→ Compare the MinSup count with Support count

Step 2:-

Itemset	Support count	Itemset	Support count
Bread, Jelly	1	Bread, Jelly	1
Bread, Peanut Butter	2	Bread, Peanut Butter	2
Jelly, Peanut Butter	1	Jelly, Peanut Butter	1

→ Compare

Step 3:-

Itemset	Support count	Itemset	Support count
Bread, Jelly, Peanut Butter	1	Bread, Jelly, Peanut Butter	1

→ Compare

$$L_1 = \{ \{B\}, \{J\}, \{PB\}, \{B, J\}, \{B, PB\}, \{J, PB\}, \{B, J, PB\} \}$$

Perform the Apriori Algorithm for T_3, T_4, T_5 .

Step 1:-

Scan D
for count of each candidate

Itemset	Support count
Bread	2
Milk	2
Beer	2
Peanut Butter	1

Compare

Itemset	Support count
Bread	2
Milk	2
Peanut Butter	1
Beer	2

Step 2:-

Itemset	Support Count
Bread, Milk	1
Bread, PNB	1
Bread, Beer	1
Milk, PNB	1
Milk, Beer	1
PNB, Beer	0

Compare Support
with min support

Itemset	Support Count
Bread, Milk	1
Bread, PNB	1
Bread, Beer	1
Milk, PNB	1
Milk, Beer	1

Step 3:-

Itemset	Support count.
Bread, Milk, PNB	1
Bread, Milk, Beer	0
Bread	0
Milk, PNB, Beer	0

Compare.

Itemset	Support count
Bread, Milk, PNB	1

$$L_2 = \{ \{B\}, \{M\}, \{PNB\}, \{Beer\}, \{B, M\}, \{B, PNB\}, \{B, Beer\}, \{M, PNB\}, \{M, Beer\}, \{B, M, PNB\}, \{B, M, Beer\}, \{M, PNB, Beer\} \}$$

Measuring the quality of Association Rules:- 20/02/19

The rules for measuring the quality of an Association rules are as follows:

1. Support
2. Confidence
3. Interest (Lift)
4. Conviction
5. Chi-square

Support:-

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

Confidence:-

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

Interest (or) lift:-

$$\text{interest}(A \Rightarrow B) = \frac{P(A, B)}{P(A) \cdot P(B)}$$

Conviction (or) lift:-

$$\text{Conviction}(A \Rightarrow B) = \frac{P(A) \cdot P(\neg B)}{P(A, \neg B)}$$

If the conviction value of $\neg A, B$ is 1 then $A \& B$ are not related.

Chi-square:-

Chi-square for the itemset x is calculating

$$\chi^2 = \sum_{x \in I} \frac{(O(x) - E(x))^2}{E(x)} \text{ where } I = \{I_1, I_2, \dots, I_n\}$$

The chi-square value is zero if all the values are independent.

	B	\bar{B}	Total
A	15	10	25
\bar{A}	55	20	75
Total	70	30	100

$$e_{11} = \frac{70 \times 25}{100} = \frac{70}{4} = 17.5$$

$$e_{12} = \frac{30 \times 25}{100} = \frac{30}{4} = 7.5$$

$$e_{21} = \frac{70 \times 75}{100} = \frac{210}{4} = 52.5$$

$$e_{22} = \frac{30 \times 75}{100} = \frac{90}{4} = 22.5$$

$$\chi^2 = \frac{(15-17.5)^2}{17.5} + \frac{(10-7.5)^2}{7.5} + \frac{(55-52.5)^2}{52.5} + \frac{(20-22.5)^2}{22.5}$$

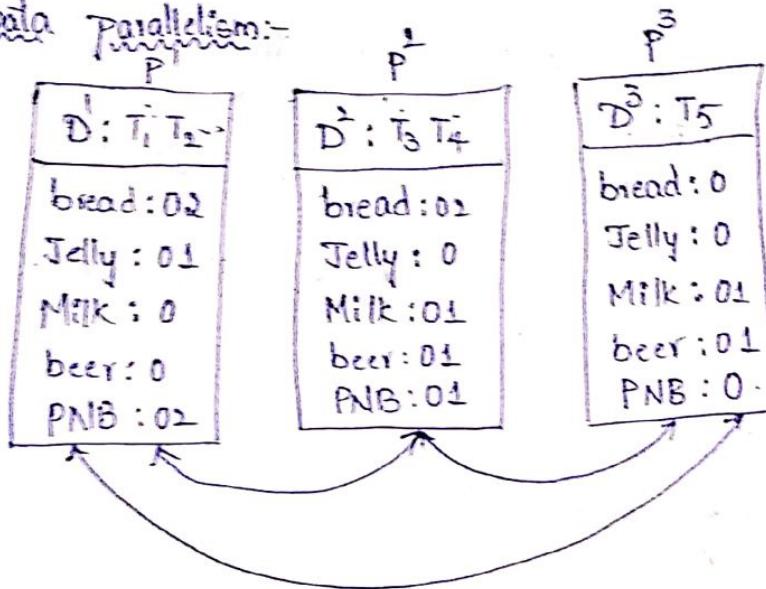
Parallel and Distributed Algorithms

25/02/19

Database:

- T₁ bread, Jelly, PNB
- T₂ bread, PNB
- T₃ bread, Milk, PNB
- T₄ beer, bread
- T₅ beer, Milk.

Data Parallelism:-



Most parallel (or) distributed association rule algorithms try to parallelise either the data known as data parallelism (or) the candidates refer to as process parallelism.

The data parallelism is also called as count distribution algorithm.

The database is divided into 3 partitions and allocated to one for each processor.

Each processor counts the candidates for its data and then broadcast its counts through all other processors.

Each processor then determines the global count. These are used to determine the large itemsets.

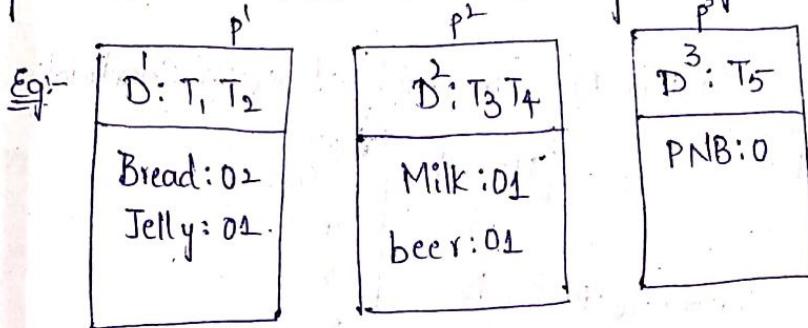
From the above diagram, There are 3 processors.

First 2 transactions are counted as P^1 and next 2 transactions at P^2 and next 1 transaction at P^3 .

When the local counts are obtained, they are broadcast to other processors, so that global counts can be generated.

Task parallelism:-

In the task parallelism, the candidates are partitioned and counted separately at each processor.



From the above diagram, There are 3 processors.

P^1 is counting bread and Jelly, P^2 is counting milk and beer

P^3 is counting PNB.

The first two Transactions are counted at P^1 and

next 2 at P^2 and next 1 at P^3 .

When the local counts are obtained, they are broadcast to other processors, so that global counts can be generated.

*Apply the Apriori Algorithm for the given set of transactions. And find what are the itemsets, 3 item sets when the minimum support is 60% and minimum confidence is 80%.

T_1 F, A, D, B
 T_2 D, A, C, E, B
 T_3 C, A, B, E
 T_4 B, A, D

Note:

Frequent items set - minimum support count

Association Rules - minimum support count
minimum confidence.

Sol:- Min Support count = 60%

$$= \frac{60}{100} \times 4 = 0.6 \times 4 = 2.4$$

Min. supcount = 2

Step 1:- Scan the dataset 'D' for count of each candidate

Itemset	Support count
A	4
B	4
C	2
D	3
E	2
F	1

Itemset	Support count
A	4
B	4
C	2
D	3
E	2

Compare the minsupcount with Support count

Step 2:- Generate C_2 from L_1

C_2	C_2	L_2
A, B	A, B	A, B
A, C	A, C	A, C
A, D	A, D	A, D
A, E	A, E	A, E
B, C	B, C	B, C
B, D	B, D	B, D
B, E	B, E	B, E
C, D	C, D	C, D
C, E	C, E	C, E
D, E	D, E	D, E

Compare

Itemset	Support count
A, B	4
A, C	2
A, D	3
A, E	2
B, C	2
B, D	3
B, E	2
C, D	1
C, E	2
D, E	1

Step3:- Generate C_3 from L_2

C_3
A, B, C
A, B, D
A, B, E
A, C, E
A, D, E
B, C, E
B, C, D
B, C, A
B, D, E

Itemset	Support Count
A, B, C	2
A, B, D	3
A, B, E	2
A, C, E	2
A, D, E	2
B, C, E	2
B, C, D	1
B, C, A	1
B, D, E	1

Itemset	Support Count
A, B, C	2
A, B, D	3
A, B, E	2
A, C, E	2
B, C, E	2

Step4:- Generate C_4 from L_3

C_4
A, B, C, D
A, B, C, E
A, B, D, E

Itemset	Support Count
A, B, C, D	1
A, B, C, E	2
A, B, D, E	1

Itemset	Support Count
A, B, C, E	2

9/03/19

Incremental Rules:-

The algorithms like Apriori, partition, data parallelism,

task parallelism, FPtree works for the static databases.

However in reality the dataset will be a dynamic in

the markets.

To address this issue, several approaches has been proposed are:-

1. Incremental Approach.

→ The incremental Approach concentrate on determining the large itemset for D and db .

When 'D' is a Database and 'db' is the updates of the database.

One Incremental Approach is fast update.

This fast update is based on the apriori algorithm. In this each iteration 'k' scans both D and db with the candidates generated from the

iteration $k-1$ based on large itemset at
that scan.

Advance Association Rule Technique:-

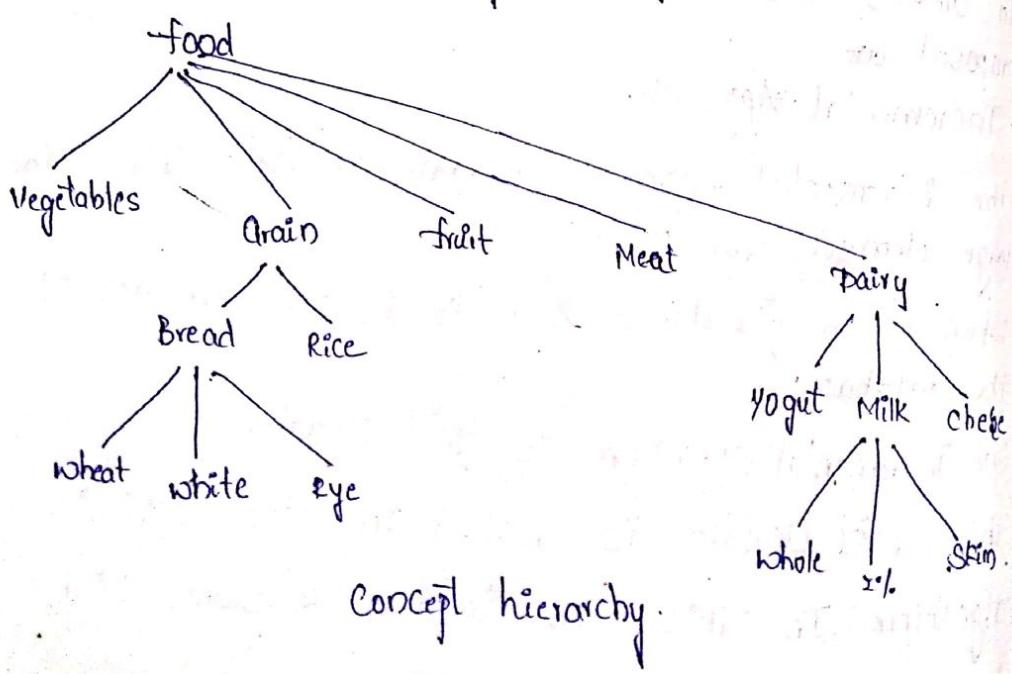
Techniques for Advanced Association Rules are
as follows:

1. Generalized Association Rule.
2. Multiple level Association Rules.
3. Quantitative Association Rules.
4. Using Multiple minimum Supports.
5. Correlation Rules.

Generalized Association Rule:-

→ A Generalized Association Rule $x \rightarrow y$ is defined
like a regular association rule with the restriction
that no item in 'y' may be above any item in 'x'.
Several algorithm have been proposed to generate
the association rules are:

→ The Simplest is Expand each transaction by adding all
the items above it in any hierarchy.



Multiple level Association Rules:-

→ Variation of Generalized Rules are Multiple level Association Rules.

In this itemset may occur from any level in the hierarchy.

Using a variation of Aprior algorithm, the concept hierarchy is traversed in top down and large itemset are formed.

The minimum support count for all the nodes in the hierarchy at the same level is identical.

If α_{i-1} is the minimum level of the level i in the hierarchy then $\alpha_{i-1} > \alpha_i$ is the level of the level $i-1$ in the hierarchy.

$$\alpha_{i-1} > \alpha_i$$

Quantitative Association Rules:-

A Quantitative Association rule is one that is Categorical and quantitative data.

Example:- A customer buys wine for between \$30 and \$50 of a bottle also buys a caviar.

Multiple minimum supports:-

Using one minimum support value for a large database can be a problem.

This problem is also called as rare item problem.

If the minimum support is too high then rules involving items that appear rarely will not be generated.

If it is too low → too many rules may be generated which are not important.

To handle this issue, one approach is partitioning the data based on the Support & generate association rules for each partition Separately.

Correlation Rule:-

A Correlation Rule is defined as a set of items which are correlated.

This can be useful for negative correlations.

formula:- Correlation of ($A \rightarrow B$) = $\frac{P(A, B)}{P(A) P(B)}$ if the

Correlation value < 1 then it indicates a negative correlation between $A \& B$.

Classification:-

Examples of classification include image and pattern recognition, medical diagnosis, loan approval, detecting faults in industry application.

Classification can be done by using different algorithms:

1. Statistical based Algorithms:-

(i) Regression , (ii) Bayes

2. Distance based Algorithms:-

K-nearest neighbour.

3. Decision Tree based Algorithms:-

(ID3, C4.5, CART)

4. Neural network based Algorithms:-

5. Rule based Algorithms:-

6. Combining Techniques.

Statistical based Algorithms:-

Baye's Algorithm:-

Consider the given weather dataset with attributes outlook, temperature, windy, humidity and classlabelplay, with values yes (or) no.

Find person that he can play (or) not with weather conditions.

outlook=Sunny , temp=cool , humidity=high , windy=strong

The dataset for weather is:

<u>Outlook</u>	<u>Temperature</u>	<u>Humidity</u>	<u>Windy</u>	<u>Play</u>
1. Sunny	hot	high	weak	No
2. Sunny	hot	high	strong	No
3. Overcast	hot	high	weak	Yes
4. Rainy	mild	high	weak	Yes
5. Rainy	cool	normal	weak	No
6. Rainy	cool	normal	strong	No
7. Overcast	cool	normal	strong	Yes
8. Sunny	mild	high	weak	No
9. Sunny	cool	normal	weak	Yes
10. Rainy	mild	normal	weak	Yes
11. Sunny	mild	normal	strong	Yes
12. Overcast	mild	high	strong	Yes
13. Overcast	hot	normal	weak	Yes
14. Rainy	mild	high	weak	No
			strong	

$$\text{probability of 'yes'} = 9/14 = P(\text{Yes})$$

$$\text{probability of 'no'} = 5/14 = P(\text{No})$$

<u>Outlook.</u>	<u>yes</u>	<u>No</u>	<u>P(Yes)</u>	<u>P(No)</u>
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0
Rainy	3	2	3/9	2/5

<u>Temperature</u>	<u>yes</u>	<u>No</u>	<u>P(Yes)</u>	<u>P(No)</u>
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cool	3	1	3/9	1/5

<u>Humidity</u>	<u>Yes</u>	<u>No</u>	<u>P(Yes)</u>	<u>P(No)</u>
High	3	4	3/9	4/5
Nominal	6	1	6/9	1/5
<u>Windy</u>	<u>Yes</u>	<u>No</u>	<u>P(Yes)</u>	<u>P(No)</u>
weak.	6	2	6/9	2/5
Strong	3	3	3/9	3/5

probability that the person can play the game.

$$P(X|play = yes) \cdot P(play = yes)$$

$$= \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}$$

$$= 0.0053$$

probability that the person can't play the game.

$$P(X|play = No) \cdot P(play = No)$$

$$= \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14}$$

$$= 0.0206$$

$$P(C) = 0.0053 + 0.0206 = 0.0259$$

$$P(A|C) = \frac{P(C/A) \cdot P(A)}{P(C)}$$

$$P(play = yes/x) = 0.0053 / 0.0259 = 0.206$$

$$P(play = No/x) = 0.0206 / 0.0259 = 0.795$$

The person.

$\therefore P(No)$ is having the maximum values when it is compared with $P(Yes)$.

Hence the person cannot go outside to play based on the given sample i.e., weather condition outlook = sunny, Temperature = cool, Windy = strong and humidity = high.

Baye's Classification:-

13/05/14

For the given dataset patient income with the attributes cold, running nose, headache, fever and the class label flu. Find a patient having flu or not with the given symptoms cold = yes, runny nose = no, headache = mild, fever = yes.

<u>Cold</u>	<u>runny nose</u>	<u>headache</u>	<u>fever</u>	<u>flu</u>
1. Yes	No	Mild	Yes	No
2. Yes	Yes	No	No	Yes
3. Yes	No	Strong	Yes	Yes
4. No	Yes	Mild	No	No
5. No	No	No	Yes	Yes
6. No	Yes	Strong	No	No
7. No	Yes	Strong	Yes	Yes
8. Yes	Yes	Mild	Yes	No

Sol:-

$$\text{probability of yes} = 5/8$$

$$\text{probability of no} = 3/8$$

<u>Cold</u>	<u>yes</u>	<u>No</u>	<u>P(yes)</u>	<u>P(No)</u>
Yes	3	1	3/5	1/3
No	2	2	2/5	2/3

<u>runny nose</u>	<u>yes</u>	<u>No</u>	<u>P(yes)</u>	<u>P(No)</u>
Yes	4	1	4/5	1/3
No	1	2	1/5	2/3

<u>headache</u>	<u>yes</u>	<u>No</u>	<u>P(yes)</u>	<u>P(No)</u>
Mild	2	1	2/5	1/3
No	1	1	1/5	4/3
Strong	2	1	2/5	1/3

<u>fever</u>	<u>yes</u>	<u>No</u>	<u>P(yes)</u>	<u>P(No)</u>
Yes	4	1	4/5	1/3
No	1	2	1/5	2/3

Probability that a patient can have the flu

$$P(x|\text{play} = \text{yes}) \cdot P(\text{play} = \text{yes}) = \frac{5}{8} \times \frac{1}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{3}{5}$$

$$= 0.024$$

Probability that a patient cannot have the flu

$$P(x|\text{flu} = \text{no}) \cdot P(\text{flu} = \text{no}) = \frac{3}{8} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}$$

$$= 0.0009$$

$$P(c) = 0.024 + 0.0009 = 0.0249$$

$$P(A|c) = \frac{P(c|A) \cdot P(A)}{P(c)}$$

$$P(\text{flu} = \text{yes}|x) = 0.024 | 0.0249 = 0.963855$$

$$P(\text{flu} = \text{No}|x) = 0.0009 | 0.0249 = 0.036$$

\therefore Hence $P(\text{yes})$ has the maximum value when it is compared with $P(\text{No})$

The patient have flu with the given symptoms.

According to Baye's theorem.

But, from the database the person is having no flu.

Hence the baye's theorem classified as false positive.

Pros of Naïve Baye's:-

- It is easy and fast to predict a class of dataset.
- Naïve Baye's perform better prediction compare to other models.

It performs well in the case of Categorical input variables compared to the numerical variables.

Naïve Baye's is also known as bad estimator.

Applications:-

- | | |
|------------------------|-----------------------|
| 1. Medical Analysis | 4. Sentiment Analysis |
| 2. Text classification | |
| 3. Spam filtering | |

14/3/19

Distance based Algorithms:-
K-nearest Neighbours:-

- 1. K is an integer value.
- 2. Find the value of k (must be a positive integer)
- 3. The given data point that needs to be classify
The algorithm computes the distance between i
the given data points and all other datapoints
of given datasets.
- 4. Find the 'k' nearest neighbours based on the
distance (the least distance value is the first
neighbour for a given data point).
- 5. The given data points will be placed in the
given data for majority of the classes with
respect to neighbours.

The dataset contains two attributes, namely weight and height and one class variable.

Find a person 'x' belongs to which class.
based on the given values weight = 57kgs
and height = 170cm.

Weight	height	Class	distance
51	167	Underweight	6.7
62	182	Normal	13
69	176	Normal	13.4
64	173	Normal	7.6
65	172	Normal	8.2
56	174	Underweight	4.1
58	169	Normal	1.4 (N1)
57	173	Normal	3 (N3)
55	170	Normal	2 (N2)

$$d_1 = \sqrt{(57-51)^2 + (170-169)^2} = \sqrt{6^2 + 3^2} = \sqrt{36+9} = 6.7$$

$$d_2 = \sqrt{(57-62)^2 + (170-182)^2} = 13$$

$$d_3 = \sqrt{(57-69)^2 + (170-176)^2} = \sqrt{12^2 + 6^2} = 13.4$$

$$d_4 = \sqrt{(57-64)^2 + (170-173)^2} = \sqrt{7^2 + 3^2} = \sqrt{49+9} = \sqrt{58} = 7.6$$

$$d_5 = \sqrt{(57-65)^2 + (170-179)^2} = \sqrt{8^2 + 9^2} = \sqrt{64+81} = \sqrt{145} = 8.2$$

$$d_6 = \sqrt{(57-56)^2 + (170-174)^2} = 4.1$$

$$d_7 = \sqrt{(57-58)^2 + (170-169)^2} = 1.4$$

$$d_8 = \sqrt{(57-57)^2 + (170-173)^2} = 3$$

$$d_9 = \sqrt{(57-55)^2 + (170-170)^2} = 2$$

The given data point ' x ' belongs to normal class (Consider

$k=3$)

If $k=1$ (x belongs to normal class) then find the distance which is having between all the points w.r.t ' x ' and select a datapoint which is having a least distance value and place the given datapoint in the obtained class.

Remark for k :

1. k must be a +ve integer.

2. k must not be a multiple of given number of class multiples.

* Find the person ' x ' from following with the given values age = 20 and gender = male. The given dataset is:

name	age	gender	Class following	distance
A	32	1	x_1	12
B	40	1	x_1	20
C	16	0	x_2	4.123 (N_1)
D	14	0	x_1	6.082
E	55	1	x_2	35
F	40	1	x_1	20
G	20	0	x_2	1 (N_1)
H	15	1	x_2	5
I	55	0	x_1	35.01
J	15	1	x_2	5

Sol:- Assume male = 1 and female = 0

$$A = \sqrt{(20-32)^2 + (1-1)^2} = \sqrt{12^2 + 0} = \sqrt{144} = 12$$

$$B = \sqrt{(20-40)^2 + (1-1)^2} = \sqrt{20^2 + 0} = 20$$

$$C = \sqrt{(20-16)^2 + (1-0)^2} = \sqrt{4^2 + 1^2} = \sqrt{16+1} = \sqrt{17} = 4.123$$

$$D = \sqrt{(20-14)^2 + (1-0)^2} = \sqrt{6^2 + 1^2} = \sqrt{36+1} = \sqrt{37} = 6.082$$

$$E = \sqrt{(20-55)^2 + (1-1)^2} = \sqrt{35^2 + 0} = \sqrt{35^2} = 35$$

$$F = \sqrt{(20-40)^2 + (1-1)^2} = \sqrt{20^2 + 0} = 20$$

$$G = \sqrt{(20-20)^2 + (1-0)^2} = \sqrt{0+1^2} = 1$$

$$H = \sqrt{(20-15)^2 + (1-1)^2} = \sqrt{5^2 + 0} = 5$$

$$I = \sqrt{(20-55)^2 + (1-0)^2} = \sqrt{35^2 + 1^2} = \sqrt{35^2 + 1} = \sqrt{35^2 + 1} = 35.01$$

$$J = \sqrt{(20-15)^2 + (1-1)^2} = \sqrt{5^2 + 0} = \sqrt{5^2} = 5$$

∴ The persons belong to x_2 .

Decision Tree:-

Each internal node is labelled with attribute a_i .
 Each arc is labelled with a predicate that can be applied to the attribute associated with the parent.

Each leaf node is labelled with a class c_j .

Construct a decision tree for the weather dataset.
 → Select the root node for which attribute is having the highest gain value.

Gain of attribute (A) = Information Gain - Entropy(A)

$$\text{Information gain}(P, N) = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{Entropy}(A) = \sum_{i=1}^r \frac{P_i + n_i}{P+N} I_G(A).$$

Let $P = Yes = 9$

$N = No = 5$

$$I_G(9, 5) = \frac{-9}{9+5} \log_2 \left(\frac{9}{9+5} \right) - \frac{5}{9+5} \log_2 \left(\frac{5}{9+5} \right)$$

$$= 0.41 + 0.54$$

$$= \underline{\underline{0.95}}$$

Entropy for Outlook:-

	P_i	n_i	$I_G(P_i, n_i)$
Sunny	2	3	0.970
Overcast	4	0	0
Rainy	3	2	0.940

$$I_G(2, 3) = \frac{-2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right)$$

$$= 0.52 + 0.44 = 0.96$$

$$\begin{aligned}
 IG(4,0) &= -\frac{1}{4} \log_2(1/4) - \frac{1}{4} \log_2(4/4) \\
 &= -1 \log_2(1) - 1 \log_2(1) \\
 &= (-1)0 - 0 = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Entropy(outlook)} &= \frac{2+3}{14} (0.94) + 0 + \frac{5}{14} (0.94) \\
 &= \underline{0.69}
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain(outlook)} &= IG - E(\text{outlook}) \\
 &= 0.94 - 0.69 = \underline{0.25}
 \end{aligned}$$

Entropy for temperature:-

	P_i	n_i	$IG(P_i, n_i)$
Hot	2	2	1
Mild	4	2	0.91
Cool	3	1	0.81

$$\begin{aligned}
 IG(2,2) &= -\frac{1}{2} \log_2(1/2) - \frac{1}{2} \log_2(1/2) = 1 \\
 &= -1/2 \log_2(1/2) - 1/2 \log_2(1/2) = 1
 \end{aligned}$$

$$\begin{aligned}
 IG(4,2) &= -\frac{1}{6} \log_2(4/6) - \frac{2}{6} \log_2(2/6) \\
 &= 0.38 + 0.52 = \cancel{0.88} = 0.91
 \end{aligned}$$

$$IG(3,1) = -\frac{3}{4} \log_2(3/4) - \frac{1}{4} \log_2(1/4) = 0.81$$

$$\begin{aligned}
 \text{Entropy(temperature)} &= \frac{2+2}{14} (1) + \frac{4+2}{14} (0.91) + \frac{3+1}{14} (0.81) \\
 &\approx 0.90
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain(temperature)} &= 0.94 - 0.90 \\
 &= 0.04
 \end{aligned}$$

Entropy for humidity:-

	P_i	n_i	$IG(P_i, n_i)$
High	3	4	0.98
Normal	6	1	0.59

$$IG(3,4) = \frac{-3}{7} \log_2(3/7) - \frac{4}{7} \log_2(4/7) = 0.98$$

$$IG(6,1) = \frac{-6}{7} \log_2(6/7) - \frac{1}{7} \log_2(1/7) = 0.59.$$

$$\text{Entropy (humidity)} = 0.78$$

$$\text{Gain (Humidity)} = 0.94 - 0.78 = 0.16.$$

Entropy for Windy:-

	P_i	n_i	$IG(P_i, n_i)$
weak	6	2	0.311
Strong	3	3	1

$$IG(6,2) = \frac{-6}{8} \log_2(6/8) - \frac{2}{8} \log_2(2/8) \\ = 0.311$$

$$\text{Entropy (windy)} = \frac{6+2}{14}(0.31) + \frac{3+3}{14}(1) = \frac{4}{7}(0.31) + \frac{3}{7}(1) =$$

$$\text{Gain (windy)} = 0.048.$$

Decision Tree Construction:-

Select the rootnode for which is having the highest Gain.

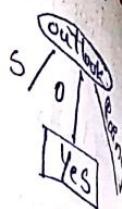
$$\text{Gain (outlook)} = 0.25 (\checkmark) \text{ (Rootnode)}$$

$$\text{Gain (Temperature)} = 0.04$$

$$\text{Gain (Humidity)} = 0.16$$

$$\text{Gain (windy)} = 0.048$$

<u>Outlook</u>	<u>temperature</u>	<u>humidity</u>	<u>windy</u>	<u>class</u>
Sunny	hot	high	weak	No
Sunny	hot	high	Strong	No
Sunny	mild	high	weak	No
Sunny	cool	high	weak	Yes
Sunny	mild	nominal	Strong	Yes
Sunny	mild	nominal	Strong	Yes



$$P = \text{Yes} = 2$$

$$N = N_{\text{No}} = 3$$

$$IG(3,2) = 0.970.$$

Entropy for temperature:-

	P_i	n_i	$IG(P_i, n_i)$
hot	0	2	0
mild	1	1	1
cool	1	0	0

$$\text{Entropy (temperature)} = \frac{1+1}{2+3}(1) = \frac{2}{5} = 0.4$$

$$\text{Gain (temperature)} = IG - \text{Entropy (temperature)}$$

$$= 0.970 - 0.4 = \underline{\underline{0.57}}$$

Entropy for humidity:-

	$P_i * n_i$	$IG(P_i, n_i)$
high	0 * 3	0
nominal	2 * 0	0

$$\text{Entropy (humidity)} = 0$$

$$\text{Gain (humidity)} = IG - \text{Entropy (humidity)}$$

$$= 0.970 - 0 = \underline{\underline{0.970}}$$

Entropy for Windy:-

	P_i	n_i	$IG(P_i, n_i)$
weak	1	2	0.91
strong	1	1	1

$$IG_{\text{windy}}(1, 2) = \frac{1}{3}(\log_2(1/3)) - \frac{2}{3}\log_2(2/3) = 0.91$$

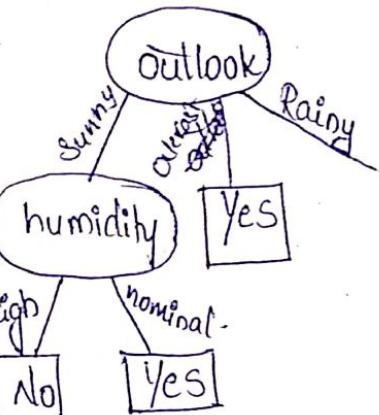
$$\text{Gain}_{\text{windy}} = IG_{\text{windy}} - \text{Entropy}_{\text{windy}} = \frac{3}{5}(0.91) + \frac{2}{5}(1) = 0.94$$

$$\text{Gain}_{\text{windy}} = IG_{\text{windy}} - \text{Entropy}_{\text{windy}}$$

$$= 0.970 - 0.94 = 0.03$$

Data set for Rainy:-

Outlook	temperature	humidity	windy	class
rainy	mild	high	weak	yes
rainy	cool	nominal	weak	yes
rainy	cool	nominal	strong	no
rainy	mild	nominal	weak	yes
rainy	wild	high	strong	no



$$P=3, N=2 \quad IG(3, 2) = 0.970.$$

Entropy for temperature:-

	P_i	n_i	$IG(P_i, n_i)$
wild	2	1	0.91
cool	2	1	1

$$\text{Entropy}(\text{temperature}) = \frac{3}{5}(0.91) + \frac{2}{5}(1) = 0.97$$

$$\text{Gain}(\text{temperature}) = 0.970 - 0.94 = 0.03$$

Entropy for humidity:-

	P_i	n_i	$IG(P_i, n_i)$
high	1	1	1
nominal	2	1	0.91

$$\text{E(humidity)} = 0.94$$

$$\text{G(h)} = 0.97 - 0.94 \\ = 0.03$$

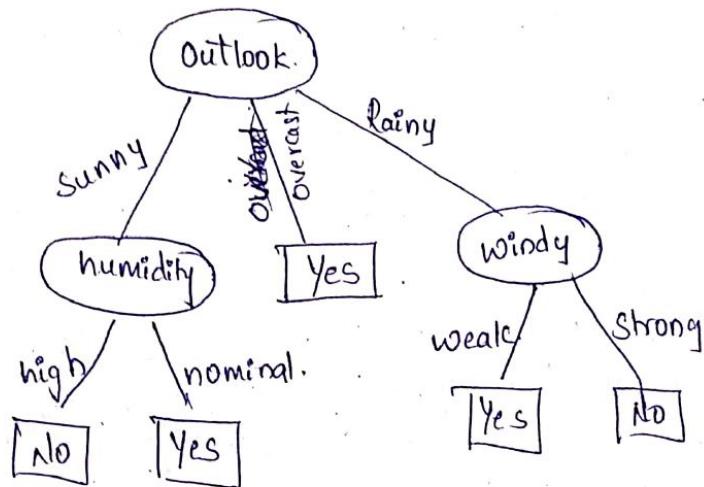
Entropy for windy:-

	P_i	n_i	$I_G(P_i, n_i)$
Weak	3	0	0
Strong	0	2	0

$$\text{Entropy(windy)} = 0$$

$$C_1(\text{windy}) = 0.97$$

$$= 0.97$$



1. Decision Tree for weather dataset

26/03/19

*Construct a decision tree with the attribute gender and height and class label

Gender	Height	Class
F	1.6m	Short
M	2m	Tall
F	1.9m	Medium
F	1.88m	Medium
F	1.7m	Short
M	1.85m	Medium
F	1.6m	Short
M	1.7m	Short
M	1.82m	Tall
M	2.1m	Tall
F	1.8m	Medium
M	1.7m	Medium
M	1.95m	Medium
F	1.9m	Medium
F	1.8m	Medium

Short = 4

Medium = 8

Tall = 3

$$IG(S, M, T) = \frac{S}{S+M+T} \log_2 \left(\frac{S}{S+M+T} \right) - \frac{M}{S+M+T} \log_2 \left(\frac{M}{S+M+T} \right) - \frac{T}{S+M+T} \log_2 \left(\frac{T}{S+M+T} \right)$$

$$IG(S, M, T) = \frac{4}{15} \log_2 \left(\frac{4}{15} \right) - \frac{8}{15} \log_2 \left(\frac{8}{15} \right) - \frac{3}{15} \log_2 \left(\frac{3}{15} \right)$$

$$= 1.456$$

Entropy for Gender:-

	S	M	T	$IG(S, M, T)$
Female	3	6	0	0.918
Male	1	2	3	1.459

$$IG(3, 6, 0) = \frac{-3}{9} \log_2 \left(\frac{3}{9} \right) - \frac{6}{9} \log_2 \left(\frac{6}{9} \right) - 0$$

$$= 0.918$$

$$IG(1, 2, 3) = \frac{-1}{6} \log_2 \left(\frac{1}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right)$$

$$= 1.459$$

$$\text{Entropy(Gender)} = \frac{\sum_{i=1}^2 \frac{S_i + M_i + T_i}{S+M+T} IG(M_i)}{S+M+T} + \frac{S+M+T}{S+M+T} IG(F)$$

$$= \frac{1+2+3}{4+8+3} (1.459) + \frac{3+6+0}{15} (0.918)$$

$$= \frac{6}{15} (1.459) + \frac{9}{15} (0.918) = 1.1344 (1.12)$$

$$Gain(\text{Gender}) = IG - \text{Entropy(Gender)}$$

$$= 1.456 - 1.12$$

$$= \underline{\underline{0.32}}$$

Entropy for Height:-

	S.	M	T	$I(G(S, M, T))$
≤ 1.7	3	0	0	0
$> 1.7 - < 2$	0	8	0	0
> 2	0	0	4	0

$$I_G(3,0,0) = \frac{-3}{3} \log_2 \left(\frac{3}{3}\right) = 0$$

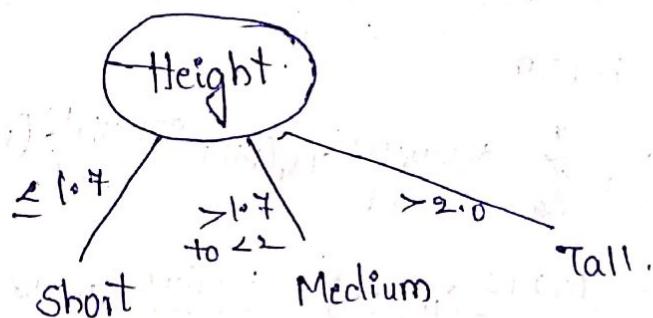
$$I_G(0,8,0) = \frac{-8}{8} \log_2 \left(\frac{8}{8}\right) = 0$$

$$I_G(0,0,4) = \frac{-4}{4} \log_2 \left(\frac{4}{4}\right) = 0$$

$$\begin{aligned} \text{entropy(height)} &= \frac{3}{15} I_G(\leq 1.7) + \frac{8}{15} I_G(> 1.7 - < 2) + \frac{2}{15} I_G(\\ &= \frac{3}{15}(0) + \frac{8}{15}(0) + \frac{2}{15}(0) = 0. \end{aligned}$$

$$\text{Gain(height)} = I_G - \text{entropy(height)}$$

$$= 1.456 - 0 = 1.456.$$



Proning:-

Once a tree is constructed some modifications to the tree might be required to improve the performance of the key tree during the classification.

The pruning phase removes the redundant Comparisons (or) removes subtrees to achieve better performance.

Neural Networks:-

→ Neural Network is a set of input and output nodes, and it is connected by giving some weights for the edges.

No. of source Nodes:-
Generally, the attribute names can be used as the source nodes (or) input nodes.

Determining which attributes to use as inputs is an issue-

No. of hidden layers:-

In the simplest case, there is only one hidden layer.

The layer which is in between source node (Input nodes) and Destination nodes (Output nodes) is called hidden layer.

No. of hidden nodes:-

The nodes available in the hidden layer can be called as hidden nodes.

Choosing the best no. of hidden nodes for hidden layer is one of the most difficult problem in neural networks.

At most $(n+1)$ hidden nodes can be used, where 'n' is no. of input nodes.

No. of sinks (or) Destination nodes:-

Usually, the no. of output nodes is the same as the no. of classes.

Interconnections:-

Each node is connected to all the nodes in the next level.

Wait:-
The wait assigned to an edge indicates the relative wait between the two nodes.

Activation functions:-

Many different types of activation functions can be used (Sigmoid, relu).

28/03/19.

Running Technique:-

The Technique for adjusting the weights is called Running Technique.

The most commonly used approach is back propagation.

The running may stop when all the training tuples has propagated through the network. It may be based on time or may be based on error rate.

Advantages of Neural Networks:-

1. Neural Networks are more robust than the decision trees because of the weights.
2. NN improves the performance by learning.
3. NN performs/contains no error rate and have more accuracy.

Neural Networks are more robust than the decision tree (or) in noisy environment.

Disadvantages:-

1. NN are difficult to understand.
2. Generating rules from the NN is not a straight forward.

Input attribute values must be numeric.

- Neural Network may be quite expensive to use.

Perceptron:-

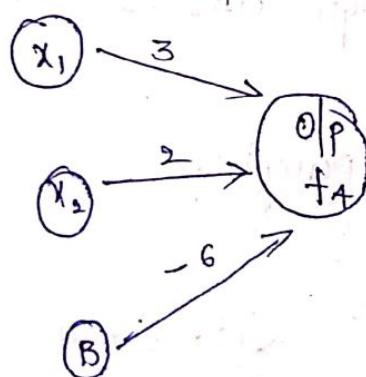
The Simplest neural network is called as a perceptron.

A perceptron is a single neuron with multiple inputs and output.

A simple perceptron can be used for binary classification.

Using a unipolar activation function a binary classification can be performed.

Example for a simple Perceptron:-



$$S = 3x_1 + 2x_2 - 6$$

$$S > 0$$

$$f_4 = \begin{cases} 1 & S > 0 \\ 0 & S \leq 0 \end{cases}$$

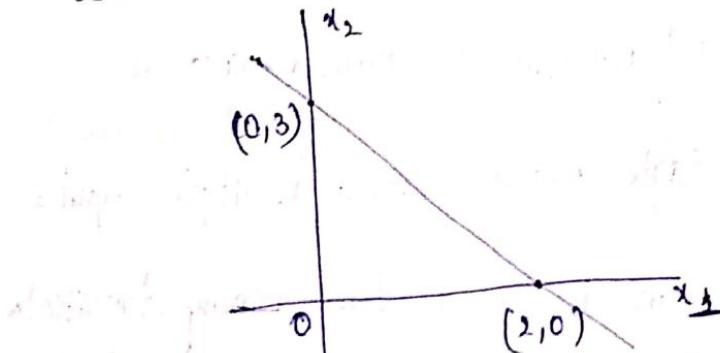
1 - Yes

0 - No.

<u>x_1</u>	<u>x_2</u>	<u>class</u>	<u>S</u>	<u>Predicted (f_4)</u>
0	1	Yes	-4	No (FN)
2	4	No	8	Yes (FP)
3	6	No	15	Yes (FP)
6	7	yes	26	Yes (TP)
2	0	No	0	Yes (FP)

$$\text{Accuracy} = \frac{1}{5} = 20\%$$

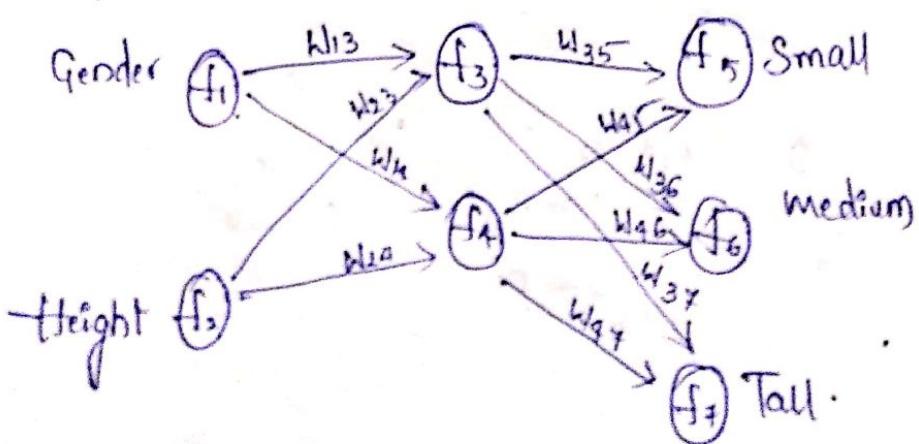
Alternative view:-



classification perception.

If only one layer exist (output) is called perception. If multiple layers exists is called multi-layer perception.

Example for Multi-layer perception.



Back Propagation:-

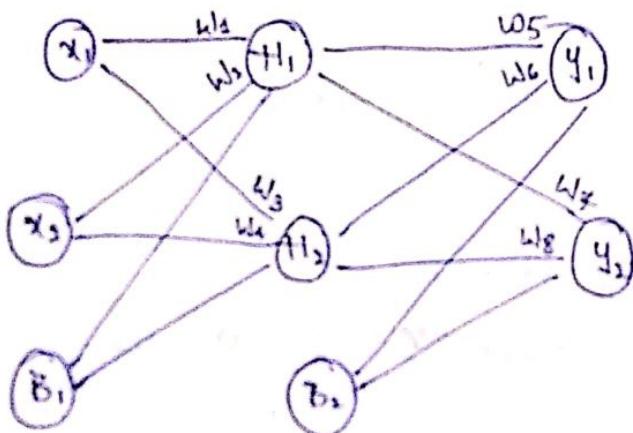
29/03/19.

Back Propagation learns by iterating processing a dataset of training tuples, comparing the network predictions for each tuple with the actual Node target value.

For each tuple, the weights are modified so as to minimize the mean squared error between the networks prediction and the actual target value.

These modifications are made in the backward direction (from the output layer) through each hidden layer down to the first hidden layer. Hence the name is Back Propagation.

* Find the output values for the output nodes and find the error between the predicted and the actual for a given network.



Target

$$x_1 = 0.05$$

$$T_1 \quad T_2$$

$$x_2 = 0.10$$

$$0.01 \quad 0.99$$

$$B_1 = 0.35$$

$$W_1 = 0.15 \quad W_5 = 0.40$$

$$B_2 = 0.60$$

$$W_2 = 0.20 \quad W_6 = 0.45$$

$$W_3 = 0.25 \quad W_7 = 0.50$$

$$W_4 = 0.30 \quad W_8 = 0.55$$

$$H_1 = x_1 w_1 + x_2 w_2 + B_1$$

$$= (0.05 \times 0.15) + (0.10 \times 0.20) + 0.35$$

$$H_1 = 0.375.$$

Activation function used is sigmoidal function is

$$\sigma = \frac{1}{1+e^{-x}}.$$

$$\text{Output of } H_1 = \frac{1}{1+e^{-H_1}} = \frac{1}{1+e^{-0.375}} = 0.592.$$

$$H_2 = x_1 w_3 + x_2 w_4 + B_2$$

$$H_2 = 0.05 \times 0.25 + 0.10 \times 0.30 + 0.35 = 0.3925$$

$$\text{Output of } H_2 = \frac{1}{1+e^{-0.3925}} = 0.5967.$$

For y_1 :

$$y = H_1 w_5 + H_2 w_6 + B_2$$

$$y = 0.592 \times 0.40 + 0.5967 \times 0.45 + 0.60 = 0.925.$$

$$\text{Output of } y_1 = \frac{1}{1+e^{-y_1}} = \frac{1}{1+e^{-0.925}} = 0.71.$$

For y_2 :

$$y_2 = H_1 w_7 + H_2 w_8 + B_2$$

$$y_2 = 0.592 \times 0.50 + 0.5967 \times 0.55 + 0.60 = 1.224.$$

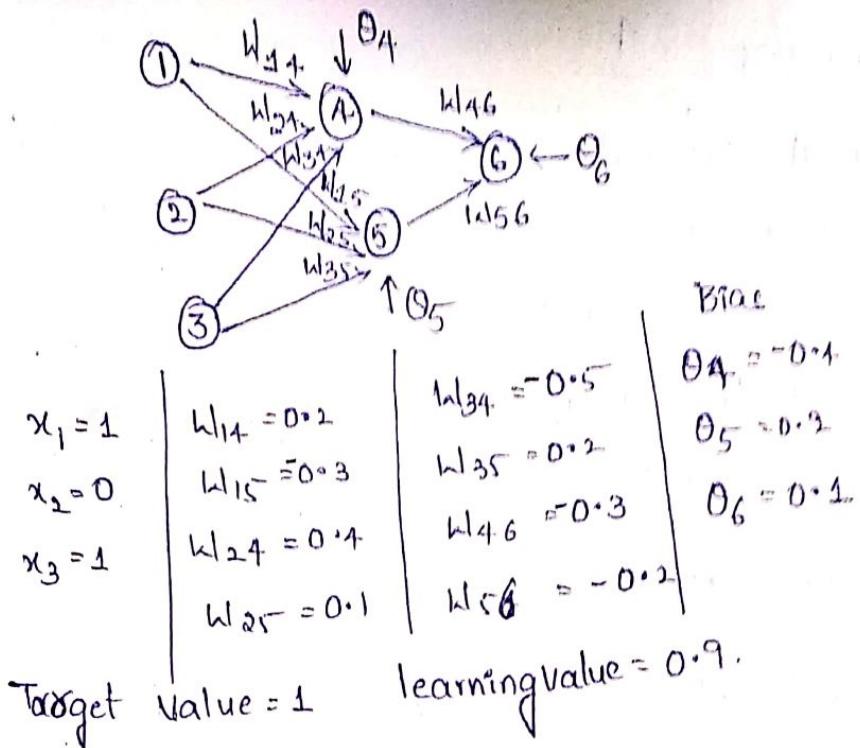
$$\text{Output of } y_2 = \frac{1}{1+e^{-y_2}} = \frac{1}{1+e^{-1.224}} = 0.772.$$

Error prediction is $\sum_{i=1}^m \frac{(\text{Target} - \text{predicted})}{m}$

$$\text{error} = \frac{1}{2} (T_1 - \text{output } y_1)^2 + \frac{1}{2} (T_2 - \text{output } y_2)^2$$

$$= \frac{1}{2} (0.01 - 0.73)^2 + \frac{1}{2} (0.99 - 0.775)^2 = 0.297$$

- For a given neural network. Find the output for
 - the hidden nodes & output nodes, find the errors
 of hidden nodes and output nodes and change
 the weights for the edges. For the bias values.



At node 4:-

$$\begin{aligned} \text{node } 4 &= x_1 w_{14} + x_2 w_{24} + x_3 w_{34} + \theta_4 \\ &= 1 \times 0.2 + (0)(0.4) + (1)(-0.5) - 0.4 \\ &= 0.2 + 0 - 0.5 - 0.4 \\ &= 0.2 - 0.9 = -0.7 \end{aligned}$$

$$\theta_4 = \frac{1}{1+e^{-(-0.7)}} = \frac{1}{1+e^{0.7}} = 0.33$$

At node 5:-

$$\begin{aligned} \text{node } 5 &= x_1 w_{15} + x_2 w_{25} + x_3 w_{35} + \theta_5 \\ &= -0.3 + 0 + 0.2 + 0.2 = 0.1 \end{aligned}$$

$$\theta_5 = \frac{1}{1+e^{0.1}} = 0.52$$

All nodes :-

$$\text{node } 6 = \frac{0.7(0.4)(0.4)(0.4)}{1 + 0.7(0.4)(0.4)(0.4)} = 0.476$$

$$O_6 = \frac{1}{1 + e^{-0.476}} = 0.476$$

The output of node 6 will equal to the target value. Hence there is a zero error for node 6 as well as bias units.

To calculate the error of the output node 6:

$$E_{116} = O_6(1-O_6)(T-O_6)$$

Error for hidden nodes :-

$$E_{115} = O_5(1-O_5)E_{116}w_{56}$$

Where w_{56} is the weight of the connection from the node 5 to node 6.

Error at output node 6:-

$$E_{116} = O_6(1-O_6)(T-O_6)$$

$$E_{116} = (0.476)(1-0.476)(1-0.476)$$

$$\boxed{E_{116} = 0.13\%}$$

Error at node 5:-

$$E_{115} = O_5(1-O_5)(1+0.5) E_{116} w_{56}$$

$$= (0.52)(1-0.52)(0.13)(-0.2)$$

$$\boxed{E_{115} = -0.006}$$

Error at node 4

$$E_{114} = O_4(1-O_4) E_{116} w_{46}$$
$$= (0.33)(1-0.33)(0.13)(-0.3)$$

$$E_{114} = -0.008$$

Formula for adjusting the weights is

$$w_{ij} = w_{ij} + (\lambda) E_{11j} \cdot O_i$$

Formula for adjusting the bias values is

$$\theta_j = \theta_j + (\lambda) E_{11j}$$

λ learning rate.

$$w_{14} = w_{14} + (\lambda) E_{114} \cdot O_i$$

$$w_{146} = w_{146} + (0.9) E_{116} \cdot O_4$$
$$= -0.3 + (0.9)(0.13)(0.55)$$

$$w_{146} = -0.26$$

$$w_{56} = w_{56} + (0.9) E_{116} \cdot O_5$$

$$= -0.2 + (0.9)(-0.008)(0.13)(0.52)$$
$$= -0.13$$

$$w_{14} = w_{14} + (\lambda) E_{114} \cdot O_1$$

$$= 0.2 + (0.9)(-0.008) \cdot O_1$$

$$= 0.19$$

$$w_{24} = -0.15$$

$$w_{15} = -0.306$$

$$w_{25} = 0.1$$

$$w_{34} = -0.508$$

$$w_{35} = 0.194$$

$$\theta_4 = \theta_4 + (\lambda) E_{114}$$

$$= (-0.4) + (0.9)(-0.008)$$

$$= -0.408$$

$$\begin{aligned}\theta_5 &= \theta_5 + (1) \text{err}_5 \\ &= 0.2 + (0.9)(-0.006) \\ &= 0.1946\end{aligned}$$

$$\begin{aligned}\theta_6 &= \theta_6 + (1) \text{err}_6 \\ &= 0.1 + (0.9)(0.13) \\ &= 0.217\end{aligned}$$

At Node 4:-

$$\begin{aligned}\text{node}_4 &= x_1 w_{14} + x_2 w_{24} + x_3 w_{34} + \theta_4 \\ &= 1 \times 0.19 + 0 + 1 \times (0.508) + (-0.408) \\ &= -0.726\end{aligned}$$

$$\theta_4 = \frac{1}{1 + e^{-(-0.726)}} = 0.32.$$

At Node 5:-

$$\begin{aligned}\text{node}_5 &= x_1 w_{15} + x_2 w_{25} + x_3 w_{35} + \theta_5 \\ &= 1 \times (-0.306) + 0 + 1 \times 0.194 + 0.2 \\ &= 0.079.\end{aligned}$$

$$\theta_5 = \frac{1}{1 + e^{0.079}} = 0.51.$$

Rule Based classification:-

Rule based classification is a set of "if-then" rules.

Syntax:- IF(condition) THEN Conclusion

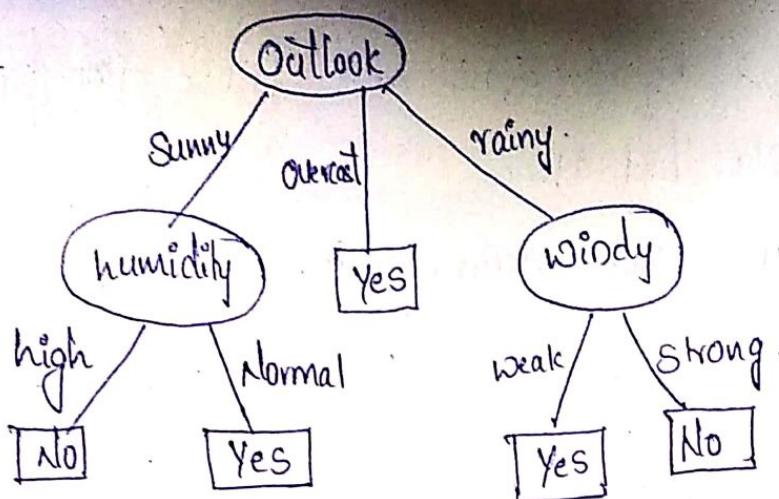
The 'IF' part is called as rule antecedent and 'THEN' part is called as rule Consequent.

Ex:- IF(outlook=sunny) AND(humidity=high) then play="no"

Extracting IF-THEN Rules from a DT:-

For a given DT the no. of if-then rules can be derived based on the no. of leaf nodes.

Ex:-



Rule 1:- IF (outlook = Sunny) and (humidity = high) THEN (play = No)

Rule 2:- IF (outlook = Sunny) and (humidity = Normal) THEN (play = Yes)

Rule 3:- IF (outlook = overcast) THEN (Play = Yes)

Rule 4:- IF (outlook = rainy) and (windy = weak) THEN (play = yes)

Rule 5:- IF (outlook = rainy) and (windy = strong) THEN (play = No)

IF-THEN rules can be assessed by 1. Coverage and Accuracy. (parameters of IF-THEN)

Let 'n_{covers}' be the no. of tuples covered by rule 'i' for a given data set

Let 'n_{correct}' be the no. of tuples correctly classified by the rule 'i'.

Let |D| be the size of the dataset.

$$1. \text{ Coverage} = \frac{n_{\text{covers}}}{|D|}$$

$$2. \text{ Accuracy} = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

From the above rules, rule 1 covers 3 tuples and those three tuples are correctly classified. Hence Coverage = $3/4$ and Accuracy = $3/3 = 100\%$. Rule 2 covers 2 tuples (9, 11) and those two tuples are correctly classified. Hence Coverage = $2/4$ and Accuracy = $2/2 = 100\%$.

Extracting IF-THEN rules from a neural network.

Rule Extraction Algorithms:

Input:-

D II. Training data

A II. Initial neural network.

Output:-

R II. Derived rules.

PX Algorithm:-

Cluster output node activation values;

Cluster hidden node activation values;

generate rules that describe the output values in terms of the hidden activation values;

generate rules that describe hidden output values in terms of inputs;

combine two set of rules;

Combining Techniques:-

4/03/19

Given a classification problem, one classification technique does not yield best results.

There have been some proposals for combining techniques.

There are two approaches for Combining Technique

1. Synthesis Approach
2. Multiple independent Approach.

Synthesis Approach:-

It takes multiple Techniques and combine them into new approach.

Clustering

Methods:-

1. Partitioning (K means)
2. Hierarchical (Agglomerative or divisive) approaches
3. Density based (DBSCAN)
4. Grid based.

K-means clustering:-

Let $k=2$

2, 4, 6, 7, 18, 20, 23, 30, 26, 35

Select 2 elements from the given

Let $m_1 = 7$

$m_2 = 30$

$$k_1 = \{2, 4, 6, 7, 18\}$$

$$k_2 = \{20, 23, 30, 26, 35\}$$

$$m_1 = 2+4+6+7+18 \\ = 7 \cdot 4 \approx 8$$

$$m_2 = 20+23+30+26+35$$

$$m_2 = 26 \cdot 8 \approx 27$$

2:

$$k_1 = \{2, 4, 6, 7\}$$

$$k_2 = \{18, 20, 23, 30, 26, 35\}$$

$$m_1 = 2+4+6+7$$

$$m_2 = 25 \cdot 3 (26)$$

$$m_1 = 4 \cdot 7 \approx 5$$

$$k_1 = \{2, 4, 6, 7\}$$

$$k_2 = \{18, 20, 23, 30, 26, 35\}$$

$$m_1 = 5$$

$$m_2 = 26$$

Here 2 previous clusters are same. So we can stop it.

K-mean Algorithm:-

I/P:-

D || Dataset

K || no. of clusters

O/P:- \rightarrow set of K clusters.

Method:- Arbitrarily choose k objects from D .

repeat

Assign each object to the cluster to which the object is most similar based on the mean value.

Update the cluster means.

until no change.

• k -Medoid: This is also known as PAM (Partitioning Around Medoids)

Ex:- 1, 2, 3, 8, 9, 10, 25. Perform the k -medoid algorithm for given dataset. The given pairs are:

(1, 2), (1, 3), (1, 8), (1, 9), (1, 10), (2, 3), (2, 8), (2, 9), (2, 10), (3, 8), (3, 9), (3, 10), (8, 9), (8, 10), (9, 10).

Sol:- let $k=2$

distance (3, 4) distance from (7, 4)

If $(x_1, y_1) \neq (x_2, y_2)$

Manhattan distance = $|x_1 - x_2| + |y_1 - y_2|$

$$|2-3| + |6-4| = 3$$

$$|3-3| + |8-4| = 4$$

$$|4-3| + |7-4| = 4$$

$$|6-3| + |2-4| = 5$$

$$|6-4| + |4-4| = 2$$

$$|7-3| + |3-4| = 5$$

$$|8-3| + |5-4| = 6$$

$$|7-3| + |6-4| = 6$$

$$|2-7| + |6-4| = 7$$

$$|3-7| + |8-4| = 8$$

$$|4-7| + |7-4| = 6$$

$$|6-7| + |2-4| = 3$$

$$|6-7| + |4-4| = 1$$

$$|7-7| + |3-4| = 1$$

$$|8-7| + |5-4| = 2$$

$$|7-7| + |6-4| = 2$$

Total cost = $3+4+4+3+1+1+2+2 = 20$

Cluster 1: $\{(3, 4), (2, 6), (3, 8), (4, 7)\}$

Cluster 2: $\{(1, 4), (6, 2), (6, 4), (7, 3), (8, 5), (7, 6)\}$

Select the next 2 points $(3, 4) \times (7, 3)$

distance $(3, 4)$

distance $(7, 3)$

$$(2, 6) |2-3| + |6-4| = 3 \quad |2-7| + |6-3| = 8$$

$(3, 4)$

$$(3, 8) |3-3| + |8-4| = 4 \quad |3-7| + |8-3| = 9$$

$$(4, 7) |4-3| + |7-4| = 4 \quad |4-7| + |7-3| = 7$$

$$(6, 2) |6-3| + |2-4| = 5 \quad |6-7| + |2-3| = 2$$

$$(6, 4) |6-4| + |4-4| = 2 \quad |6-7| + |4-3| = 2$$

$$(7, 3) |7-3| + |3-4| = 5 \quad |7-7| + |4-3| = 1$$

$$(8, 5) |8-3| + |5-4| = 6 \quad |8-7| + |5-3| = 3$$

$$(7, 6) |7-3| + |6-4| = 6 \quad |7-7| + |6-3| = 3$$

$$\text{Total cost} = 3+4+4+2+2+1+3+3 = 22.$$

$$22 > 20$$

$$(3, 4) (7, 3) > (3, 4) (7, 4)$$

\therefore Select the previous clusters as answers

\therefore Answer is

$$\text{Cluster for } (3, 4) = \{(3, 4), (2, 6), (3, 8), (7, 3)\}$$

$$\text{Cluster for } (7, 4) = \{(7, 4), (6, 2), (6, 4), (7, 3), (8, 5), (7, 6)\}$$

k-medoid is more robust than k-mean in the presence of noise and outlier because a medoid is less influenced by outliers (or) other extreme values than a mean.

The Complexity of each iteration in k-medoid algorithm is $O(k(n-k)^2)$

where n is no. of records (or) points (or) size of data points.

$k = \text{no. of k medoid} / \text{no. of clusters}$.

A typical k medoid (PAM) works well for small datasets but doesn't scale well for large datasets.

To deal with large dataset Sampling based methods called CLARA (clustering large applications) can be used.

CLARA:

Instead of taking whole dataset into consideration, CLARA uses a random sample of dataset. The PAM algorithm is then applied to compute best medoid for the sample.

→ CLARA builds clustering from the multiple random samples and returns the best clustering as the output.

The Complexity of Computing the medoid on a random sample is $O(ks^2 + k(n-k))$

where s = size of sample

k = no. of clusters

n = size of the total objects.

Effectiveness of CLARA depends upon the samplesize.

→ The Time Complexity of k-means to find the clusters is $O(n)$.

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

k-medoid Algorithm:-

I/P:-

K = no. of clusters

D = dataset

O/P:-

A set of K clusters.

Method:-

Orbitrarily choose k objects in D as the initial representative Objects.

Repeat

Assign each remaining object to the clusters with the nearest rep object.

Randomly select a non-representative object O_{random}

Compute the total cost, s , of swapping rep object O_j with O_{random} .

If $s < 0$ then swap O_j with O_{random} to form the new set of k rep objects.

Until no change.

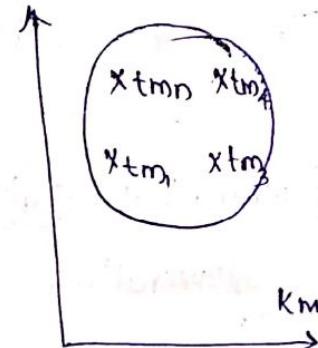
Similarity and distance measure:-

Given a cluster K_m of n points $\{t_{m1}, t_{m2}, \dots, t_{mn}\}$

$$\text{Centroid} = C_m = \frac{\sum_{i=1}^n t_{mi}}{N}$$

$$\text{radius} = R_m = \sqrt{\frac{\sum_{i=1}^n (t_{mi} - C_m)^2}{N}}$$

$$\text{Diameter } D_m = \sqrt{\frac{\sum_{j=1}^n \sum_{i=1}^j (t_{mi} - t_{mj})^2}{N(N-1)}}$$

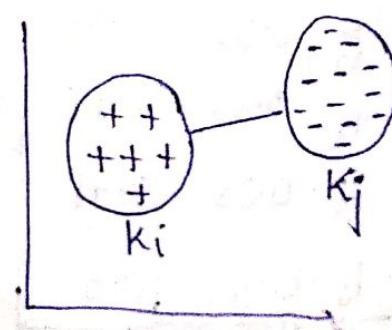


Centroid is the middle of the clusters.

Given clusters k_i & k_j there are several standard alternatives to calculate the distance between clusters.

→ The representative list is

1. Single link
2. Complete link
3. Average link



Single link:-

The smallest distance between an element in one cluster and an element in other cluster.

$$\text{Thus, the distance of } (k_i, k_j) = \min(\text{distance}(t_{i1}, t_{j1}), \dots)$$

$\forall t_{i1} \in k_i \notin k_j \&$

$\forall t_{j1} \in k_j \notin k_i$

Complete link:-

The largest distance between an element in one cluster and an element in other.

$$\text{Thus, the distance of } (k_i, k_j) = \max(\text{distance}(t_{i1}, t_{j1}), \dots)$$

Average link:-

An average distance between in one cluster an element in other cluster.

$$\text{Thus, the distance of } (k_i, k_j) = \text{mean}(\text{distance}(t_{i1}, t_{j1}), \dots)$$

16/04/19

Hierarchical clustering:-

1. Agglomerative.

2. Divisive.

Perform Agglomerative clustering using Single link.

	x	y
P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.3

Step 1:- Calculate the distance matrix using Euclidean distance.

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0					
P ₂	0.23	0				
P ₃	0.21	0.14	0			
P ₄	0.36	0.20	0.15	0		
P ₅	0.34	0.14	0.28	0.29	0	
P ₆	0.28	0.25	0.11	0.22	0.39	0

P₁ :-

$$(0.22, 0.38), (0.40, 0.53)$$

$$D = \sqrt{(0.40 - 0.22)^2 + (0.53 - 0.38)^2} = 0.23$$

:-

$$(0.35, 0.32), (0.40, 0.53)$$

$$D = \sqrt{(0.40 - 0.35)^2 + (0.53 - 0.32)^2} = 0.21$$

$$(0.35, 0.32), (0.22, 0.38)$$

$$D = \sqrt{(0.22 - 0.35)^2 + (0.38 - 0.32)^2} = 0.14$$

$$(0.26, 0.19), (0.40, 0.53)$$

$$D = \sqrt{(0.40 - 0.26)^2 + (0.53 - 0.19)^2} = 0.36$$

$$(0.26, 0.19), (0.22, 0.38)$$

$$D = \sqrt{(0.22 - 0.26)^2 + (0.38 - 0.19)^2} = 0.20$$

P₄P₁:

$$(0.28, 0.19), (0.35, 0.39)$$

$$ED = \sqrt{(0.35 - 0.28)^2 + (0.39 - 0.19)^2} = 0.15$$

P₅P₁:-(0.08, 0.41), (0.40, 0.53)

$$ED = \sqrt{(0.40 - 0.08)^2 + (0.53 - 0.41)^2} = 0.34$$

P₅P₂:-(0.08, 0.41), (0.22, 0.58)

$$ED = \sqrt{(0.22 - 0.08)^2 + (0.58 - 0.41)^2} = 0.14$$

P₅P₃:-(0.08, 0.41), (0.35, 0.32)

$$ED = \sqrt{(0.35 - 0.08)^2 + (0.32 - 0.41)^2} = 0.27$$

P₅P₄:-(0.08, 0.41), (0.26, 0.19)

$$ED = \sqrt{(0.26 - 0.08)^2 + (0.19 - 0.41)^2} = 0.29$$

P₆P₁:-(0.45, 0.3), (0.40, 0.53)

$$ED = \sqrt{(0.40 - 0.45)^2 + (0.53 - 0.3)^2} = 0.53$$

P₆P₂:-(0.45, 0.3), (0.22, 0.58)

$$ED = \sqrt{(0.22 - 0.45)^2 + (0.58 - 0.3)^2} = 0.25$$

P₆P₃:-(0.45, 0.3), (0.35, 0.32)

$$ED = \sqrt{(0.35 - 0.45)^2 + (0.32 - 0.3)^2} = 0.11$$

P₆P₄:-(0.45, 0.3), (0.26, 0.19)

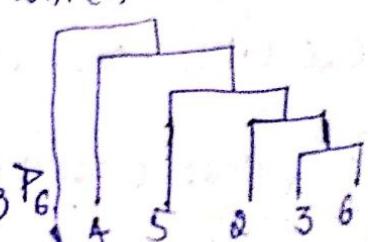
$$ED = \sqrt{(0.26 - 0.45)^2 + (0.19 - 0.3)^2} = 0.22$$

P₆P₅:-(0.45, 0.3), (0.08, 0.41)

$$ED = \sqrt{(0.08 - 0.45)^2 + (0.41 - 0.3)^2} = 0.39$$

Step:- Select the point in which is having
the minimum distance.

∴ Min distance point is P₃P₆

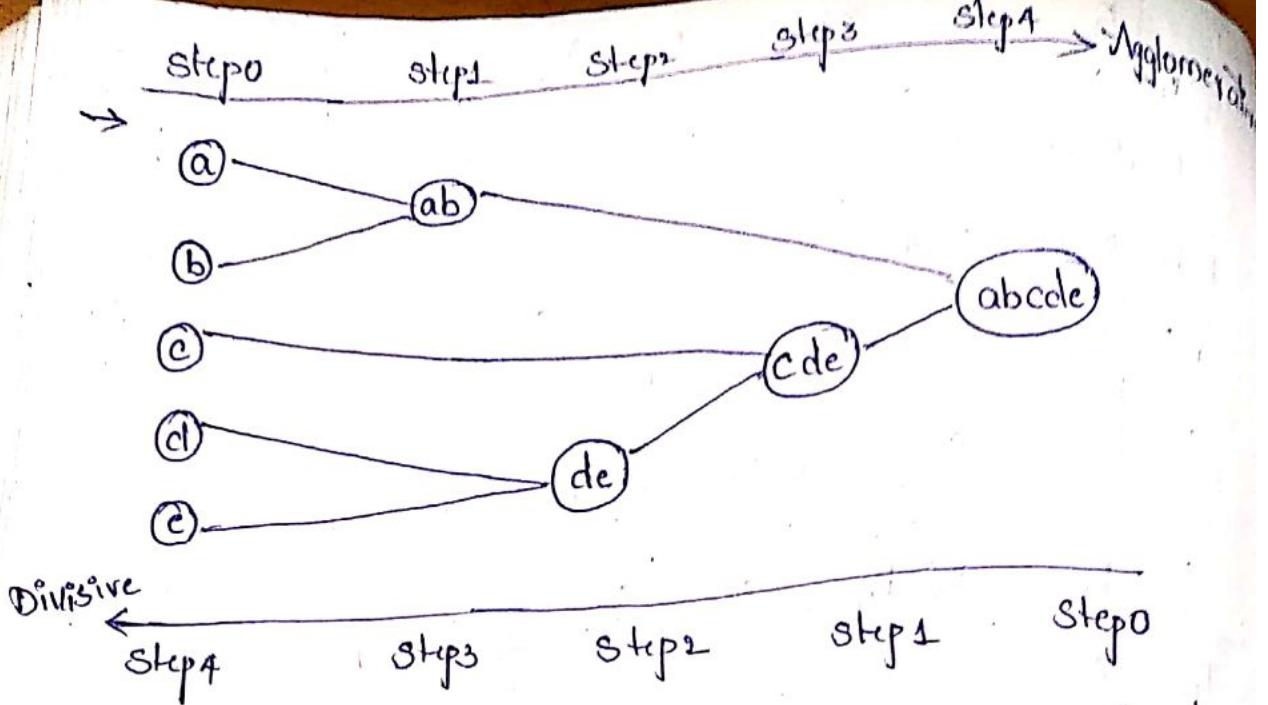


	P_1	P_2	$P_3 P_6$	P_4	P_5
P_1	0				
P_2	0.23	0			
$P_3 P_6$	0.21	0.14	0		
P_4	0.36	0.20	0.15	0	
P_5	0.34	0.14	0.28	0.29	0

	P_1	$P_3 P_6 P_4$	P_4	P_5
P_1	0			
$P_3 P_6 P_4$	0.21	0		
P_4	0.36	0.15	0	
P_5	0.34	0.14	0.29	0

	P_1	$P_3 P_6 P_2 P_5$	P_4
P_1	0		
$P_3 P_6 P_2 P_5$	0.21	0	
P_4	0.36	0.15	0

	P_1	$P_3 P_6 P_2 P_5 P_4$
P_1	0	
$P_3 P_6 P_2 P_5 P_4$	0.21	0



Clustering with large Databases:-

BIRCH → Balanced

20/04/19

clustering With large Databases:

- BIRCH → Balance Iterative Reducing & clustering with Hierarchies.
- ~~DBSCAN~~ Combination of hierarchical & Partitioning Algorithms.
- 2 Difficulties in Agglomerative Overcomes
 - 1.) Scalability
 - 2.) Undo Operation (for previous step).
- applicable for numeric data.

notations:

CF tree (clustering feature) \rightarrow contains tuples (n, LS, SS)

CF Tree (height balanced tree)
 $n \rightarrow$ no. of points.
 $LS \rightarrow$ Linear sum.

$$CF_1 = (n_1, LS_1, SS_1)$$

$SS \rightarrow$ Squared sum.

$$CF_2 = (n_2, LS_2, SS_2)$$

$$x = \{x_1, x_2, x_3, \dots, x_n\}$$

$$CF_1 + CF_2 =$$

$$LS = \sum_{i=1}^n x_i$$

$$(n_1+n_2, LS_1+LS_2, SS_1+SS_2). \quad SS = \sum_{i=1}^n x_i^2$$

When cluster C_1 contains the elements

Suppose there are 3 points

$$(2, 5) (3, 2) (4, 3) \text{ in cluster } C_1$$

$$(4, 6) (6, 7) (7, 8) \text{ in cluster } C_2 \dots \quad LS_1 = (2+3+4, 5+2+3)$$

$$CF_1 = (3, (9, 10), (29, 38))$$

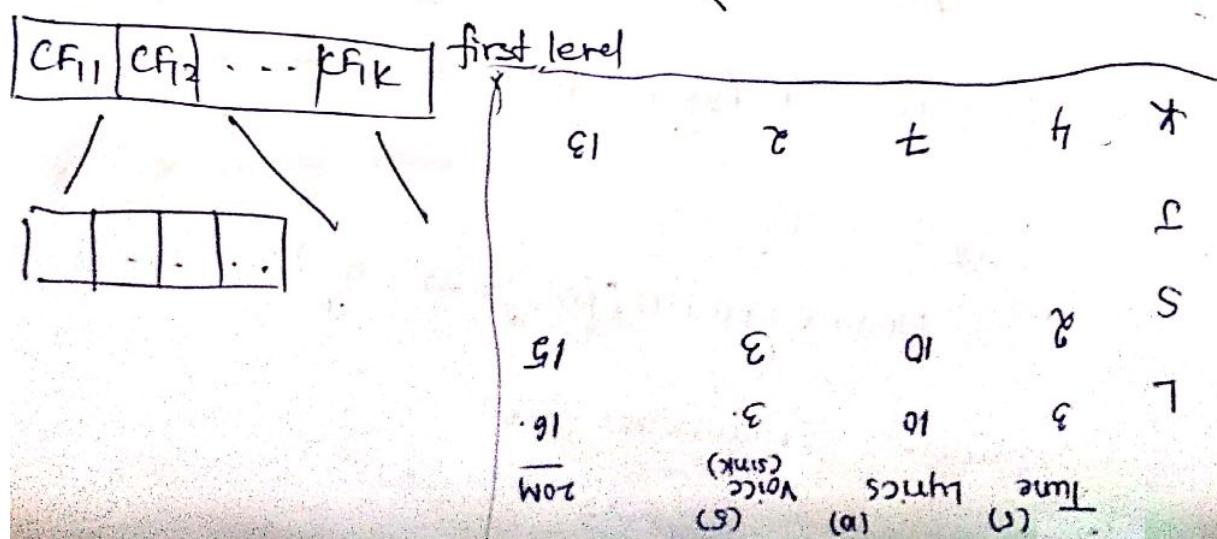
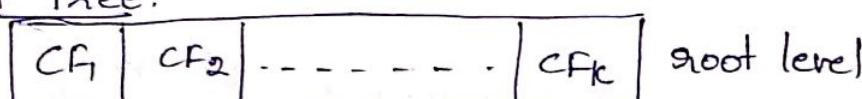
$$LS_1 = (9, 10) \\ 2+3+4+7, 5+2+3$$

$$CF_2 = (3, (17, 21), (101, 149))$$

$$= (29, 38)$$

$$CF_1 + CF_2 = (6, (26, 31), (130, 187))$$

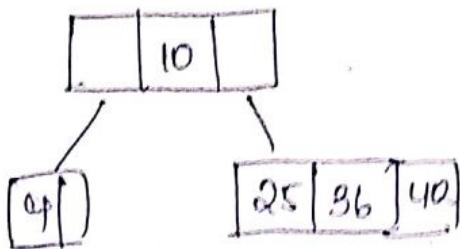
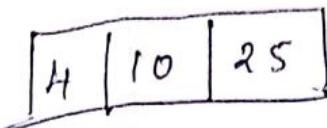
CF Tree:



It is a balanced Tree and It is also a B-Tree.

4, 10, 25, 36, 40

$$m=3$$

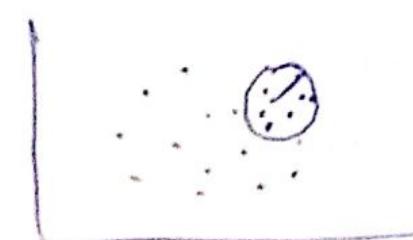


(2) DBSCAN: Density based Spatial Clustering Applications with NOISE.

- Used for Spherical shape of clusters for finding by Hierarchical & Partitioning algorithms.
- S-shaped & Oval shaped clusters by using DBSCAN.

3 Imp points

- Density
- Core point
- Boundary Point



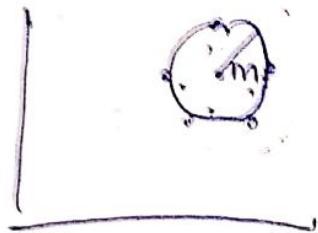
ϵ (Epsilon) $> 0 \Rightarrow$ +ve Integer $R = 2$ (Radius).
with Radius-2 we need to derive a circle by selecting a random point.

Density :- no. of points covered within a radius of a circle is called Density.

$$\therefore \text{Density} = 7.$$

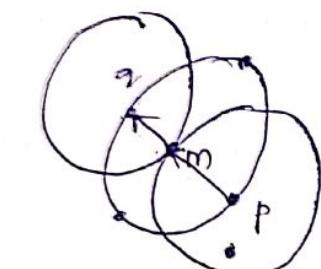
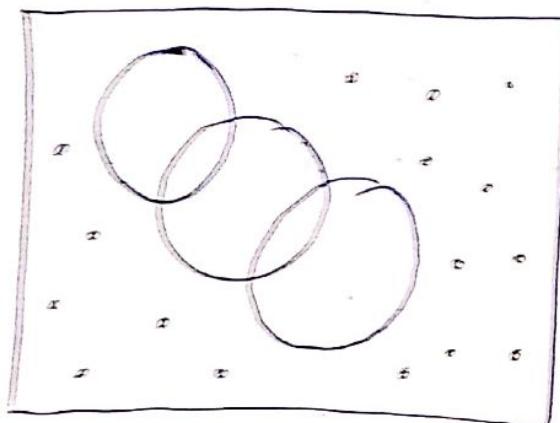
Core point:

Boundary point: Points that lie on the Boundary of a circle is called Boundary point.



Core points:

From the given point, if we can traverse any of the 3 points in that circle. Given min points = 3.



m → core point

If min point = 3; then the density should be = (6) \Rightarrow its density then the

p → core point

m → core point

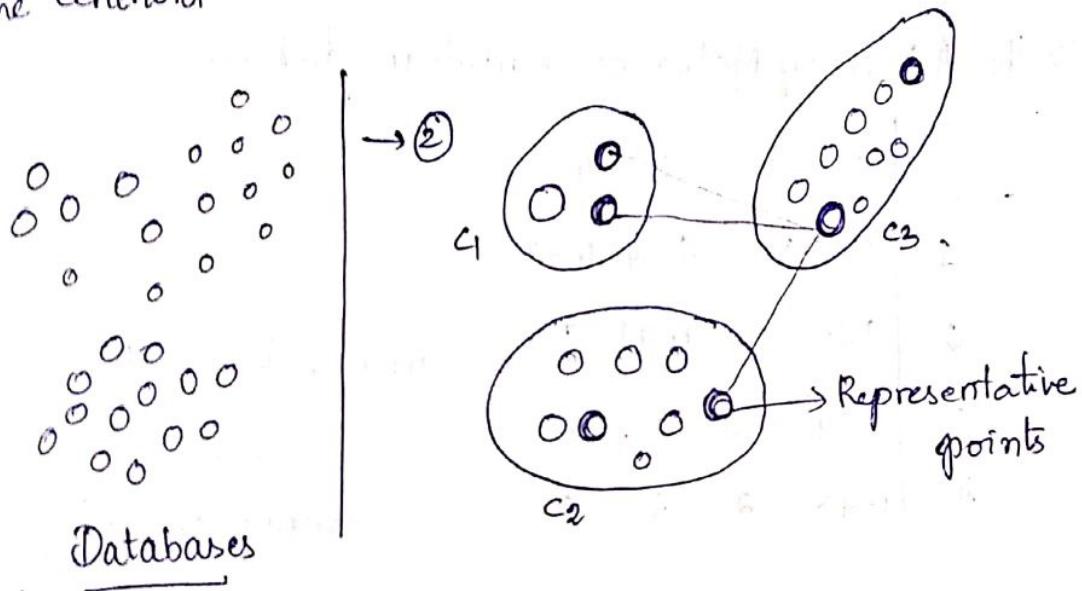
q → core point (iff and only if it traverses from p → q).

From the above diagram p & m are directly core point, q is also a core point bcoz, q can be reached from $P \rightarrow m$ & $m \rightarrow q$. Hence 'q' is a Core point.

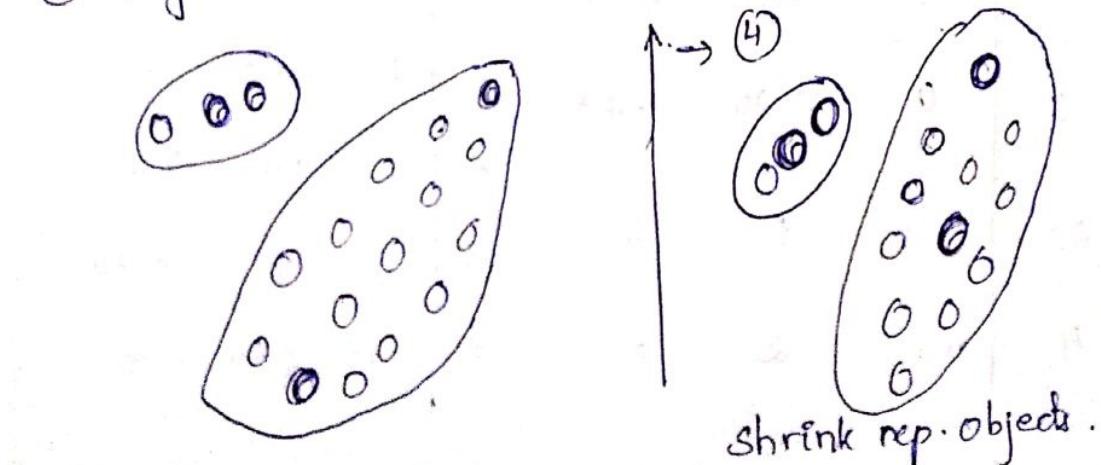
DBSCAN finds core objects i.e; the Objects that have dense neighbourhoods, it connects core objects & their neighbourhoods to form dense regions as clusters.

b) CURE :- Clustering Using Representative

- 1) For a given dataset find the no. of clusters that can be formed
- 2) Identify the Representative points
- 3) From the Representative points, find the Distance b/w the points from one cluster to another cluster.
- 4) Merge the cluster which is having the min. Distance b/w the Representative points.
- 5) After Merging the Representative points | shrinking towards the Centroid.



⑤ Merge c_2 & c_3 (min distance b/w 2 rep. point)



Clustering with Categorical Attributes:

Ex:

→ {book, water, sun, sand, swim, ready}

→ There are 4 documents

→ 1 - {books}

→ 2 - {water, sun, sand, swim}

→ 3 - {water, sun, swim, ready}

→ 4 - {read, sand}

boolean Matrix

1 (1, 0, 0, 0, 0, 0)

2 (0, 1, 1, 1, 1, 0)

3 (0, 1, 1, 0, 1, 1)

4 (0, 0, 0, 1, 0, 1)

Write the Adjacency Matrix as Euclidean distance

	1	2	3	4
1	0	0.24	0.24	1.73
2	0.24	0	1.41	2
3	0.24	1.41	0	2
4	1.73	2	2	0

min value = 1.41

= (2, 3).

→ average of 2 & 3 Matrix

k=2 → no. of clusters.

→ 2 & 3 (0, 1, 1, 0.5, 1, 0.5)

	1	23	4
1	0	0.24	1.73
23	0.24	0	2
4	1.73	2	0

	1	23	4
1	0	2.12	<u>1.73</u>
23	2.12	0	2.34
4	1.73	2.34	0

k-cluster = (2, 3)(1, 4)

2) ROCK Algorithm:

It can be used to find the clusters for both boolean data and Categorical data.

- ✓ In this Identifying, the similarity based on no. of links b/w the items.
- ✓ Similarity of items: A pair of items are said to be Neighbour if their similarity exceeds some threshold.
- ✓ The no. of links b/w the items is defined as no. of common neighbours that they have.
- ✓ Instead of using Euclidean distance, use the Jaccard Coefficient $JC = \text{simi}(t_i, t_j) = \frac{q}{q+r+s}$.

given threshold $\neq 0.6$.

$q \rightarrow 1$

$r \rightarrow 1$

$s \rightarrow 0$

$t \rightarrow 0$

	1	2	3	4
1	1	0	0	0
2	0	1	0.6	0.2
3	0	0.6	1	0.2
4	0	0.2	0.2	1

$$\text{Sim}(1,2) = 0/$$

$$\text{Simi}(1,3) = 0$$

$$\text{Simi}(1,4) = 0$$

$$\text{Simi}(2,3) = \frac{3}{3+1+1} = \frac{3}{5} = \frac{3}{5} = 0.6$$

The points that satisfy the threshold are $(2,3)(2,4)$ & $(3,4)$

The points that exceeds the threshold are $(1,1)(1,2)(3,1)(4,1)$

- Agglomerative Algorithms may also called as Agnes
Agglomerative Nesting. (AGNES)
- Divisive clustering is called as DIANA (Divisive Analysis)