

Data Reduction :-

It is a technique to reduce the size of data but provides the same analytical results if you work on original data is called as data reduction.

Advantages :-

- To improve the performance of data mining technique.
- To reduce the network bandwidth to transfer.
- Reduce the space.

Types of data Reduction :-

- * Dimensionality Reduction
- * Numerosity Reduction
- * Data compression

→ Dimensionality Reduction :

It is the process of reducing the number of random variables or attributes under consideration.

Methods in dimensional Reduction are:-

- ↳ wavelet transforms
- ↳ principal component Analysis (PCA) ↳ - Transf
original data to sample data
(smaller space)
- ↳ Attribute Subset selection → irrelevant,
weakly relevant- or redundant attribut
(or) dimensions are deleted (or) removed

→ Numerosity Reduction :-

It is the technique to replace the original data size by alternative smaller forms of data representation.

There are two types in numerosity reduction

- parametric Numerosity Reduction
- Non parametric Numerosity Reduction

Parametric: In this model, we used to estimate the data so that, the required parameters needs to be stored in

memory, value from actual data.

Techniques:- Regression

- log linear models

Non-parametric :- It is used to storing reduce representation of data

Techniques: Histogram

- clustering
- Sampling
- Data cube aggregation.

Data compression:

In this, transformations are applied to obtain reduced or compressed representation of data.

It is a reduction in the no. of bits need to be represent the data.

There are two types → loss less (If the original data can be reconstructed from the compressed data without any information loss)

→ lossy. (Instead of getting original data we can reconstruct only an approximation of the original data)

Techniques in loss less < string compression
All the dimensionality & numerosity reductions.

Wavelet Transformation:-

used for data Analysis to make more accurate data. The diff b/w Fourier transformation and wavelet transformation is, The fourier transformation will not produce accurate image/ data compare to the wavelet transformation

- It is not easy to work on continuous signal but in wavelet transformation signal.
- In wavelet transformation, both areas of +ve & -ve will have the common factors then we can apply the transformation only one side

- There are 2 types of wavelet transformation

- continuous wavelet

- discrete wavelet

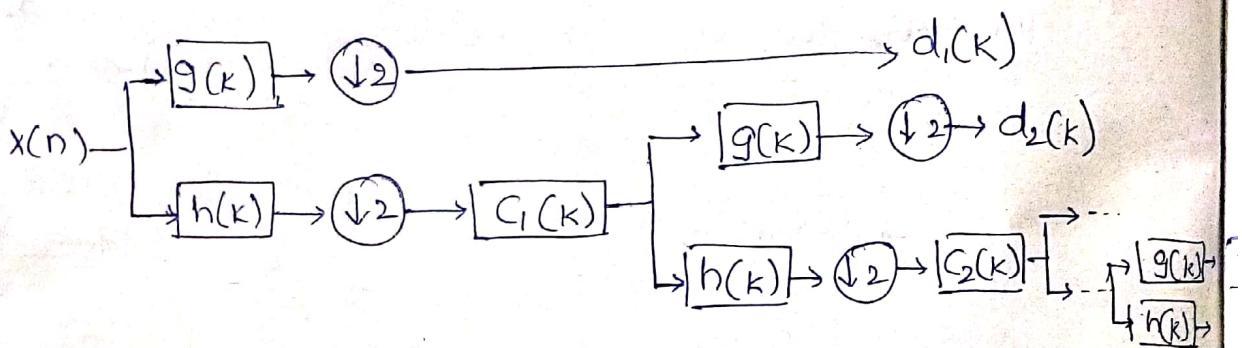
for $\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$ $\frac{1}{16}$

Discrete wavelet Transformation :-

In DWT, we given a data of size n and divided into $x = (x_1, x_2, x_3 \dots x_n)$ and each data part contains same length, we call this data parts as wavelet coefficients.

In DWT, There are many Algorithms used to analyse the data and one among is pyramid algo.

Pyramid Algorithm :-

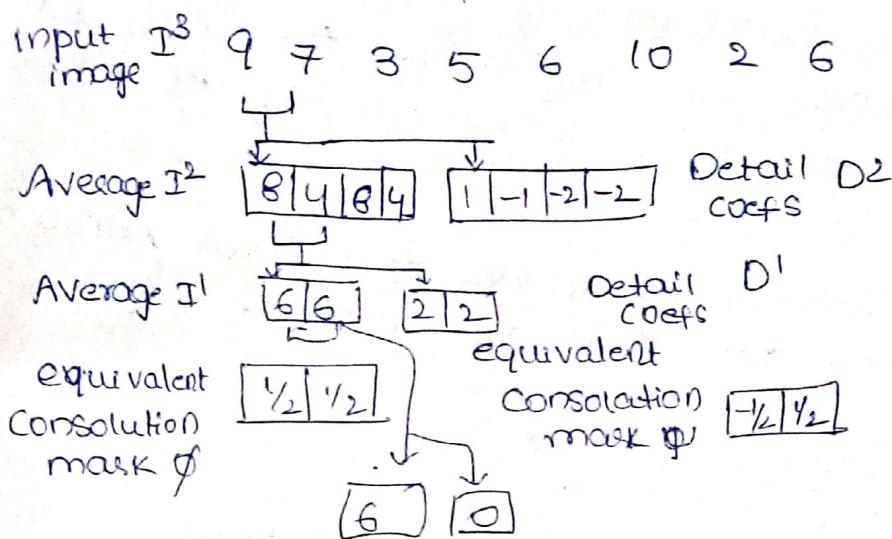


- * The length l of input data must be an integer power (2). This condition can be padding zeros to the data, if it is not sufficient in the power(2).
- * Each transformation involves applying 2 functions. The first applies some data smoothing such as sum or aggregation.

The second performs waited diff which acts to bring out the detail features of data.

- * The two functions are applied to pair of data points in x which is the length $1/2$.
- * The 2 functions are recursively applied to the data sets obtained in the previous loop until the result data sets which is transferred data at the power (2)

The selected values from the datasets obtained in previous iterations are decide the wavelength Coefficients of the transferred data.



Wavelet transformed image: 6 0 2 2 1 -1 -2 -2

Principal component Analysis :-

Ex:-	point	x	y	$x-\bar{x}$	$y-\bar{y}$	$(x-\bar{x})(y-\bar{y})$	$\frac{\Sigma \cdot 00}{3 \cdot 83}$	$\frac{3 \cdot 83}{1 \cdot 17}$
$\boxed{A(1)}$	A	1	1	-3.17	-2.83	8.97	6.00	
	B	2	1	-2.17	-2.83	6.14		$\frac{3 \cdot 83}{2 \cdot 17}$
	C	4	5	-0.17	1.14	0.19		
	D	5	5	0.83	1.14	0.97		
	E	5	6	0.83	-1.17	2.17	1.80	
	F	8	5	3.83	1.17	4.48		
		$\bar{x} = 4.16$	$\bar{y} = \frac{23}{6}$			$\Sigma = 22.55$		
				$= 3.83$				

$$\text{co-variance} = \frac{\Sigma(x-\bar{x})(y-\bar{y})}{n}$$

$$= \frac{22.55}{6} = 3.75$$

PCA is also called as K-L-method (Karhunen Loeve). It summarizes the data with less characteristic that represents the entire data set.

- PCA reduces the dimensionality of data containing large set of variables by transforming the initial variable into new set of variables without losing the important information of data.
- These new variables corresponds to a linear combination of the original data and are called principal components.

Usage of PCA:-

- To identify the hidden patterns in a dataset
- To reduce the dimensionality of data by removing noise and redundancy in the data.
- It is used to identify correlative variables

Algorithm for PCA:-

- The ip data are normalized so that each attribute fall within the same range.
- PCA computes or calculates K orthogonal vectors that provides a basis for normalized ip data.
- These are unit vector that each point in a direction + far to other one.
- The principal components are stored in ascending order. It is provided a new set of axis for data providing imp information about variance.
- The components are sorted in descending order of significance the data size can be reduced by eliminating the weaker components.

Ex:- $A \begin{pmatrix} x & y \\ 2 & 7 \\ 3 & 2 \\ 4 & 3 \end{pmatrix}$

$\frac{2+3+4}{3} \rightarrow \frac{3}{3} \quad \frac{7+2+3}{3} \rightarrow \frac{12}{3}$

$$C = A - M \rightarrow \text{mean} \quad \xrightarrow{\text{normalized}}$$

$$C = \begin{pmatrix} -1 & 3 \\ 0 & -2 \\ 1 & -1 \end{pmatrix}$$

characteristic equation $A\bar{V} = \lambda\bar{V}$

$$A\bar{V} - \lambda\bar{V} = 0$$

$$(A - \lambda I)\bar{V} = 0$$

$$(A - \lambda I)\bar{V} = 0$$

$$A - \lambda I = 0$$

first find covariance

$$\text{covariance } (C) = AT \cdot A / \frac{1}{n-1}$$

$$= \frac{1}{2} \begin{pmatrix} -1 & 0 & 1 \\ 0 & -2 & -1 \end{pmatrix} \begin{pmatrix} -1 & 3 \\ 0 & -2 \\ 1 & -1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & -2 \\ -2 & 7 \end{pmatrix}$$

characteristic eq $A\bar{V} = \lambda\bar{V}$

$$A\bar{V} - \lambda\bar{V} = 0$$

$$(A - \lambda I)\bar{V} = 0$$

$$(A - \lambda I)\bar{V} = 0$$

$$\begin{pmatrix} 1 & -2 \\ -2 & 7 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0$$

$$|ad - bc| = 0$$

$$\begin{pmatrix} 1-\lambda & -2-0 \\ -2-0 & 7-\lambda \end{pmatrix} = 0$$

$$(1-\lambda)(7-\lambda) - (-2)(-2) = 0$$

$$\lambda^2 - 8\lambda + 3 = 0$$

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{8 \pm \sqrt{64 - 4(1)(3)}}{2 \times 1}$$

$$= \frac{8 \pm \sqrt{52}}{2} = \frac{8 \pm 7.2}{2} = \frac{8+7.2}{2}, \frac{8-7.2}{2}$$

$$\lambda = 7.6$$

$$\lambda = \frac{15.2}{2}, \frac{0.8}{2} \Rightarrow 7.6, 0.4.$$

$$\left[\begin{pmatrix} 1 & -2 \\ -2 & 7 \end{pmatrix} - 7.6 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right] \begin{pmatrix} x \\ y \end{pmatrix} = 0$$

$$\begin{pmatrix} 1-7 \cdot 6 & -2 \\ -2 & 7-7 \cdot 6 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = 0$$

$$\begin{pmatrix} -6 \cdot 6 & -2 \\ -2 & -0 \cdot 6 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = 0$$

$$-6 \cdot 6x - 2y = 0$$

$$\frac{x}{2} = \frac{y}{-6 \cdot 6} \Rightarrow (2, -6 \cdot 6)$$

$$-2x - 0 \cdot 6y = 0$$

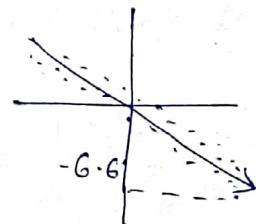
$$\frac{x}{0 \cdot 6} = \frac{y}{-2} \Rightarrow (0 \cdot 6, -2)$$

Eigen vector

Decreasing Order :-

$$(2, -6 \cdot 6) \checkmark$$

$$(0 \cdot 6, -2) \times$$



$$\lambda = 0 \cdot 4$$

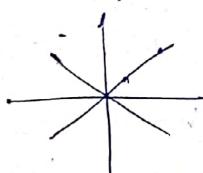
$$\left[\begin{pmatrix} 1 & -2 \\ -2 & 7 \end{pmatrix} - 0 \cdot 4 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right] \begin{pmatrix} x \\ y \end{pmatrix} = 0$$

$$0 \cdot 6x - 2y = 0 \Rightarrow (2, 0 \cdot 6)$$

$$-2x + 6 \cdot 6y = 0 \Rightarrow (6 \cdot 6, 2)$$

Decreasing order:

$$(6 \cdot 6, 2) (2, 0 \cdot 6) (2, -6 \cdot 6) (0 \cdot 6, -2)$$



Attribute subset selection (Feature selection)

The attribute subset selection also called as feature selection. It reduces the dataset size by removing irrelevant or redundant attributes or dimensions.

- The goal of attribute feature selection to find a minimum set of attributes that resulting probability distribution of data & classes is as close as possible to original data set.

Advantages:-

- It enables the machine learning algorithms to train very faster.
- It reduces the complexity of model and makes it easier to interpret.
- It improves accuracy of a model.

Feature Selection

Filter method

- Pearson's correlation
- LDA (linear discriminant analysis)
- ANOVA (Analysis of variables)
- chi-square

Wrapper method

- Forward selection
- Backward selection
- Combination of FS & BS
- Decision tree induction
- Recursive feature combination

Embedded method

- LASSO Regression
- RIDGE regression

Forward selection :-

This procedure starts with an empty set of attributes as a reduced set. The best of the original attributes is determined & added to the reduced set. At each subsequent iteration, the best of the remaining original attributes is added to the set.

Backward Selection: Elimination

The procedure starts with full set of attributes. At each step it removes the worst set of attributes remaining in the set.

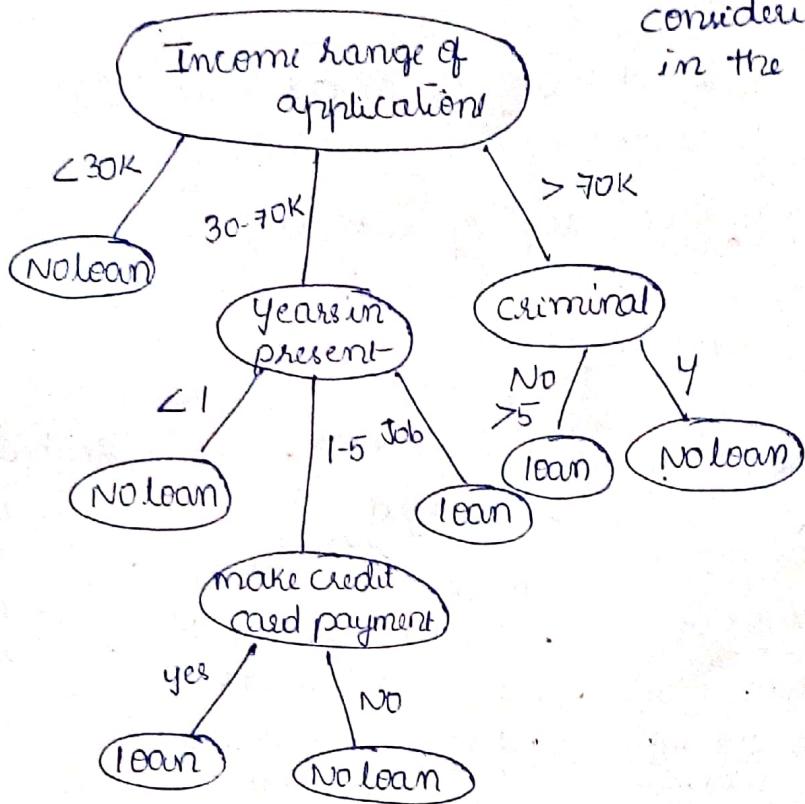
combination of forward selection & backward elimination

The step-wise forward selection & backward elimination methods can be combined so that each step the procedure select the best attribute & remove worst attribute in the set.

Decision tree elimination :-

- ID.3 C4.5 & CART

Decision tree is used to reduce the dimension by removing or not considering the attribute in the original dataset.



Parametric data reduction Technique :-

↳ Regression

↳ Log linear method

we can also called as numerosity reduction

- These two models used to reduce the size of data.
- In linear regression, the data or model to fit a st.line.
 $y = mx + b$ where x is independent variable
or predictor variable
 y is Response variable
- The data which is not fit in the line is called as outlier and this will be deleted from original dataset.

Loglinear model :-

This is used to approximate discrete, multidimensional probability distribution.

- estimate the probability of each point in a multi dimensional space for a set of discrete attributes based on a smaller subset of dimensional combination. $F = ae^{bx}$
where a, b - loglinear coeff

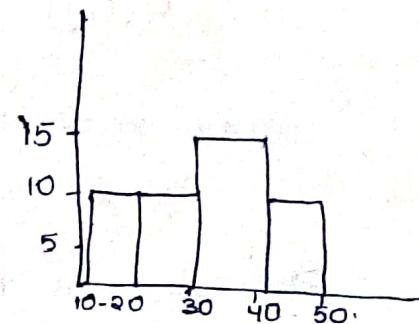
Non parametric Data Reduction:-

- Histogram - The graphical representation of frequency distribution of a continuous series.

Range 10-20 20-30 30-40 40-50

No. of persons 10 10 15 10

(Singleton bucket)



Similarity Measurements :-

The similarity between the two Obj will indicates the relationship b/w the two objects.

$$\text{Similarity of } (i,j) = 1 - \text{dissimilarity of } (i,j)$$

Similarity = 0 means there is no relation b/w obj
also indicated by dissimilarity = 1.

Sim = 1 means more similarity b/w data objects also seem as dissim = 0.

- To represent the similarity b/w Objects, we use 2 data structures - Data Matrix

- Dissimilarity matrix

✓ Ex

$$\begin{bmatrix} 0 \\ d_{2,1} \\ d_{3,1} d_{3,2} \\ d_{4,1} d_{4,2} d_{4,3} \\ 0 \end{bmatrix}$$

nominal - there is no order, the val may contain more than one.

order is Imp

Ex:-

Object identifier	Test-1 (nominal)	Test-2 (ordinal)	Test-3 (numerical)
1	codeA	Excellent	45
2	codeB	fair	22
3	codeC	good	64
4	codeA	Excellent	28

Proximity measures b/w nominal attributes :-

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

where m is no of matches
 'p' is number of attributes describe the object

From above Ex., $P=1$ which shows the
 no. of col to nominal

Obj test 1

1 codeA } No So, $m=0$
 2 codeB } matches

If codeA, codeA
 then $m=1$

$$\text{Now, } d(2,1) = \frac{P-m}{P}$$

$$= \frac{1-0}{1} = 1$$

$$\text{For } d(3,1) = \frac{P-m}{P}$$

1 codeA
 3 codec

$$= \frac{1-0}{1} = 1$$

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

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

$$d(4,2) = 1$$

$$d(4,3) = 1$$

Now dissimilar matrix,

As dissimilarity = 0

so, similarity for $d(4,1)$

is high.

0	1	0	
1	0		
0	1	0	
0	1	1	0

Four columns
 for four
 attributes

Proximity measures for numeric value :-

To find out the similarity b/w numeric values we use
 three methods :- Euclidean distance

$$d(i,j) = \sqrt{(x_{i1}-x_{j1})^2 + (x_{i2}-x_{j2})^2 + \dots + (x_{ip}-x_{jp})^2}$$

- Manhattan distance (or) city block distance :-

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

Ex:- $x_1 = (1, 2), x_2 = (3, 5)$
 $x_{i1} x_{i2}$ $x_{j1} x_{j2}$

Euclidean :-

$$\begin{aligned} d(i,j) &= \sqrt{(1-3)^2 + (2-5)^2} \\ &= \sqrt{4+9} = \sqrt{13} = 3.61 \end{aligned}$$

For weight Euclidean distance

Just multiply with w
values

$$\sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2}$$

Manhattan :-

$$d(i,j) = |1-3| + |2-5| = 5$$

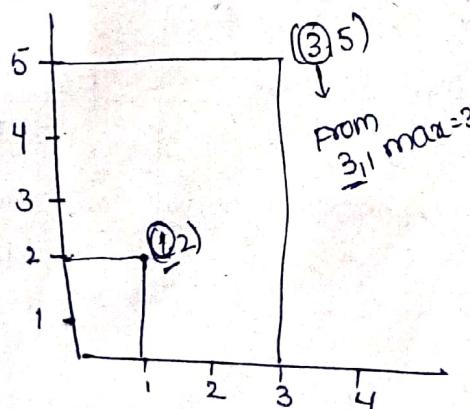
- Minkowski distance :-

$$d(i,j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

where h is a real number it represents the distance
of the (manhattan) when $(h=1)$. and for $(h=2)$, it represents
 $\frac{\text{Euclidean}}{\text{Euclidean}}$ distance. $\frac{\text{manhattan}}{\text{manhattan}}$ $\frac{\text{Euclidean}}{\text{Euclidean}}$

- Supreme Distance :- (or) L_{max} (or) L_∞ norm (or) Chebychev
distance

$$\begin{aligned} d(i,j) &= \lim_{h \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{1/h} \\ &= \max_f |x_{if} - x_{jf}| \end{aligned}$$



Ex:- $x_1 = (1, 2)$
 $x_2 = (3, 5)$

Supreme distance = $\max_f (x_{if} - x_{jf})$

$$= |(3)-5| = 2$$

max

To find out dissimilarities for numerical values

$$d_{ij}^f = \frac{|x_{if} - x_{jf}|}{\max_h x_{hf} - \min_h x_{hf}} \quad \text{where } f \text{ is numeric}$$

From above Ex, table

type 3

$$d(2,1) = \frac{|22 - 45|}{\frac{64 - 22}{\substack{\text{max value} \\ \text{from entire} \\ \text{table}} \text{ min value}}} = 0.55 \quad \begin{matrix} 1 & 45 \\ 2 & 22 \end{matrix}$$

$$d(3,1) = \frac{|64 - 45|}{42} = \frac{19}{42} = 0.452$$

$$d(3,2) = \frac{|64 - 22|}{42} = \frac{42}{42} = 1$$

$$d(4,1) = \frac{28 - 45}{42} = \frac{17}{42} = 0.404$$

$$d(4,2) = \frac{28 - 22}{42} = \frac{6}{42} = 0.14$$

$$d(4,3) = \frac{|28 - 64|}{42} = \frac{36}{42} = 0.85$$

1.00
0.14
0.86

dissimilarity

0	
0.55	0
0.45	1
0.40	0.14
	0.860

For this
Sim=0

Similarity = 1 - dissimi

so, consider
this

Similarity = 0.9

1-0.14

for proximity measures for Ordinal

From above ex:-

we provide ranks 4 refers to rank $\frac{r_{if} - 1}{m_f - 1}$ max rank	(3) Excellent (1) (1) fair (0) (2) good (0.5) (3) Excellent (1)	} values got from formula
---	--	------------------------------

dissimilarity matrix

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

For $d(3,1)$

check 3 value = 0.5 diff = 0.5
 1 value 1

Proximity measures for Binary Attributes :-

The Binary attribute values may be 0 or 1, +ve or no

'0' indicates obj is absent '1' refers to obj is present.

$$P = q + r + s + t$$

$$q = \sum_{\substack{i=1 \\ j=1}}^{} q_{ij}, \quad r = \sum_{\substack{i=1 \\ j=0}}^{} q_{ij}$$

$$s = \sum_{\substack{i=0 \\ j=1}}^{} q_{ij}, \quad t = \sum_{\substack{i=0 \\ j=0}}^{} q_{ij}$$

$$4, P_i = \frac{1}{N} = 0$$

	Name	gender	fever	cough	t-1	t-2	t-3	t-4
Jack	m	y	n	p	n	n	n	n
Jim	m	y	y	n	n	n	n	n
Mary	f	y	n	p	n	p	n	n

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

For q $i=1$
 $j=0$

check for P or q
 Jack, Jim N or y

$$\text{dissimilarity} (Jack, Jim) = \frac{1+1}{1+1+1} = \frac{2}{3} = 0.67$$

$$d(Jim, Mary) = \frac{1+2}{1+1+2} = \frac{3}{5} = 0.6$$

$$\text{Similarity more} d(Jack, Mary) = \frac{1+1+2}{0+1/2+1} = \frac{1}{3} = 0.33$$

SL302G
 31/02/2019
 9

Proximity measure for mixed attributes

Nominal

0		
1	0	
1	1	0
0	1	1

ordinal

0		
1	0	
0.5	0.5	0
0	1	0.5

(Nominal) Numerical

0		
0.55	0	
0.45	1	0
0.4	0.14	0.86

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

↑ No. of attributes

distance of par position
if all values are present in par attribute then $S=1$ else $S=0$

$$\begin{aligned} S_{ij}(d_{ij}) &= \frac{1(1) + 1(1) + 1(0.55)}{1+1+1} \\ &= \frac{2+55}{3} = 0.85 \end{aligned}$$

From nominal

$$\begin{aligned} S_{ij}(d_{ij}) &= \frac{1(1)(d_{2,1})}{1(d_{2,1})} \\ &\text{from numeric} \end{aligned}$$

BJ 255 (B:
24
15)

$$\text{For } d(3,1) = \frac{1(1) + 1(0.5) + 1(0.45)}{3} = \frac{1+0.5+0.45}{3} = 0.65$$

$$d(3,2) = \frac{1(1) + 1(0.5) + 1(1)}{3} = \frac{2.5}{3} = 0.83$$

$$d(4,1) = \frac{1(0) + 1(0) + 1(0.4)}{3} = \frac{0.4}{3} = 0.133$$

15
0.88
236

$$d(4,2) = \frac{1(1) + 1(1) + 1(0.14)}{3} = \frac{2.14}{3} = 0.713$$

$$d(4,3) = \frac{1(1) + 1(0.5) + 1(0.86)}{3} = \frac{1.5 + 0.86}{3} = 0.78$$

dissimilar more	0	
	0.85	0
	0.65	0.83
similar more	0.133	0.71 0.78 0

To find out the similarity b/w 2 obj in asymmetric data we use

Jaccard Quotient

$$\text{Sim}_{ij} = \frac{q_{ij}}{q_{j+it}}$$

$$\text{dis}_{ij} = \frac{q_{it}}{q_{j+it}}$$

Cosine similarity :-

It is used to handle the sparse data (the most of the data be zero in matrix) In information retrieval system the most of the data is absent. For those documents the normal similarity measure will not produce accurate similarity.

The cosine similarity is used to handle this sparse data if the resultant value equal to 1 (angle b/w x, y is 0°). Hence, the documents are more similar. If the cosine value is nearer to zero (90°) then the documents are dissimilar.

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

Ex:-

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document 1	5	0	3	0	2	0	0	2	0	0
→ 2	3	0	2	0	1	1	0	1	0	1
3	0	7	0	2	1	0	0	3	0	0
4	0	1	0	0	1	2	2	0	3	0.

Term-frequency table or Document vector

$$\text{Sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|} = \frac{25}{(6.48)(4.12)} = 0.94 \text{ No similarity}$$

$$x \cdot y = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1 + 0 \times 0 + 0 \times 1 \\ = 25$$

$$\|x\| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 1^2} = 6.48$$

$$\|y\| = \sqrt{3^2 + 0^2 + 2^2 + 1^2 + 1^2 + 1^2 + 1^2} = 4.12$$

- If the document attributes are binary attributes then the

$$\text{Sim}(x, y) = \frac{x \cdot y}{x \cdot x + y \cdot y - x \cdot y}$$

Data Transformation :-

- Types:-
- Smoothing
 - Attribute construction
 - Aggregation
 - Generalization
 - Normalization

Data are transferred or consolidated into the req format so that, the resultant mining process may be more efficient and the patterns are easier to understand.

* Normalization :-

In this normalization, the attribute data are scaled within a specified range $[1, 1]$ or $[0, 1]$

- Types of normalization - min-max normalization
- z-score normalization
 - Decimal scaling

→ min-max normalization :-

$$v' = \frac{v - \text{min}}{\text{max} - \text{min}} (\text{new max} - \text{new min}) + \text{new min}$$

Ex:- marks 8 10 15 20

$$\text{new max} = 1$$

$$\text{new min} = 0$$

*This is
numerical
value*

$$v' = \frac{8-8}{20-8} (1-0) + 0 = 0$$

$$v' = \frac{20-8}{20-8} (1-0) + 0 = 1$$

$$v' = \frac{10-8}{20-8} (1-0) + 0 = \frac{1}{6} = 0.1$$

$$[0, 0.1, 0.5, 1]$$

$$v' = \frac{15-8}{20-8} (1-0) + 0 = \frac{7}{12} = 0.5$$

→ Z-score normalization :-

$$S.D = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$Z\text{-Score} = \frac{x - \mu}{\sigma}$$

Ex:- marks 8 10 15 20

$$\bar{x} = \frac{8+10+15+20}{4} = 13.25$$

For 8

$$S.D = \sqrt{\frac{(8-13.25)^2 + (10-13.25)^2 + (15-13.25)^2 + (20-13.25)^2}{3}}$$
$$= 5.37$$

$$Z\text{-Score} = \frac{8 - 13.25}{5.37} = -0.97$$

For 10

$$Z\text{-Score} \frac{10 - 13.25}{5.37} = -0.60$$

$$\text{For } 15 : \frac{15 - 13.25}{5.37} = 0.32$$

For 20 : $\frac{20 - 13.25}{5.37} = 1.25$ data must be in
[−1, 1] range but it exceeds
so we didn't consider.

→ Decimal Scaling :-

$$V' = \frac{V}{10^j} \quad j \text{ represents no. of digits}$$

10 No. of digits = 2

$$V' = \frac{8}{10^1} = \frac{8}{10} = 0.8$$

$$V' = \frac{10}{10^2} = 0.1$$

$$V' = \frac{15}{10^2} = 0.15$$

$$V' = \frac{20}{10^2} = 0.2$$

Smoothing :- which works to remove the noise from the data by using the techniques

1. clustering

2. Binning

3. Regression

Aggregation :- where the summary and aggregation operation are applied to the given data.

Ex:- Rollup & cube operation in data ware house



Generalization :-

The data where low level are primitive data are placed by the higher level concept by using hierarchy.

Customer → higher level

↓
Address → low level

↓
Street → low level

Attribute construction:-

where new attributes are constructed and added from the given set of attributes to help the mining aspects Ex:- PCA

UNIT - 3

ASSOCIATION RULES

Association rule:- This is used to findout the hidden patterns in the given data.

It indicate the probability of buying one obj if you buy another obj. In association rules we use two terminologies-

- support
- confidence

- * Support means the probability to participate in a transaction or event out of all the transactions.
- * confidence indicate the probability of two events occurred in a transaction.

Ex:- Bread \Rightarrow Butter (20%, 40%)

Here, 20% indicates out of 100 events the 20 people are going to buy the bread is called support out of 20 members 40% of the people will take butter along with bread we called as confidence.

Basic Algorithm :-

- Apriori Algorithm
- . sampling Algorithm
- . partitioning Algorithm

Apriori Algo :-

It is an algorithm for data mining to find out frequent item sets for boolean association rules. It is follows an bottom up approach where frequent subsets are entered one item at a time which is known as a candidate key generation and grouping of candidates are tested against the data.

Support - A support for an association rule $x \Rightarrow y$ is the percentage of transaction in the database that contains $x \cup y$.

$$\text{Support}(x, y) = \frac{s(x \cup y)}{n}$$

confidence (α) or strength :-

confidence for an association rule $x \rightarrow y$ is
the ratio of no. of transactions that contains $x \cup y$ to the
no. of transactions that contains ' x '.

$$\alpha = \frac{S(x \cup y)}{S(x)}$$

By using large item set or frequent set we find an
association rules for the given data. It indicates the no. of
occurrences of an item set above the threshold of S . we
also called as downward closure.

Ex:- For the following given transaction data set generates
rules using Apriori Algorithm. consider the value of support
= 50% and confidence = 75%.

Transaction id	Items purchased
1	Bread, cheese, Egg, Juice
2	Bread, cheese, Juice
3	Bread, milk, Yogurt
4	Bread, Juice, milk
5	Cheese, Juice, milk

→ Find the frequency of each item

frequent item set :-

Items	frequency	Support
Bread	4	$4/5 = 80\%$
cheese	3	$3/5 = 60\%$
Egg	1	$1/5 = 20\%$
Juice	4	$4/5 = 80\%$
milk	3	$3/5 = 60\%$
Yogurt	1	$1/5 = 20\%$

Egg & Yogurt is not Satisfy the threshold S we can
remove.

→ make 2 items as candidate set and write their frequency.

	frequency	support
{ Bread, cheese }	2	$2/5 = 40\%$
{ Bread, juice }	<u>3</u>	$3/5 = 60\% \checkmark$
{ Bread, milk }	2	$2/5 = 40\%$
{ cheese, juice }	3	$3/5 = 60\% \checkmark$
{ cheese, milk }	1	$1/5 = 20\%$
{ juice, milk }	2	$2/5 = 40\%$

$$\begin{aligned} \{ \text{Bread, juice} \} &\leftarrow \begin{array}{l} \text{Bread} \Rightarrow \text{Juice} \\ \text{Juice} \Rightarrow \text{Bread} \end{array} \quad \alpha = \frac{s(\text{Bread} \cup \text{Juice})}{s(\text{Bread})} = \frac{3}{4} = 75\%. \\ \{ \text{cheese, juice} \} &\leftarrow \begin{array}{l} \text{cheese} \Rightarrow \text{Juice} \\ \text{Juice} \Rightarrow \text{cheese} \end{array} \quad \alpha = \frac{3}{4} = 75\%. \\ & \qquad \qquad \qquad \alpha = \frac{3}{3} = 100\% \quad \text{all satisfied} \\ & \qquad \qquad \qquad \alpha = \frac{3}{4} = 75\% \quad \text{the confidence.} \end{aligned}$$

Ex :-

1	A, B, C	Support = 2%
2	A, C	
3	A, D	confidence = 50%.
4	B, E, F	

A	3	$3/6 = 50\%$	$\frac{1}{3} \times \frac{1}{10} = \frac{1}{30}$
B	2	$2/6 = 33\%$	
C	2	$2/6 = 33\%$	
D	1	$1/6 = 16\%$	
E	1	$1/6 = 16\%$	
X - F	1	$1/6 = 16\%$	

X { A, B }	1	
{ A, C }	2	$2/6 = 33\%$
X { B, C }	1	

{ A, C }	$\left\{ \begin{array}{l} A \Rightarrow C \\ C \Rightarrow A \end{array} \right.$	$\alpha = \frac{s(A \cup C)}{s(A)} = \frac{2}{3} = 66.6\% \checkmark$
		$\alpha = \frac{s(C \cup A)}{s(C)} = \frac{2}{2} = 100\% \checkmark$

Notations :-

$\alpha \rightarrow$ confidence

$S \rightarrow$ support

$D \rightarrow$ Database of transaction

$t_i \rightarrow$ Transaction in D

$x, y \rightarrow$ Item sets

$x \Rightarrow y \rightarrow$ Association rule

$L \rightarrow$ Large item set

$l \rightarrow$ Large item set in L

$C \rightarrow$ Set of candidate item sets

$P \rightarrow$ Number of partitions

Ex:- $B_1 (M, C, B)$

$B_6 (M, C, B, J)$

$B_2 (M, P, J)$

$B_7 (C, B, J)$

$B_3 (M, C, B, J)$

$B_8 (B, C)$

$B_4 (C, J)$

$B_5 (M, P, B)$

Support

$MCB \Rightarrow MC \rightarrow B$

$$= 3/3 = 1$$

$$\Rightarrow MB \rightarrow C = 3/4 = 0.75$$

$$\Rightarrow BC \rightarrow M = 3/5 = 0.6$$

M	5
C	6
B	6
P	X
J	5

3-combi

$$\begin{aligned} CBJ &\Rightarrow CB \rightarrow J = 3/5 \\ CJ \rightarrow B &= 3/4 \\ BJ \rightarrow C &= 3/3 \end{aligned}$$

M, C	3
2 Combi M, B	4
M, J	3
C, B	5
C, J	4
B, J	3

MCB - 3

MCJ - 2 X

CBJ - 2 X

CBJ - 3

MBJ - 2 X

MCB - 3

CBJ - 3

S(MUCOB)

S(CJ)

Fp-Growth Algo :- (Frequency pattern)

- The Drawback of Apriori Algorithm is, we repeatedly scan the database and it will take high expensive when the datasets are very large.
- FP-growth means Frequency pattern growth Algo used for avoid the drawback for Apriori Algo by considering FP tree where each node represent the support count of itemset and the branch b/w the nodes will indicate the association of the itemsets.

Ex:- Generate FP-tree for the following transaction dataset
(min support = 30%).

Transaction-id	Items
1	E, A, D, B
2	D, A, C, E, B
3	C, A, B, E
4	B, A, D
5	D
6	D, B
7	A, D, E
8	B, E

Items	frequency	priority
A	5	3
B	6	1
C	3	5
D	6	2
E	4	4

frequencies
 All are greater than 3 which is Support
 If value is less than support then remove it.

order set \Rightarrow B, D, A, E, C
of frequency item

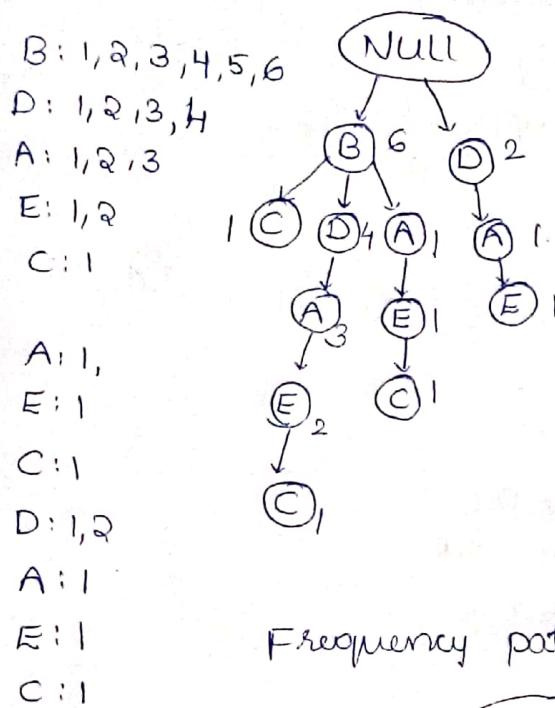
* Order the items According to priority

Transid:	Items	ordered items
1	E A D B	B D A E
2	D A C E B	B D A E C
3.	C A B E	B A E C
4	B A D	B D A
5	D	D

By preferring this order

6	D, B	B, D
7	A, D, E	DAE
8	B, C	B, C

FP-tree :



conditional item sets

B	NULL
D	$\langle B:4 \rangle$, NULL
A	$\langle D:1 \rangle$, $\langle B:3 \rangle$, $\langle B:1 \rangle$, $\langle D:A:1 \rangle$, $\langle BDA:1 \rangle$
E	$\{ A \leftarrow E \}$, $\{ B \leftarrow C \}$, $\langle B:1 \rangle$, $\langle BDAE:1 \rangle$
C	$\langle BAE:1 \rangle$

Common letters in all

Frequency pattern



Ex-2 :-

Transid	itemsets	Support $\rightarrow 3$
1	MONKEY	
2	DONKEY	
3	MAKE	
4	MOCKY	
5	COOKIE	

item frequency priority

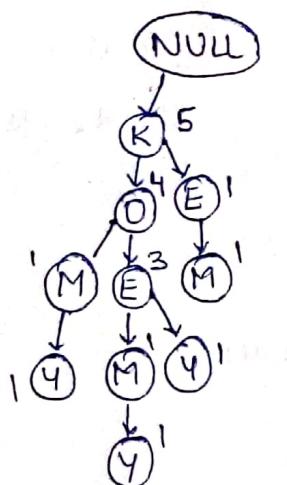
I	1	X	support 8	
M	3	X	3	
O	5		4	
N	2	X	support 6	
K	5		6	
E	4		1	
Y	3		3	
D	1	X	5	
A	1	X	9	
C	2	X	10	

order of frequency item

DKEMY

Transid	items	ordered items
1	MONKEY	KOEMY
2	DONKEY	KOEY
3	MAKE	KEM
4	MOCKY	KOMY
5	COOKIE	KOE

FP-tree :-



K:1,2,3,4,5
 O:1,2,3,4
 E:1,2,3
 M:1
 Y:1
 Y:1
 E:1
 M:1
 M:1
 Y:1

Conditional item sets :-

K	NULL	
O	<K:4>, A,B,C	→ K:4
X ^{A-2} X ^{B+1} X ^{C-1}	E	
K:3, O:2	M	
X ^{A-1}	Y	
K:3, O:3, E:2		
E:2		

<K:4>, ~~A,B,C~~ → K:4
 <KO:3>, <K:1> → K:4, O:3
 <KOE:1>, <KO:1>, <KE:1> → K:3, O:2, E:2
 <KOEM:1>, <KOM:1>, <KOE:1> → K:3, O:3, M:2, E:1

Frequency pattern

$$\begin{array}{cccc}
 E \rightarrow K & M \rightarrow K & Y \rightarrow O & Y \rightarrow K \\
 K \rightarrow E & K \rightarrow M & O \rightarrow Y & K \rightarrow Y
 \end{array}$$

Frequency pattern

K → Y, O → Y, KO → Y, K → M, K → O, K → E, O → E, KO → E

Ex-3

Trans-id items

1	I_1, I_3, I_4	Support - 2
2	I_2, I_3, I_4	
3	I_1, I_2, I_3, I_4	

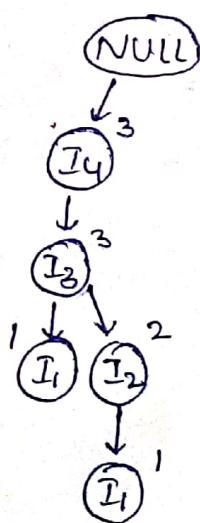
item frequency priority

I_1	2	4	Order I_4, I_3, I_2, I_1
I_2	2	3	
I_3	3	2	
I_4	3	1	

Trans-id items ordered items

1	I_1, I_3, I_4	I_4, I_3, I_1
2	I_2, I_3, I_4	I_4, I_3, I_2
3	I_1, I_2, I_3, I_4	I_4, I_3, I_2, I_1

FP-tree



$I_4: 1, 2, 3$

$I_3: 1, 2, 3$

$I_1: 1$

$I_2: 1, 2$

$I_1: 1$

conditional item sets :-

I_4 NULL

I_3 $\langle I_4: 3 \rangle$ $I_4 = 3$

I_2 $\langle I_4, I_3: 2 \rangle$ $I_4: 2, I_3: 2$

I_1 $\langle I_4, I_3: 1 \rangle, \langle I_4, I_3, I_2: 1 \rangle$ $I_4: 2$

$I_3: 2$

$I_2: 1$

Frequency pattern

$I_4 \rightarrow I_3$

$I_4 \rightarrow I_2$ $I_3 \rightarrow I_2$, $I_4, I_3 \rightarrow I_2$

$I_4 \rightarrow I_1$, $I_3 \rightarrow I_1$, $I_4, I_3 \rightarrow I_1$

2. Sampling Algorithm :-

It is used to find out the frequency itemset from given large dataset. It is difficult by using apriori algo. Instead of working on entire population we will take some sample and apply the apriori algorithm which results indicates the frequency items in entire population.

Potentially Large Frequency Itemset (PL) :- It indicates the large itemset in sample.

Negative Border (BD) :-

- * It is generalization of apriori gen algorithm.
- * The minimal set of itemsets which are not in PL, But the samples of items are present in PL.

Algorithm :-

D_S D_S = Sample of database

PL = Large itemsets in D_S using "smalls"

G = $PL \cup BD^-(PL)$

Count c in database using S

ML = Large itemsets in $BD^-(PL)$

if $ML = \emptyset$ done

else

G = repeated application of BD^- count c in db.

Example :-

Trans ID	Itemset-	
1	Bread, Jelly, PB	Support = 20%
2	Bread, PB	
3	Bread, milk, PB	Smalls = 10%
4	Beer, Bread	
5	Beer, milk	

PB - peanut Butter.

$$DS = \{T_1, T_2\}$$

Data set - select any of two transactions

$T_1 = \{ \text{Bread, Jelly, PB} \}$

T₂ = {Bread, PB}

$$P_L = \{ \{ \text{Bread} \}, \{ \text{Jelly} \}, \{ \text{PB} \}, \{ \text{Bread, Jelly} \}, \{ \text{Bread, PB} \}, \{ \text{Jelly, PB} \}, \\ \{ \text{Bread, Jelly, PB} \} \}$$

$$BO^-(PL) = \{\{Beer\}, \{milk\}\} \rightarrow \text{remaining items not repeated in PL}$$

$$\downarrow \mathcal{G} = \text{PL} \cup \text{BD}^-(\text{PL})$$

canned set = $\{\{Bread\}\{Jelly\}\{PB\}\{Bread, Jelly\}\{Bread, PB\}\{Jelly, PB\}$
 $\{Bread, Jelly, PB\}\{Beer\}\{milk\}\}$

$$\text{Support} = \frac{20}{100} \times 5 \rightarrow \text{No. of Transactions}$$

= 1 (remove the items that are not repeated atleast once)

$$C(B\bar{O}) = \left\{ \begin{array}{l} \{\text{Beer, Bread}\} \times \{\text{Beer, Jelly}\} \times \{\text{Beer, PB}\} \times \{\text{Beer, milk}\} \\ \{\text{Bread, milk}\} \times \{\text{Jelly, milk}\} \times \{\text{milk, PB}\} \end{array} \right\}$$

$$\text{Support} = \frac{20}{100} \times 5 = 1 \quad (\text{remove items that are not repeated atleast once as a above combo})$$

$$= \{\{Beer, Bread\}, \{Beer, milk\}, \{Bread, milk\}, \{milk, PB\}\}$$

$$(A, B) \cup (B, C)$$

$$\{ \{ \text{Beer}, \cancel{\text{Bread}}, \text{milk} \} \{ \text{Beer}, \cancel{\text{milk}}, \text{PB} \} \} \checkmark \text{Bread, milk, PB}$$

Again check in support = 1

$\hat{Y} = \{\text{Bread, milk, PB}\}$

partition Algorithm :-

To find out the frequency itemset for entire db in apriori Algorithm is difficult because of memory restrictions.

- * By using partition Algo, we can divide entire dataset into diff partitions & allocate each partition to the processor to find out frequency itemset of each partition. At the last combine all (partition) items generated by each partition and apply apriori gen algo to find out large frequency itemset.

Ex:-

	A ₁	A ₂	A ₃	A ₄	A ₅
S=20% P ₁	0	1	1	0	0
	1	0	1	0	0
	2	1	0	1	0
	3	0	1	0	0
	4	0	0	1	0
	5	0	1	0	0
	6	0	1	0	1
S=20% P ₂	7	0	0	1	0
					1

Support = 20%.

$$\text{Support } (P_1) = \frac{20}{100} \times 4 = 0.8 \text{ (Every state must enter atleast once)}$$

$$L^1 = \{ \{A_1\} \{A_2\} \{A_3\} \{A_5\} \{A, A_2\} \{A, A_3\} \{A, A_5\} \{A_3, A_5\} \}$$

↓

$\{A, A_3, A_5\}\}$

check in single combo If 1 present
whether it enter one time in both A₁, A₂

$$P_2 : L^2 = \{ \{A_2\} \{A_3\} \{A_4\} \{A_5\} \{A_2, A_4\} \{A_3, A_5\} \}$$

$$L = P_1 \cup P_2$$

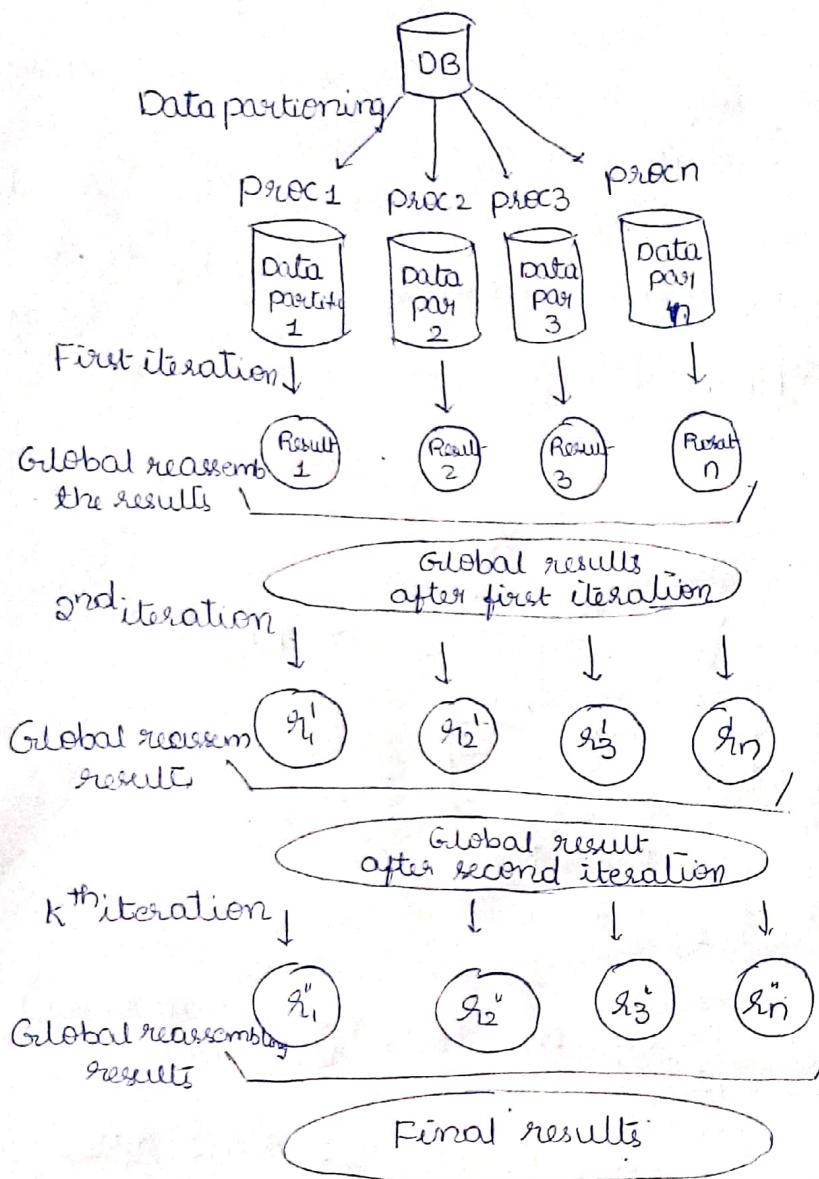
$$\text{check: } 8 \times \frac{20}{100} = 1.6 \approx 2$$

$$\therefore L = \{ \{A_1\} \{A_2\} \{A_3\} \{A_5\} \{A_3, A_5\} \}$$

Total no of 5
Transactions.

By using check, As check = 2 we need to check in above L¹, L² whether the inside Attributes repeat 2 times or more than 2 then write to L.

Parallel & Distributed Association Algo :- (Count Distribution Algorithm)



In most Association rules, the parallelism mechanism will be taken by 2 ways :- → Data parallelism
→ Task parallelism

- * In Data parallelism, The data is distributed into diff processors and in Task parallelism, the candidate sets are distributed to the diff processors.
- * In data parallelism, the Algorithms reduce communication cost over the tasks.
- * While task parallelism, not only candidates but also the local set of transactions must be broadcast to all other sides.

+ In data parallelism, algo requires memory at each processor will large enough to store all the candidate at each scan but in task parallelism, it can be avoid because the subsets of the candidates that are assign to the processor during each scan must fit into memory.

Ex:	Transid	Itemset
	1	A, B, C
	2	A, C, B
	3	B, A, A
	4	C, B, D

processor 1

1	ABC
2	ACB

processor 2

3	B A(A)
4	C B D

consider only
one in one set

$$\begin{aligned} A &\rightarrow 2 \\ B &\rightarrow 2 \\ C &\rightarrow 2 \\ D &\rightarrow 0 \end{aligned}$$

$$\begin{aligned} A &\rightarrow 1 \\ B &\rightarrow 2 \\ D &\rightarrow 1 \\ C &\rightarrow 1 \end{aligned}$$

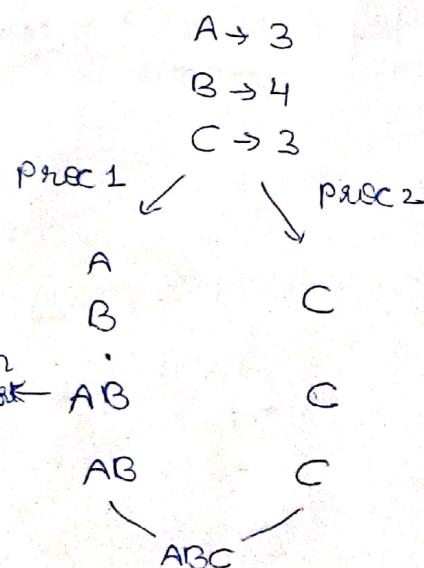
$$\text{Support} = 40\%$$

$$\frac{240}{600} \times 4 = 2$$

merging both

$$\begin{aligned} A &\rightarrow 3 \\ B &\rightarrow 4 \\ C &\rightarrow 3 \\ D &\rightarrow 1X \end{aligned}$$

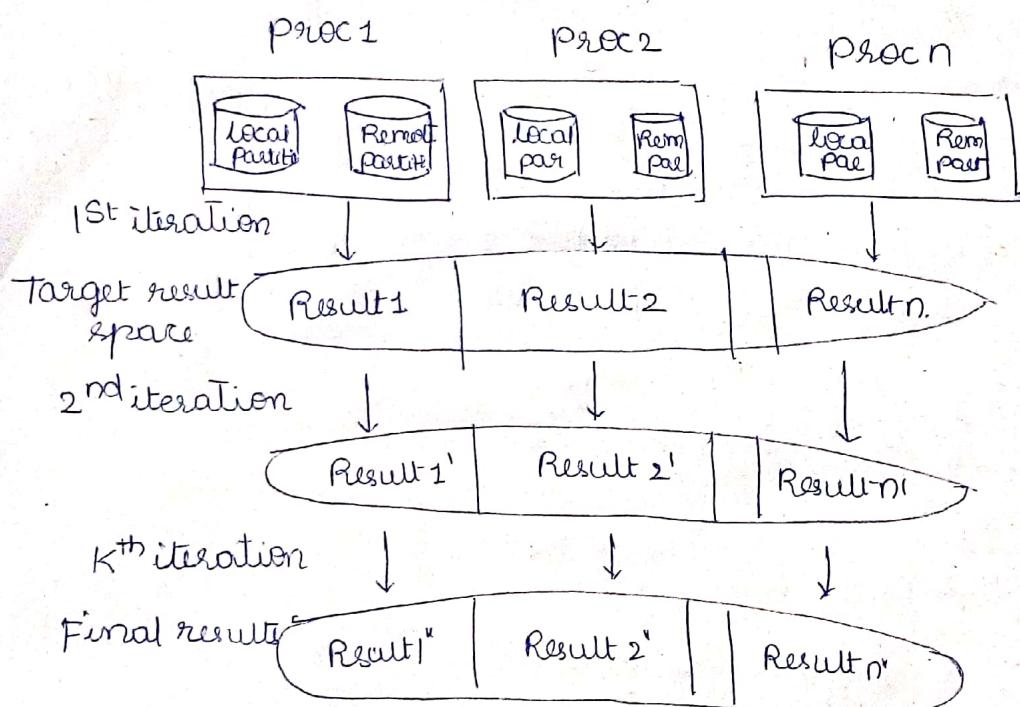
$$\begin{aligned} A &\rightarrow 3 \\ B &\rightarrow 4 \\ C &\rightarrow 3 \\ D &\rightarrow 1X \end{aligned}$$



- At the end of each iteration, since support of each candidate itemset in each processor is incomplete, each processor will need to redistribute the count to all processor. Hence term Count distribution is used.

Task parallelism :-

- * It is also called as data distribution Algo.
- * The candidates in task parallelism is partitioned among the processors.
- * G_k^l indicates the candidate of size K , find out at processor P_l .
- * L_k^l indicate the local large K items at processor L .



Ex:- TransId Items

t_1	Br, Jelly, PB	Support = 2
t_2	Br, PB	
t_3	Br, milk, PB	
t_4	Beer, Br	
t_5	Beer, milk	

P_1 [Bread
Beer]

P_2 [Jelly
milk]

$P_3 - PB$.

P¹: t₁t₂

P²: t₃t₄

P³: t₅

Bread = 2

Jelly = 0

Beer = 0

Milk = 1

PB = 0

Broadcast data partition to diff processes.

P¹: t₃t₄

P²: t₅

P³: t₁t₂

Br: 2+2

Jelly: 0+0

PB: 0+2

Be: 0+1

Milk: 1+1

P¹: t₅

P²: t₁, t₂

P³: t₃t₄

Br: 4

Milk: 2

PB: 3

Be: 2

~~Jelly: 1~~

support < 2

Q-item sets :-

{Br, Milk} {Be, Milk} {Br, B} {m, PB} {Br, PB} {Be, PB}

Advance Association rule mining / Techniques :-

In advance association rule technique we are

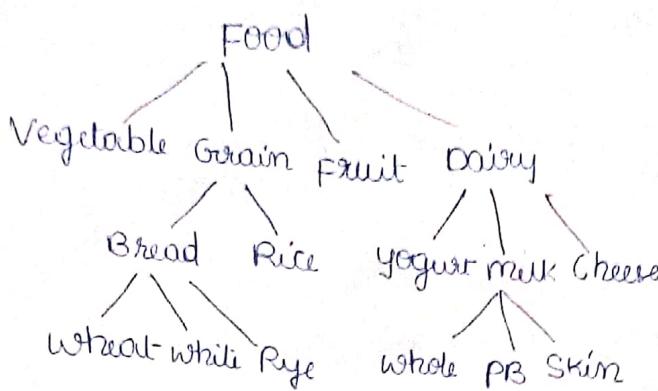
addressing -

- How to handle the Quantitative dataset
- How to define handle the support value for the given dataset
- How to find out the frequency itemsets in Transactional db

To Address the above problems we use following association rule Techniques:-

- Generalized association rules
- Multilevel association rule
- Quantitative association rule
- Multiple minimum Support Techniques
- Correlation rules.

Generalized Association rules:-



By using the concept of hierarchy that shows the relationship b/w the diff items, generalized association rules.

- A generalized association rule $x \Rightarrow y$ indicates that no item in 'y' may be above any item in x

$x \Rightarrow y$ means we can't generate association rules from lower level to higher level

Ex:- $Bread \Rightarrow PB$

$Grain \Rightarrow PB$

Out of these two association rules $Grain \Rightarrow PB$ has highest support value compared to $Bread \Rightarrow PB$ because Grain having higher level than bread.

Multilevel Association rules:-

A variation of generalised association rule is called as multilevel Association rule

In multilevel association rule the items may occur in any level in hierarchy

- Using a variation of apriori algo the concept hierarchy is traversed in top down manner to generate large items.
- The minimum support in all nodes in hierarchy at same level is identical.
- If α_i is min support for level i in hierarchy α_{i-1} is min support for level $i-1$ then $\alpha_{i-1} > \alpha_i$

$Vegetable \rightarrow Grain$

$Grain \rightarrow Yogurt$

Quantitative Association rules :-

In Quantitative Association rule used to find the Association rules when the Database contains the quantity data.

By using the discretization method called as equal width, we divide the given quantity data into equal width.

- Apply the support after discretization to find out Association rules.

Ex:-

Cid	Age	married	cars
1	21	No	0
2	23	Yes	1
3	29	No	1
4	30	Y	2
5	34	Y	2

✓

Age	
21	20:24
23	20:24
29	25:29
30	30:34
34	30:34



Cid	Age	married	cars	
1	20:24	N	0	0:1
2	20:24	Y	1	0:1
3	25:29	N	1	0:1
4	30:34	Y	2	2
5	30:34	Y	2	2

itemset
Age

$\langle \text{Age: 20:24} \rangle$

$\langle \text{Age: 25:29} \rangle$

$\langle \text{Age: 30:34} \rangle$

$\langle \text{married: Yes} \rangle$

$\langle \text{married: No} \rangle$

$\langle \text{cars: 0:1} \rangle$

$\langle \text{cars: 2:3} \rangle$

Support
Age 20:24 appears
twice Support = 2

② → two times in above

1 - X

2

3

2

3

2

$\langle \text{Age: 30:34} \rangle \langle \text{married: Yes} \rangle$

\downarrow
 $\langle \text{cars: 2:3} \rangle$

X $\langle \text{married: Yes} \rangle \Rightarrow \langle \text{cars: 2} \rangle$

→ Multiple min support —
Rare item problem :-

If the support is too low, then more no. of items
Consider for frequent itemset which are not useful.

- If support value is too high then less no. of frequent itemset in this we can lose frequent items.
- To solve the problem, we use two techniques.
 - partition (clustering) data sets & Apply Association rule.
 - MIS Apriori

Ex:- MIS Apriori

Consider itemset $\{A, B, C\}$

$$mis(A) = 20\%$$

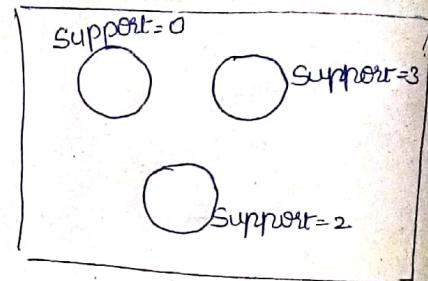
$$mis(B) = 3\%$$

$$mis(C) = 4\%$$

for 2 combos

$$\begin{aligned} AB &= \min\{mis(A), mis(B)\} \\ AC &= 4\% = \min\{20\%, 3\%\} \\ &= 3\% \end{aligned}$$

clustering



Correlation Techniques:-

The correlation used to find out the relation b/w two frequent itemsets.

If the value of correlation = 1, then there is a strong association b/w two objects.

$$cor(A \rightarrow B) = \frac{P(A, B)}{P(A) P(B)}$$

Incremental Algorithm:-

This Algorithm is used to find out the frequent itemset, when the db is dynamic.

Rules to find out Association rule :-

- * Find out the frequency of each item in new db, If it is already existed in old db.
- * The frequent item is not presented in T-old (old db) but present in new dataset, then ignore the item.
- * Find out the promoted border set.

Promoted border set :-

If an (frequent) item is not frequent in old db but after adding the new transaction to the db, the item is frequent, Then the set of frequent items is called as promoted border set.

TransID	Itemset	Support = 2
Told	1 B, J, P	
	2 B, m, P	
	3 m, P	
	4 C, D	
Tnew	5 C, J	Twhole = Told \cup Tnew
	6 B, C, (E) X	
	7 C, B	

Item frequency	
B	2 + 2 = 4
J	1 + 1 = 2
P	3
m	2
C	1 + 3 = 4
D	1 X

Tnew
B 4
J 2
P 3
m 2
C 4

Promoted border set

Promoted border set :-

$$\{\{J\}, \{C\} \{X\}, \{J, C\} \{J, X\} \{J, C, X\} \{C, X\}\}$$

After new set, D combinations are removed
as it doesn't satisfy support.

Promoted set
 $\{\{J\} \{C\} \{J, C\}\}$

Measuring the Quality of rules :-

- * The first support & confidence are basic methods to find out quality of Association techniques.

$$1. \quad S(x \rightarrow y) = \frac{S(x \cup y)}{n}$$

$$\alpha(x \rightarrow y) = \frac{S(x \cup y)}{S(x)}$$

drawback :-

we are not considering no. of times y had approached.

2. Lift (or) Interest

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

|| or correlation

drawback :-

symmetric

$$I(A \rightarrow B) = I(B \rightarrow A)$$

3. conviction

$$(A \Rightarrow B) = \frac{P(A)(P(\sim B))}{P(A, \sim B)}$$

$$A \Rightarrow B \equiv \sim (A \wedge \sim B)$$

If conviction = 1, then A & B are not related

4. chi-square method

$$\sum \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

If value is less than chisquare value it is strongly less it is not accepted.