

Unit-2

Association Rule Mining :-

Mining Frequent Patterns :-

Frequent patterns :-

The pattern that appears frequently in dataset.

- patterns include data items, data subsequence, data substructure.

Example :- Milk and Bread.

Market Basket Analysis (MBA) Association Rule Mining :-

This is a criteria of finding (or) determining the items (or) patterns which are likely to be purchased by the customer in single transaction.

It is mainly used for sellers.

Strategies used :-

1. Placing them together.
 2. Placing them at two different ends.
- This analysis will help ^{sellers} to plan their shelf space for increased sales.
 - Frequent patterns are represented by association rules

Example :- computer and Anti-virus.

(S) (S) (S) (S) (S) (S)

(A) (A) (A) (A) (A) (A)

Example

<u>Transaction ID</u>	<u>Items</u>
1	Bread, Milk.
2	Bread, Diaper, Beer, egg
3	Milk, Diaper, Beer, Cola
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Cola

Transactions $T = \{t_1, t_2, t_3, t_4, t_5\}$

Items (Total no. of items) $I = \{i_1, i_2, i_3, i_4, i_5\}$

Item set or set of items. To find a set

↳ Set of item sets of size 2 :- {Bread, milk} - 2 item set.
 ↳ Set of item sets of size 3 :- {Bread, milk, Beer} - 3 item set.

Support count : Frequently, an item set occurs.

$\sigma(\text{Bread, Milk, Beer}) = 1$ support

$\sigma(\text{Bread, Milk}) = 3$

Support :-

The ratio of Support Count, and no. of transactions.

$$\text{Support} = \frac{\text{Support Count}}{\text{Total no. of transactions}}$$

~~support count
of each pair is 3
and total no. of transaction is 5~~

Confidence :-

The ratio of Support of $A \cup B$ and Support of A .

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

Ex :- {Bread, Milk}

$$\text{Support (Bread)} = \frac{4}{5} = 0.8$$

$$\text{Support (Milk)} = \frac{4}{5} = 0.8$$

$$\text{Support (Bread, Milk)} = \frac{3}{5} = 0.6$$

$$\text{confidence} = \frac{\text{Support (Bread, Milk)}}{\text{Support (Bread)}}$$

$$= \frac{0.6}{0.8} = \frac{6}{8} = \frac{6}{10} = 0.75$$

Association and Correlation :-

Association :-

It is a way to identifying relation between various items.

Ex: It is used to determine sales of items that are frequently purchased together.

Correlation :-

It is Mathematical technique.

It shows how strongly pair of attributes are related together.

Ex: Tall people tend to have more weight.

Mining Methods

Mining Methods :-

There are two types of mining methods. They are:

1. Apriori Algorithm

2. FP Growth Algorithm.

① Apriori Algorithm :-

- It is given by R. Agrawal and R. Srikant.
- It shows how objects are associated with each other.

Objective : To generate an association.

Example :

Minimum Support = 50%

Threshold Confidence = 70%

T. ID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Item Set = (1, 2, 3, 4, 5)

Item set Support min Support.

1 2 $\frac{2}{4} = 50\%$

2 3 $\frac{3}{4} = 75\%$

3 3 $\frac{3}{4} = 75\%$

4 1 $\frac{1}{4} = 25\% (\times)$

5 3 $\frac{3}{4} = 75\%$

\therefore our minimum support is 50%
 so, in the above table which is less than 50%, it
 is not considered for next process.

New, Item Set = $(1, 2, 3, 5)$

- form pairs

$(1, 2), (1, 3), (1, 5), (2, 3), (2, 5), (3, 5)$.

<u>Item Set</u>	<u>Support</u>	<u>minimum Support</u>
$(1, 2)$	1	$\frac{1}{4} = 25\% (X)$
$(1, 3)$	2	$\frac{2}{4} = 50\%$
$(1, 5)$	1	$\frac{1}{4} = 25\% (X)$
$(2, 3)$	2	$\frac{2}{4} = 50\%$
$(2, 5)$	3	$\frac{3}{4} = 75\%$
$(3, 5)$	2	$\frac{2}{4} = 50\%$

$(1, 2)$ and $(1, 5)$ are eliminated because it is less
 than minimum support.

New, Item set = $(1, 3), (2, 3), (2, 5), (3, 5)$

- form triples.

$(1, 2, 3), (1, 2, 5), (1, 3, 5), (2, 3, 5)$

<u>Item Set</u>	<u>Support</u>	<u>Min Support</u>
$(1, 2, 3)$	1	$\frac{1}{4} = 25\% (X)$
$(1, 2, 5)$	1	$\frac{1}{4} = 25\% (X)$
$(1, 3, 5)$	1	$\frac{1}{4} = 25\% (X)$
$(2, 3, 5)$	2	$\frac{2}{4} = 50\%$

Item set = $(2, 3, 5)$

- Now, let's calculate support and confidence.

$$\text{confidence} = \frac{\text{support}(A \cup B)}{\text{support of } A}$$

Using $(2, 3, 5)$ we can generate association rules.

Rules	Support	confidence
$(2^1 3) \rightarrow 5$	2	100%
$(3^1 5) \rightarrow 2$	2	100%
$(2^1 5) \rightarrow 3$	2	66% (X)
$2 \rightarrow (3^1 5)$	2	66% (X)
$5 \rightarrow (2^1 3)$	2	66% (X)
$3 \rightarrow (2^1 5)$	2	66% (X)
$\text{rest} = \text{None}$	2	(None)

$$(2^1 3) \rightarrow 5 - \text{confidence} = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

$$= \frac{\text{support}((2^1 3) \cup 5)}{\text{support}(2^1 3)}$$

$$= \frac{2}{2} = 1 = 100\%$$

$$(3^1 5) \rightarrow 2 - \text{confidence} = \frac{\text{support}((3^1 5) \cup 2)}{\text{support}(3^1 5)}$$

$$= \frac{2}{2} = 1 = 100\%$$

$$(2^1 5) \rightarrow 3 - \text{confidence} = \frac{\text{support}((2^1 5) \cup 3)}{\text{support}(2^1 5)}$$

$$= \frac{2}{3} = 66\%$$

$$2 \rightarrow (3 \wedge 5) \rightarrow \text{confidence} = \frac{\text{Support}(2 \cup (3 \wedge 5))}{\text{Support}(2)}$$

$$= \frac{2}{3} = 66\%$$

$$5 \rightarrow (2 \wedge 3) \rightarrow \text{confidence} = \frac{\text{Support}(5 \cup (2 \wedge 3))}{\text{Support}(5)}$$

$$= \frac{2}{3} = 66\%$$

$$3 \rightarrow (2 \wedge 5) \rightarrow \text{confidence} = \frac{\text{Support}(3 \cup (2 \wedge 5))}{\text{Support}(3)}$$

$$\text{Support} = \frac{2}{3} = 66\%$$

$(2 \wedge 3) \rightarrow 5$, $(3 \wedge 5) \rightarrow 2$ are association rules.

② FP Growth Algorithm

- FP Stands for Frequent patterns.
- It is an efficient and scalable method for mining the complete set of FP by using a tree structure.
- Tree structure stores information about FP called FP tree.

Example: (B, E, A, C, D, S) - applying FP tree

Minimum Support = 30% - frequent items

item based support bit mask

A, B, C, D

B, C, D, E

b1, b2, b3, b4

E, B, C, D, S

B, C, D, E, S

D, E, A, S

A, B, C, D

A, B, S

C, D, E

Trans. Id	Items
1	E, A, D, B
2	D, A, E, C, B
3	C, A, B, E
4	B, A, D
5	D
6	D, B
7	A, D, E
8	B, C

- List out the priorities.

Itemset = (A, B, C, D, E)

Itemset	Frequency	Priority
A	5	3
B	6	1
C	3	5
D	6	2
E	4	4

Order of priority - (B, D, A, E, C)

- Order items according to priority.

Trans. Id	Items	Ordered Items
1	E, A, D, B	B, D, A, E
2	D, A, E, C, B	B, D, A, E, C
3	C, A, B, E	B, A, E, C
4	B, A, D	B, D, A

5

D

6

D, B

四

干

A, D, E

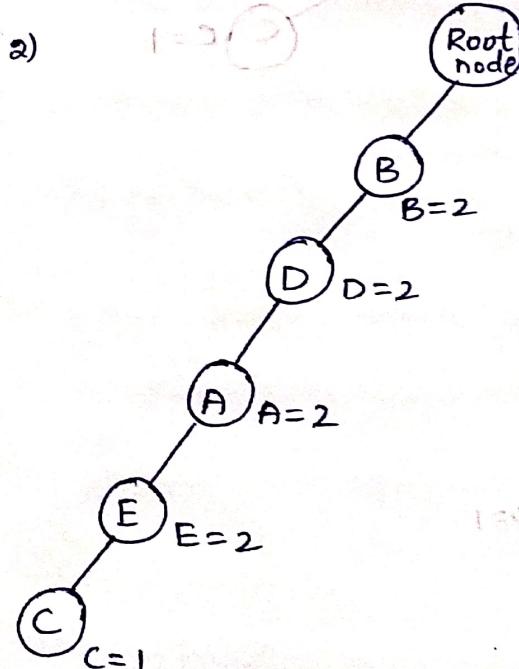
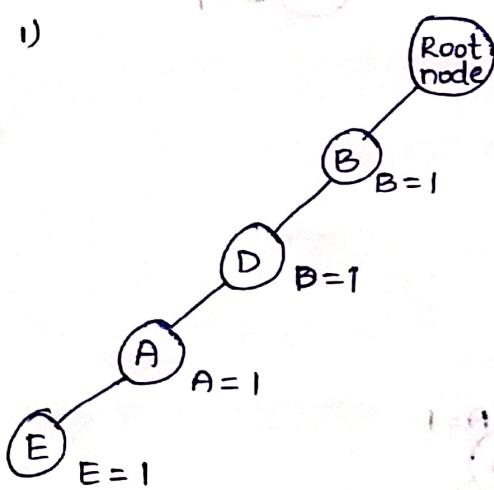
D,A,E

8

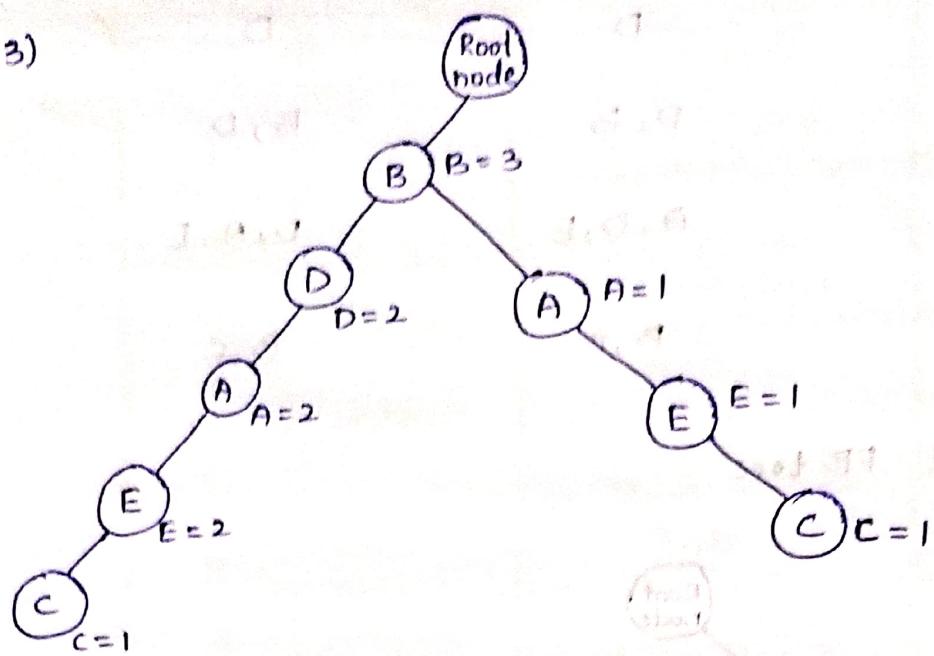
B,C

B, C

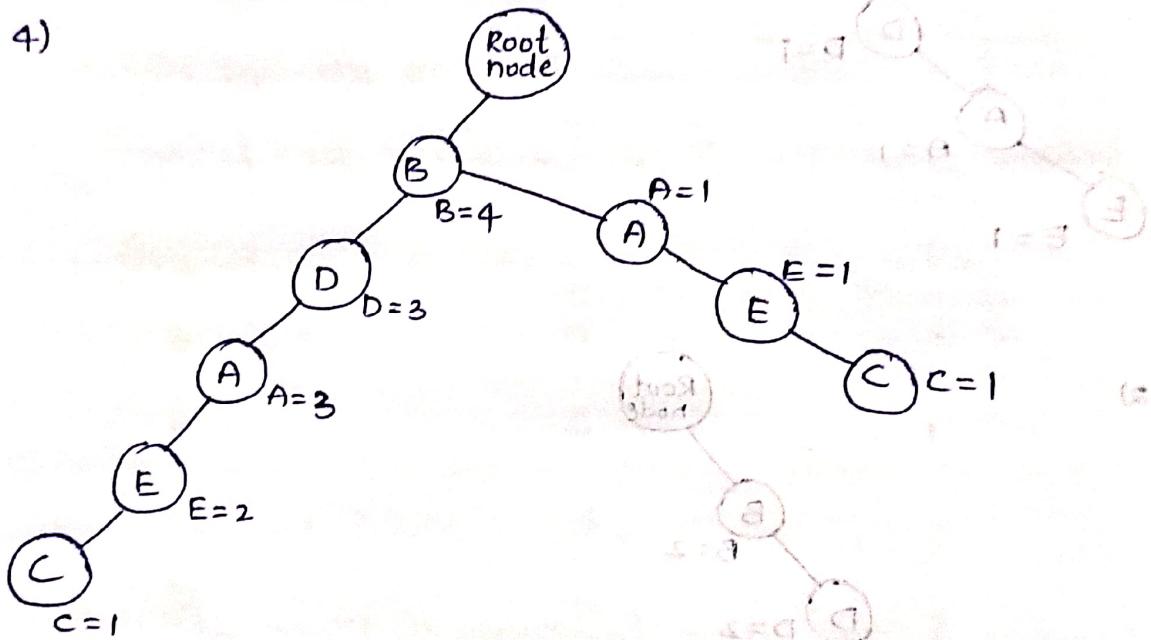
- Construct FP tree.



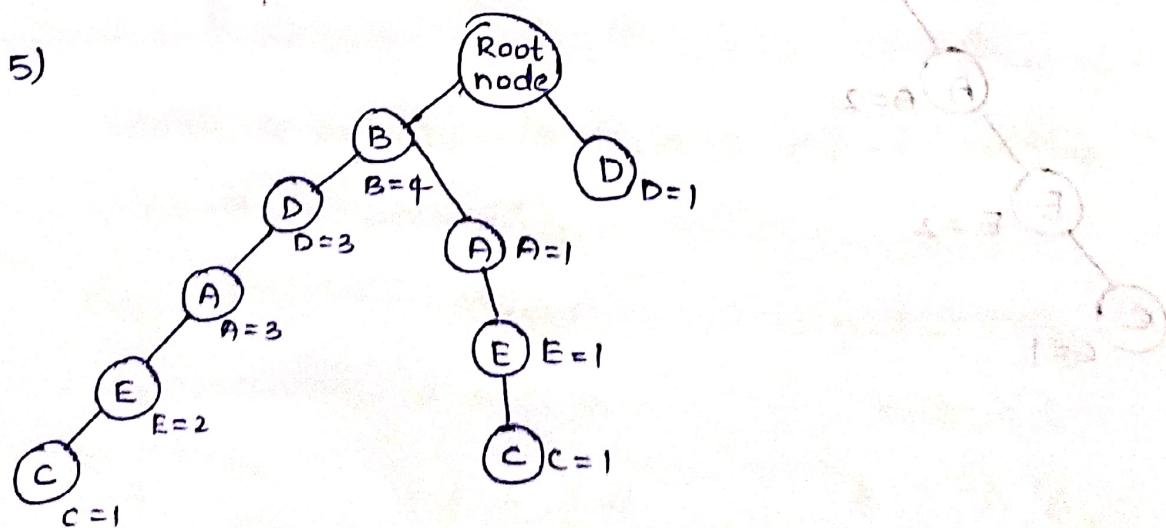
3)

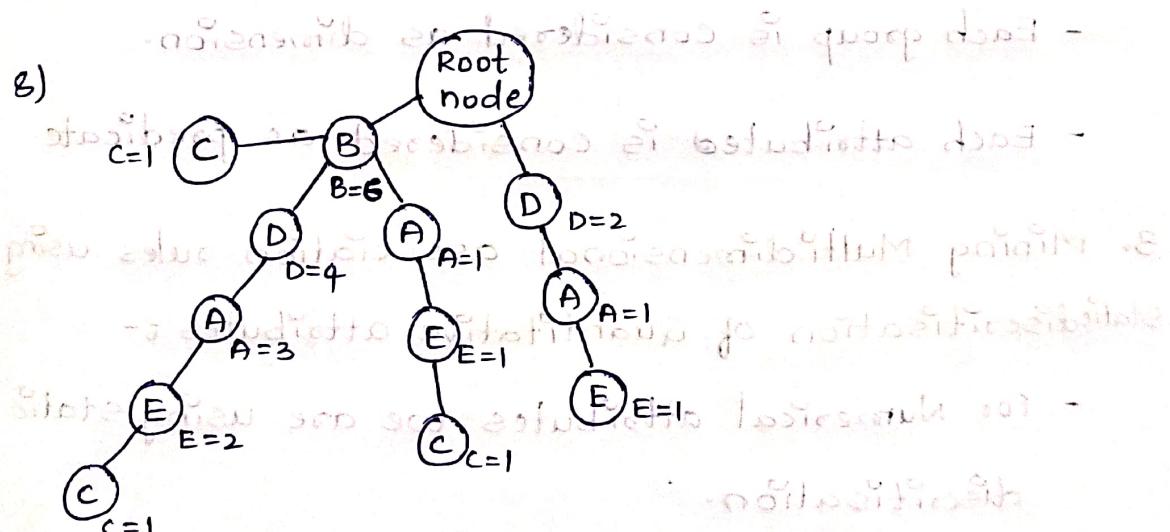
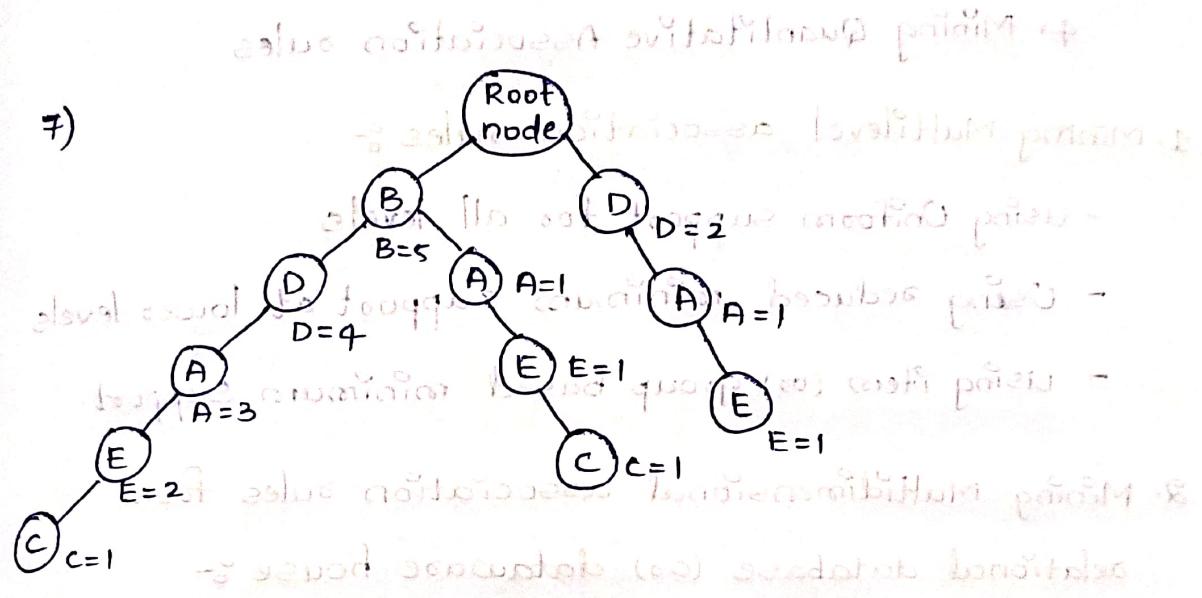
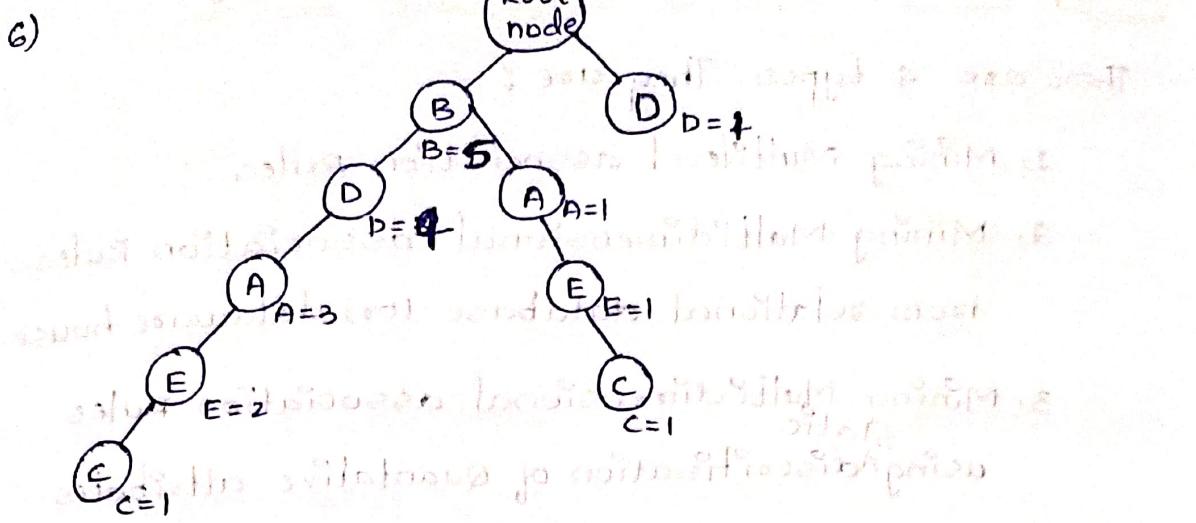


4)



5)





at the final step, we can see that the number of times each value is printed =

$B = 6 \text{ times}$

$D = 6 \text{ times}$

$A = 6 \text{ times}$

Final step $E = 4 \text{ times}$

$C = 3 \text{ times}$.

Mining Various kinds of Association Rules :-

These are 4 types. They are :-

1. Mining Multilevel association Rules.
2. Mining Multidimensional association Rules, from relational database (or) dataware house.
3. Mining Multi^{static}dimensional association Rules using discritisation of Quantitative attributes.
4. Mining Quantitative Association rules.

1. Mining Multilevel association rules :-

- Using Uniform support for all levels
- Using reduced minimum support at lower levels
- Using item (or) group based minimum support.

2. Mining Multidimensional association rules from relational database (or) dataware house :-

- Each group is considered as dimension.
- Each attribute is considered as predicate.

3. Mining Multi^{static}dimensional association rules using static discritisation of Quantitative attributes :-

- For Numerical attributes we are using static discritisation.
- static discritisation means dividing into intervals.

4. Mining Quantitative Association rules :-

- For Numerical attributes we are using dynamic discritisation
- It is more efficient.

Correlation Analysis

- Correlation Analysis is used to measure the relationship between two variables.
- One variable increases, then other variable also increases is called positive correlation.
- One variable ~~decreases~~ increases, then other variable decreases is called negative correlation.

$$\gamma_{A,B} = \frac{\sum (A - A')(B - B')}{(n-1) \sigma_A \sigma_B}$$

Where,

$\gamma_{A,B}$ = Karl Pearson Correlation Coefficient

A', B' = mean of A and B

σ_A, σ_B = standard deviation of A and B.

n = no. of tuples in d.B.

γ has 3 values (0, ±1, +1)

$\gamma \rightarrow +1 \Rightarrow$ perfect positive correlation

$\gamma \rightarrow -1 \Rightarrow$ perfect negative correlation

$\gamma \rightarrow 0 \Rightarrow$ No correlation (no dependence).

Example :

A	B
20	8
12	34
9	4

$$\gamma_{A,B} = \frac{\sum (A - A')(B - B')}{(n-1) \sigma_A \sigma_B}$$

$$n=3$$

A' = mean of A

$$A' = \frac{20+12+9}{3} = \frac{41}{3} = 13.66$$

B' = Mean of B

$$B' = \frac{8+34+4}{3} = \frac{46}{3} = 15.33$$

$$\sigma_A = \sqrt{\frac{\sum (A - A')^2}{(n-1)}}$$

$$= \sqrt{\frac{(20-13.66)^2 + (12-13.66)^2 + (9-13.66)^2}{(3-1)}}$$

$$= \sqrt{\frac{(6.34)^2 + (-1.66)^2 + (-4.66)^2}{2}}$$

$$= \sqrt{\frac{40.19 + 2.75 + 21.41}{2}}$$

$$= \sqrt{\frac{64.65}{2}}$$

$$= \sqrt{32.32}$$

$$= 5.68$$

$$\sigma_B = \sqrt{\frac{\sum (B - B')^2}{(n-1)}}$$

$$\sigma_B = \sqrt{\frac{(8-15.33)^2 + (34-15.33)^2 + (4-15.33)^2}{(3-1)}}$$

$$= \sqrt{\frac{(-7.33)^2 + (18.67)^2 + (-11.33)^2}{2}}$$

$$\text{Ansatz} = \sigma_B \sqrt{\frac{530.64 + 348.56 + 128.36}{2}}$$

$$= \sqrt{\frac{530.64}{2}} = \sqrt{265.32} = 16.128$$

$$\gamma_{A,B} = \frac{\sum [(A-A')](B-B')]}{(n-1) \sigma_A \sigma_B}$$

$$= \frac{(20-13.66)(8-15.33) + (12-13.66)(34-15.33) + (9-13.66)(4-15.33)}{16.128 \times 5.68 \times 16.128}$$

$$= \frac{6.34 \times (-7.33) + (-1.66) \times 18.67 + (-4.66) \times (-11.33)}{2 \times 5.68 \times 16.128}$$

$$= \frac{-46.47 + (-30.99) + 52.79}{184.94}$$

$$\text{Corr} = \frac{-24.67}{184.94} = -0.132$$

≤ -1 for negative correlation.

i.e. Negative Correlation.

Constraint Based Association Mining :-

- Constraint means condition.
- association rules are generated based on conditions.

Types of constraints :-

1. Knowledge type :-

Specifies the knowledge you want to do mining.
- association, correlation, Regression, etc.

2. Data Constraints :-

Specifies the data on which you generate association rules.
- only task relevant data.

3. dimension/ level constraints :-

Specifies the dimension (or) level.

- concept hierarchy.

4. Interestingness Constraints :-

Support, Confidence are used to identify.

5. Rule Constraints :-

Specifies the form of rules to be mined.

1. Metarules Guided mining

2. Constraint pushing.

Graph Patterns Mining :-

- set of tools and techniques are used to mine frequent subgraphs.
- It is used to analyze the properties of real world graphs.

- It is used to analyse the how structure of graph will effect the rules.
- There are two methods. They are
 - Apriori based Approaches
 - Pattern growth Approaches.

Algorithms used :-

- Gspan
- closed Graph.

Application :-

- In XML structures
- Anomaly detection
- Network analysis
- Control flow Analysis
- Biological structures, etc.

Sequential Pattern Mining :-

- SPM Stands for Sequential pattern mining.
- Sequence = ~~order~~ set of ordered events.
Ex :- $S = \{e_1, e_2, e_3, e_4, e_5\}$.
- SPM \rightarrow Process of finding frequent subsequences from a set of Sequence.
- Sequences are represented by $< >$

Normal Transactional data

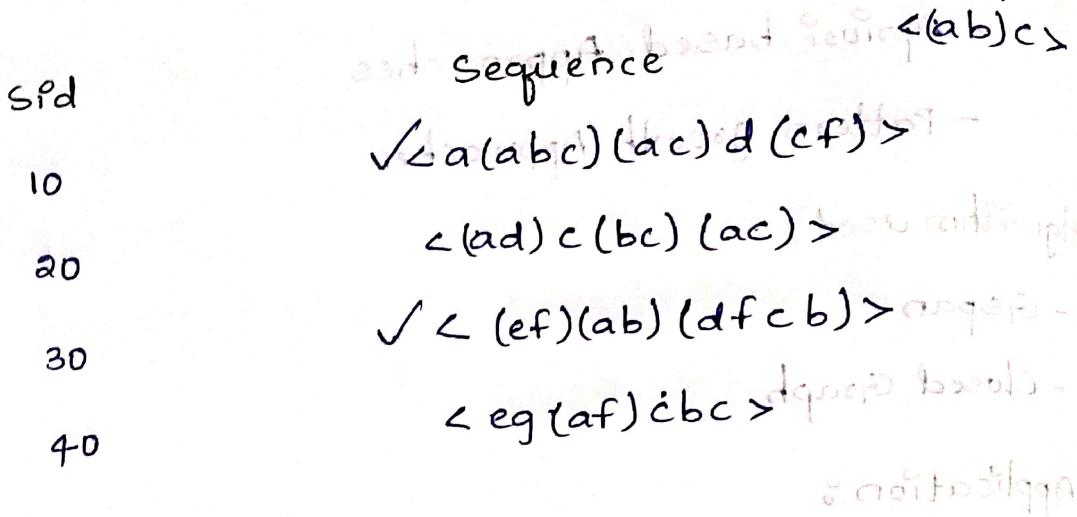
CID	TID	Transactions
1	100	a,b,e,d
3	111	a,f,d,e
1	122	d,e,p
3	133	b,f,s,a

Sequential data.

CID	Sequences
1	$\langle (abed), (d,ep) \rangle$
3	$\langle (afde), (bfsa) \rangle$

challenges in SPM :-

- finding all subsequences. min.sup = 2.



$\langle ab \rangle c$ is present in 1st sequence

$\langle ab \rangle c$ is present in 3rd sequence

$\therefore \text{min sup} = 2$.

so it is frequent

so minimum support is 2

minimum support is 2

above frequent for abcde = 3 times

$\{2, 3, 4, 5, 6, 7, 8\} = 2^7 = 128$

corresponding frequent pattern for abcde \leftarrow min

frequent for abcde

\Leftrightarrow pd. between abcde

abcde

abcde frequent

abcde frequent

abcde frequent

$\langle (abc) \rangle \langle (bcd) \rangle$

bacd frequent

$\langle (bcd) \rangle \langle (abc) \rangle$

cabd frequent

$\langle (abc) \rangle \langle (cde) \rangle$

acbed frequent

$\langle (cde) \rangle \langle (abc) \rangle$

cadbe frequent