

## Unit 4:

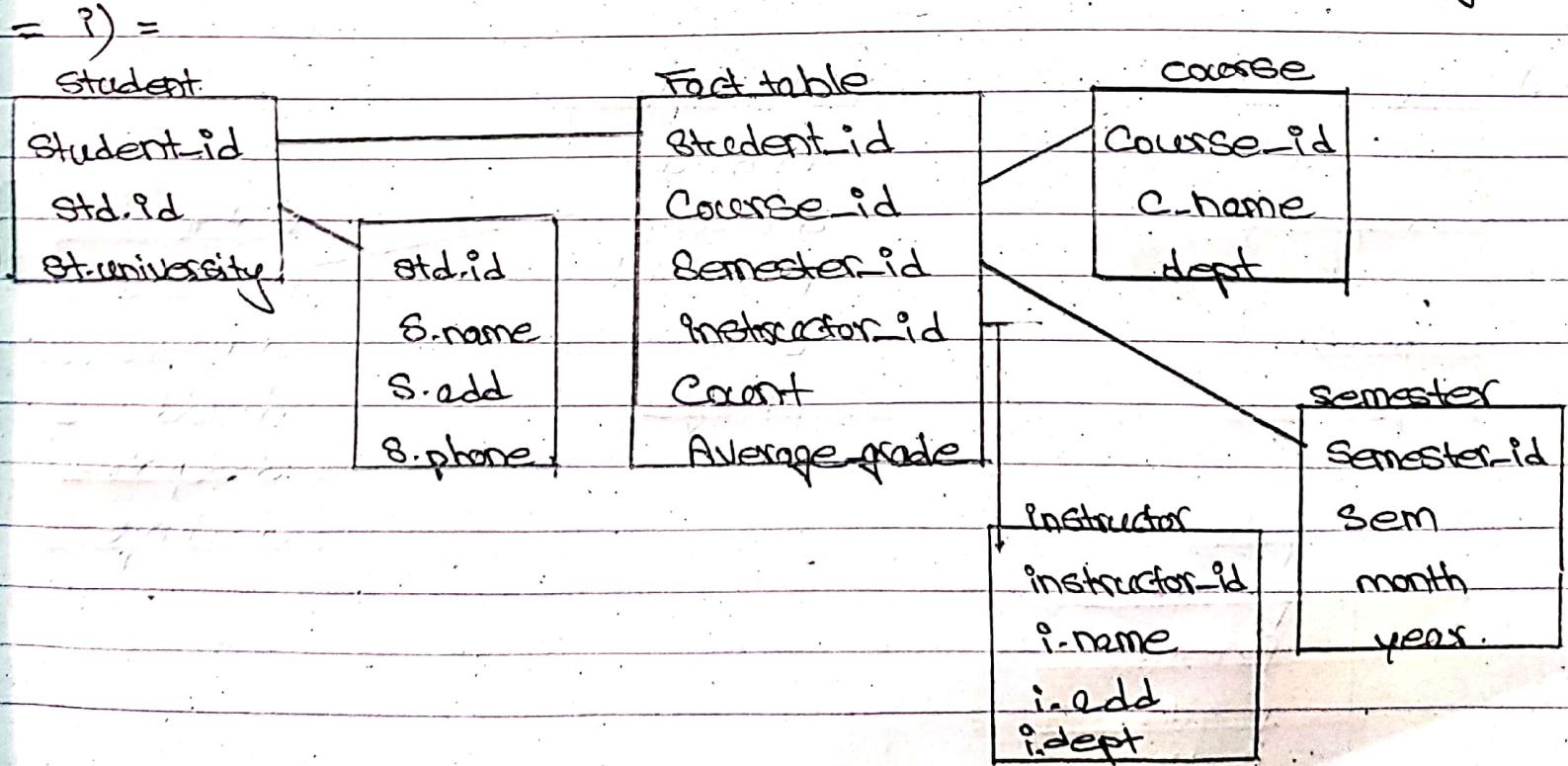
### \* Multidimensional data models and different OLAP operations:

a) Suppose that a data warehouse for big university consists of the following four dimensions:

Student, Course, Semester, an instructor and two measures: Count & average grade.

When at the lowest conceptual level (e.g. for a given Student, Course, Semester & instructor combination), the avg-grade measure stores the actual course grade of the student. At higher conceptual levels, avg-grade stores the average grade for the given combination.

i) Draw a snowflake schema diagram for the data warehouse.  
ii) Starting with the base cube (Student, Course, Semester, Instructor) what specific OLAP operations (e.g. rollup from Semester to year) should one perform in order to list the average grade of CS courses for each Big University Std



- b) i) Roll up on course from course-id to dept  
 ii) Roll up on student from student-id to university.  
 iii) Dice on course, student.  
 (course = "CS" AND University = "Big University")

- b) Suppose that a data warehouse consists of the three dimensions: time, doctor & patient and the two measures Count & charge where charge is the fee that a doctor charges a patient for visit.
- Draw a schema diagram for the above datawarehouse using star, snowflake & fact constellation Schemes.
  - Starting with the base cuboid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2004?

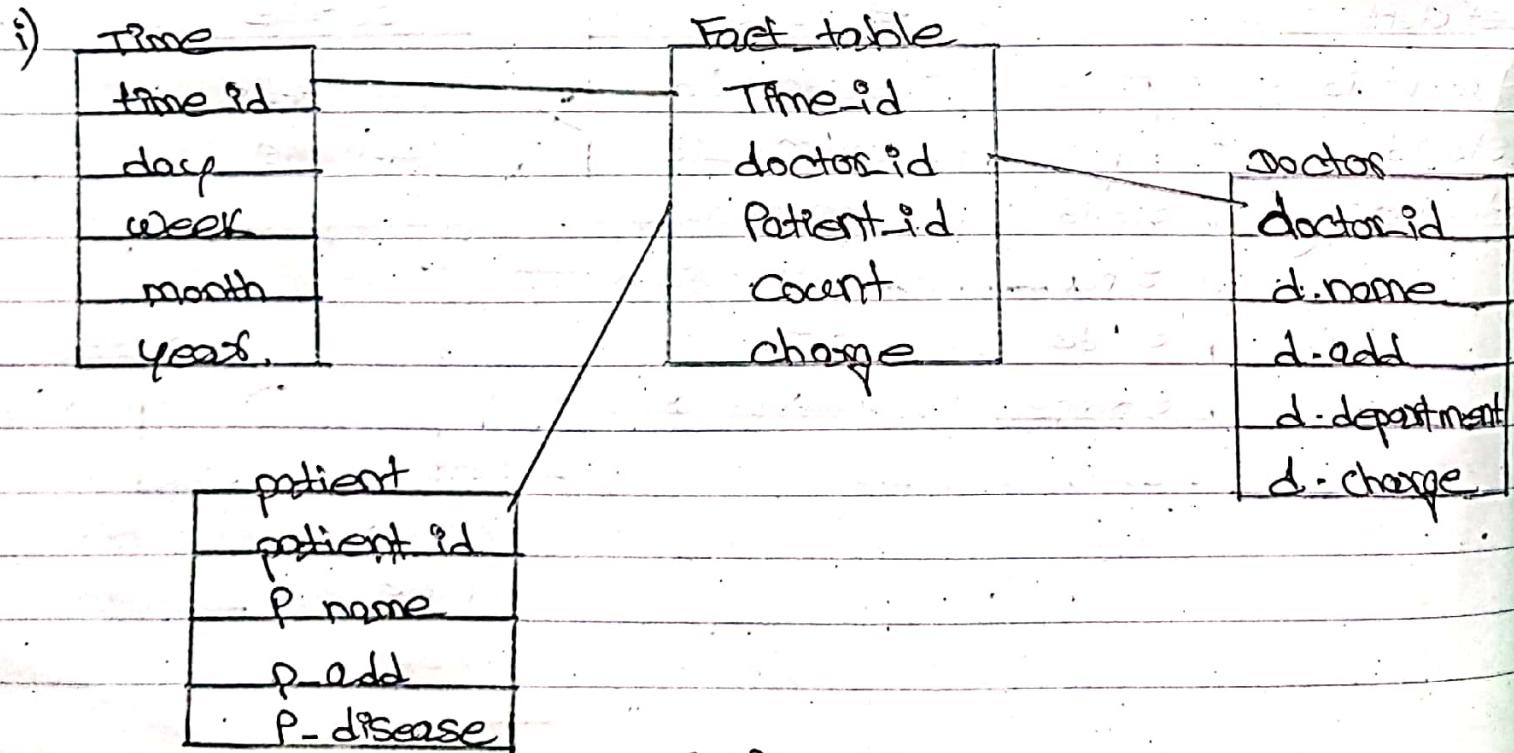
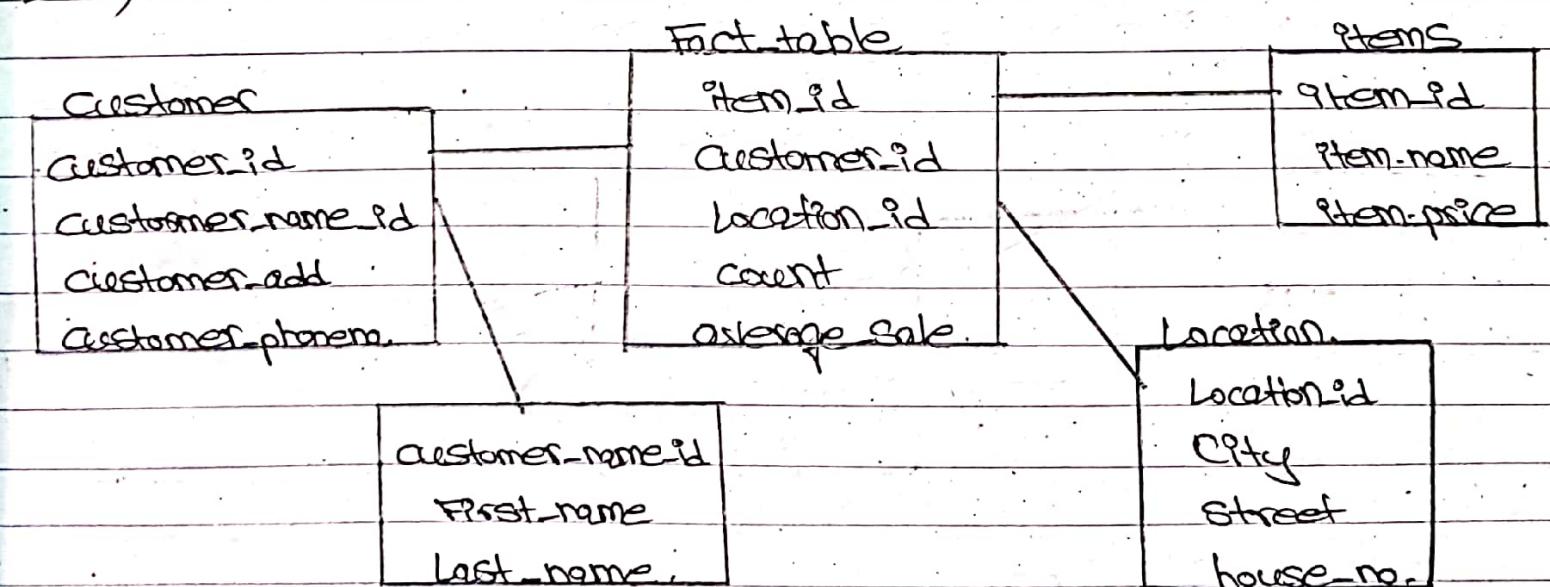


Fig:- Star Schema.

- b) = i) Roll up on time from time\_id to year.  
ii) Roll up on doctor from doctor\_id to doctor\_name.  
iii) Roll up on patient from patient\_id to patient\_name  
iv) Dice on doctor, patient, year.

- c) Suppose that a datawarehouse for one of the shopping center consists of the following three dimensions: item, customer, location & two measures count and average sale.
- Draw a snowflake schema diagram for the datawarehouse
  - Starting with the base cuboid [item, customer, location], what specific OLAP operations should one perform in order to list the average sales of items for each customer of the shopping center with reference to particular location.

= a) =



- b) = Rollup on location from location\_id to house\_no.  
Rollup on customer from customer\_id to customer\_name\_id.  
Dice on customer, items, house\_no.

## \* Data Cube Computation:

- It is a process of calculating the no. of cuboids provided by no. of dimensions.

Two approaches:

i) Dimension without hierarchy.

ii) Dimension with hierarchy.

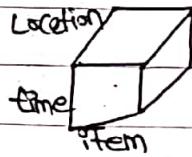
### i) Dimension without hierarchy:

- If there were no hierarchies associated with each dimension then total number of cuboids for  $n$  - dimensional data cube, is  $2^n$ .

- However, many dimensions do have hierarchies. For example, the dimension time is usually not explored at only one conceptual level, such as year but rather at multiple conceptual levels, such as in hierarchy "day < month < quarter < year".

- For e.g.:-

$$n = 3 \text{ (Item, Location, time)}$$
$$\text{cuboids} = 2^n = 8 \text{ cuboids}$$



Apex [All, ()<sub>q</sub>]

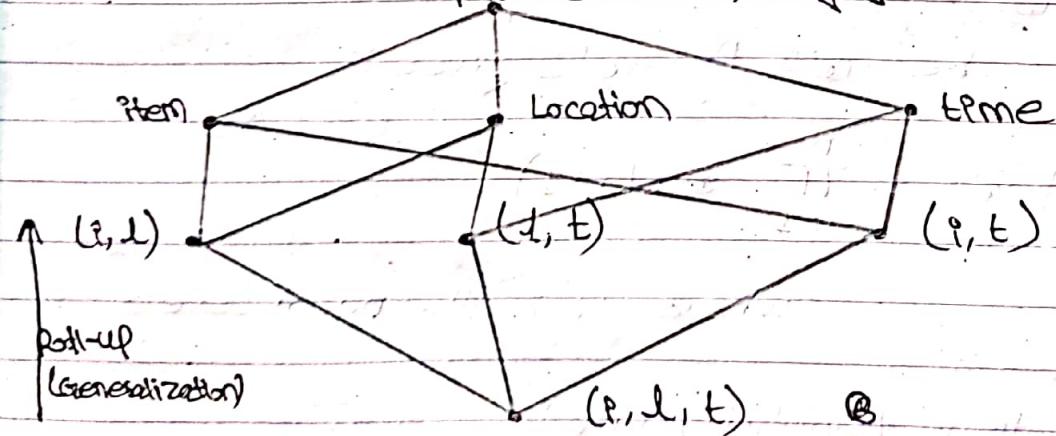
0-D

roll-down (Specialization)

1-D

- 2-D

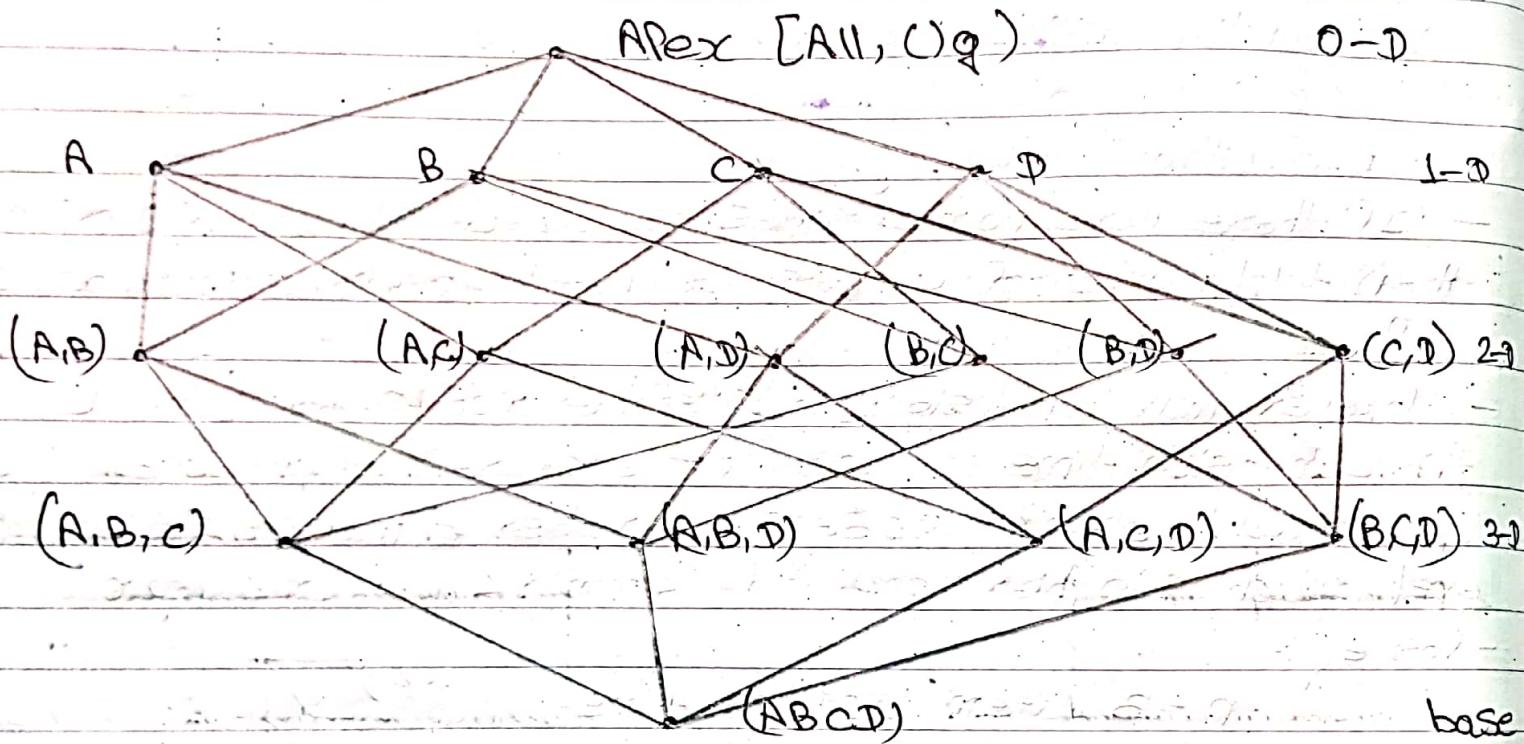
base:



\* Know about cuboids for n-dimension A, B, C, d - (dimension without hierarchy)

$$n = 4$$

$$\text{Cuboids} = 2^n = 2^4 = 16$$



### ii) Dimension with hierarchy:

- n-dimensional cube with L-levels gives:

$$\text{no. of cuboids} = \prod_{i=1}^n (L_i + 1)$$

where,  $L_i$  is the no. of levels.

For eg:-  $n = 3$  (Item, Location, time)

where, Item & Locations has Level 1 & Time has Level - 3 i.e. days  $\rightarrow$  months  $\rightarrow$  years

$L_1$        $L_2$        $L_3$

$$\text{Nb. of cuboids} = 2^d \times 2^t \times 4$$

i      d      t

Nb. of cuboids = 16

### \* Curse of dimensionality:

- space complexity  $\uparrow$  with respect to dimension  $\uparrow$
- dealing with the problem (partial materialization)

### \* Bitmap Index:

- There are major two categories of indexing:-

i) Dense

ii) Sparse.

Keys	B <sub>1</sub>
Keys	B <sub>2</sub>
:	:
Keys	B <sub>n</sub>

Fig:- Dense-indexing.

Keys	B <sub>1</sub> - B <sub>4</sub>
Keys	B <sub>5</sub> - B <sub>8</sub>
:	:
Keys	B <sub>n-4</sub> - B <sub>n</sub>

Fig:- Multi-level indexing  
sparse - Indexing.

- It is used to reduce the no. of disk access.

For e.g:-

i) Block = 512 B

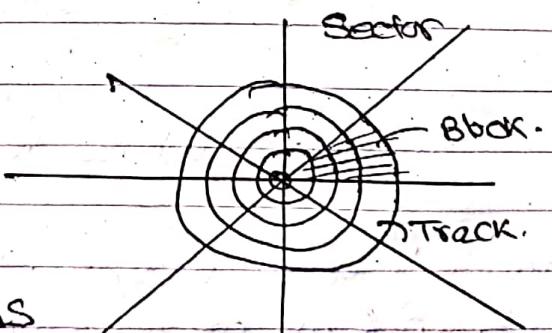
ii) no. of records that can be

$$\text{Kept in a block} = \frac{512}{128} = 4$$

iii) no. of blocks required for 100 records

$$= \frac{1}{4} \times 100 = 25$$

i.e. If no indexing is carried out the total block access is 25.



e.g:-

Table 1:

ID	Name	Sex
1	A	M
2	B	F
3	A	M
4	C	F

ID	Name	M	F
1	A	1	0
2	B	0	1
3	A	1	0
4	C	0	1

Table 2:

T.ID	I.ID	Item bought
1	1	Y
2	2	Y
3	1	M
4	3	Y

T.ID	I.ID	Item bought
1	1	Y
2	2	Y
3	1	1
4	3	0

① Publish those names who is female & have bought an item?

= We need to use ~~Bitmap~~ Bitmap index join.

so,

Female = (B, C) having id (2, 4)

bought-item = (1, 2, 3)

Bitmap index join.

So, Those whose is female & have bought an item is B.

② Publish those names who is female or have bought an item?

= being & an or operation include all names who are female or have bought an item.

i.e., (A, B, C, D).

## Unit 6:

### \* Association Rule Mining:

#### ① Frequent pattern: (Apriori Algorithm)

A pattern (a set of items, subsequences, substrings, etc.) that occurs frequently in a data set.

for e.g:-

Transaction_id	Items bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F

i) Find all the rules  $X \rightarrow Y$  with minimum support & confidence.

ii)  $\text{min-sup} = 50\%$  and  $\text{min-conf} = 50\%$ .

Sol

Support Count = min-sup of total transaction

$$= 50\% \text{ of } 5$$

$$= \frac{50}{100} \times 5$$

$$= 2.5 \text{ let it be } (3)$$

now,

$C_1$ =	Item_Set	Support
A	3 ✓	
B	3 ✓	
C	2 ✗	
D	4 ✓	
E	3 ✓	
F	2 ✗	

$L_1$ =	Item_Set	Support
A	3	
B	3	
D	4	
E	3	

$C_2$ = Item Set	Support	$L_2$ = Item Set	Support
AB	1 X	AD	3
AD	3 ✓		
AE	1 X		
BD	2 X		
BE	2 ✓		
DE	2 ✓		

so, The Association rules we get are,

$$A \rightarrow D$$

$$D \rightarrow A$$

i)  $A \rightarrow D$

Support = probability that a transaction contain AUB (for  $A \rightarrow B$ )

$$= \frac{3}{5} \times 100\% = 60\%$$

confidence =

$$\frac{\text{AUB}}{\text{occurrence of A}} \times 100\%$$

$$= \frac{3}{3} \times 100$$

$$= 100\%$$

ii)  $D \rightarrow A$

$$\text{Support} = \frac{3}{5} \times 100\% = 60\%$$

$$\text{confidence} = \frac{3}{4} \times 100\% = 75\%$$

## Why pruning:

= Combination of not frequent itemset will never generate frequent item set. So neglect itemset having support less than the min-support.

Also, Combination of frequent itemset may or may not generate another frequent itemset. So, we need to check every possible rules.

for eg:-

TID	items
1	A, B, C, D, E, F
2	B, C, D, E, F, G
3	A, D, E, H
4	A, D, F, I, J
5	B, D, E, K

Assume support = 60 %. Generate frequent item sets using Apriori algorithm. Also generate the Association rules with confidence = 80 %.

$$= \frac{60}{100} \times 5$$

$$\text{Support\_Count} = \frac{60 \times 5}{100}$$

$$= 3$$

Also, we have confidence = 80 %.

$C_1 = \text{Itemset}$	$\text{Support}$
A	3
B	3
C	2 $\times$
D	5
E	4
F	3
G	1 $\times$
H	1 $\times$
I	1 $\times$
J	1 $\times$
K	1 $\times$

$L_1 = \text{Itemset}$	$\text{Support}$
A	3
B	3
D	5
E	4
F	3

$C_2 = \text{Itemset}$	$\text{Support}$
SA, B $\exists$	1 $\times$
SA, D $\exists$	3
SA, E $\exists$	2 $\times$
SA, F $\exists$	2 $\times$
SB, D $\exists$	3
SB, E $\exists$	3
SB, F $\exists$	2 $\times$
SD, E $\exists$	4
SD, F $\exists$	3
SE, F $\exists$	2 $\times$

$L_2 = \text{Itemset}$	$\text{Support}$
SA, D $\exists$	3
SB, D $\exists$	3
SB, E $\exists$	3
SD, E $\exists$	4
SD, F $\exists$	3

$C_3 = \text{Itemset}$	Support	$L_3 =$	Itemset	Support
S, B, D, E, F	3		S, B, D, E, F	3

Association rules:

$A \rightarrow D$	confidence = $\frac{3}{3} = 100\%$	✓
$D \rightarrow A$	confidence = $\frac{3}{5} = 60\%$	✗
$B \rightarrow D$	confidence = $\frac{3}{5} = 60\%$	✓
$D \rightarrow B$	confidence = $\frac{3}{5} = 60\%$	✗
$B \rightarrow E$	confidence = $\frac{3}{3} = 100\%$	✓
$E \rightarrow B$	confidence = $\frac{3}{3} = 100\%$	✗
$D \rightarrow E$	confidence = $\frac{4}{4} = 100\%$	✓
$E \rightarrow D$	confidence = $\frac{4}{5} = 80\%$	✓
$D \rightarrow F$	confidence = $\frac{3}{4} = 75\%$	✗
$F \rightarrow D$	confidence = $\frac{3}{4} = 75\%$	✓
$B \rightarrow D, E$	confidence = $\frac{3}{3} = 100\%$	✗
$D \rightarrow B, E$	confidence = $\frac{3}{5} = 60\%$	✗
$E \rightarrow B, D$	confidence = $\frac{3}{5} = 60\%$	✗
$D, E \rightarrow B$	confidence = $\frac{3}{4} = 75\%$	✗
$B, E \rightarrow D$	confidence = $\frac{3}{3} = 100\%$	✓
$B, D \rightarrow E$	confidence = $\frac{3}{3} = 100\%$	✓

### \* F-P Growth (Frequent Pattern):

Q.no.1.	TID	Items
	1	A, B, C, D, E, F
	2	B, C, D, E, F, G
	3	A, D, E, H
	4	A, D, F, I, J
	5	B, D, E, K

Support\_Count = 3

=Sol

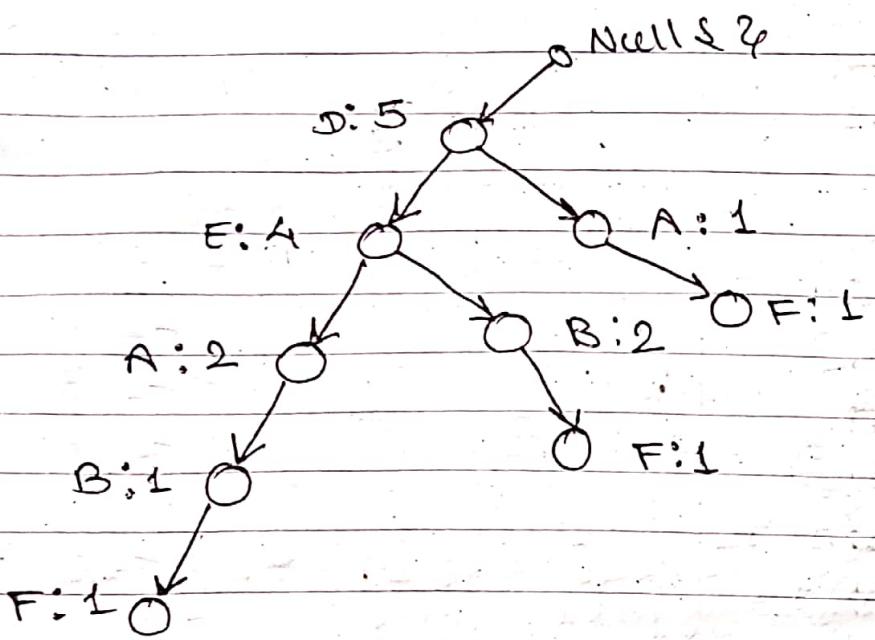
$C_1$	Itemset	Support
A	3 ✓	
B	3 ✓	
C	2 ✗	
D	5 ✓	
E	4 ✓	
F	3 ✓	
G	1 ✗	
H	1 ✗	
I	1 ✗	
J	1 ✗	
K	1 ✗	

Sorting in  
descending  
order.

$L_1$	Itemset	Support
D	5	
E	4	
A	3	
B	3	
F	3	

now, for each transaction, the respective ordered-item set is:

TID	Items	ordered-item set
1	A,B,C,D,E,F	D,E,A,B,F
2	B,C,D,E,F,G	D,E,B,F
3	A,D,E,H	D,E,A
4	A,D,F,I,J	D,A,F
5	B,D,E,K	D,E,B



Items	Conditional Pattern Base	conditional FP-tree	Frequent Patterns
F	$S(D, E, A; B: 1), (D, E, B: 1),$ $C, D, A, F: 1)$	$S D: 3\gamma$	$S D F: 3\gamma$
B	$S(D, E, A: 1), (D, E: 2)$	$S (DE: 3)\gamma$	$S D B: 3\gamma, S E, B: 3\gamma$
A	$S(D, E: 2), (D: 1)$	$S (D: 3)\gamma$	$S D, E, B: 3\gamma$ $S D, A: 3\gamma$
E	$S(D: 4)\gamma$	$S (D: 5)\gamma$	$S D, E: 5\gamma$
D	Null		

Q.No.2: TID

Items

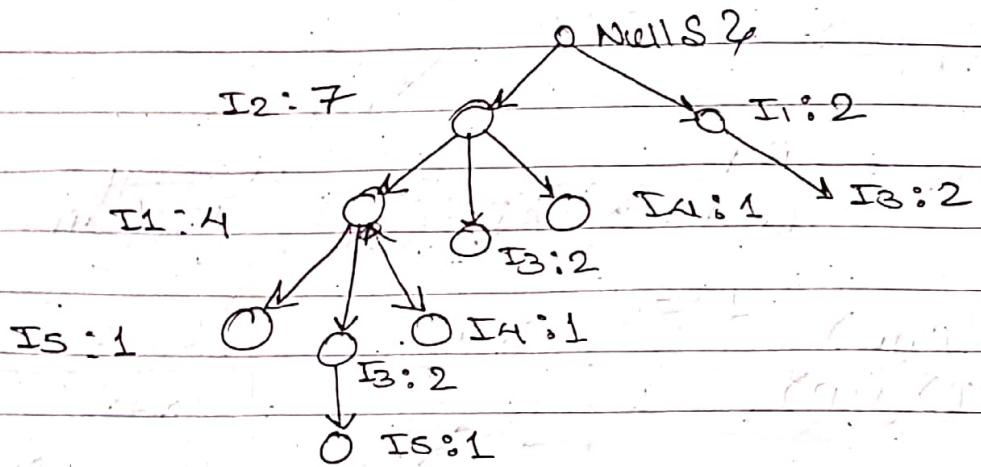
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I1, I2, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I2, I3

↓ L1

Itemset	Support		Itemset	Support
I1	6	sorting in descending order.	I2	7
I2	7		I1	6
I3	6		I3	6
I4	2		I4	2
I5	2		I5	2

now, for each transaction the respective item-set is:

TID	Items	ordered Item Set
T1	I1, I2, I5	I2, I1, I5
T2	I2, I4	I2, I4
T3	I2, I3	I2, I3
T4	I1, I2, I4	I2, I1, I4
T5	I1, I3	I1, I3
T6	I2, I3	I2, I3
T7	I1, I3	I1, I3
T8	I1, I2, I3, I5	I2, I1, I3, I5
T9	I1, I2, I3	I2, I1, I3



Items	conditional pattern	conditional FP-tree	frequent pattern
I5	$S(I_2, I_1:1), (I_2, I_1, I_3:1)$	$S(I_2, I_1:2)$	$(I_2, I_5:2)$
I4	$S(I_2, I_1, I_3:1), (I_2:1)$	$S(I_2:2)$	$S(I_2, I_4:2)$
I3	$S(I_2, I_1:2), (I_2:2), (I_1:2)$	$S(I_2:4, I_1:2), (I_1:2)$ <small><math>\text{#}_2 &gt; 4 \text{ branch}</math></small>	$S(I_2, I_3:4)$ $(I_1, I_3:4)$ $(I_2, I_1, I_3:2)$
I1	$S(I_2:4)$	$S(I_2:4)$	$S(I_1, I_2:4)$

## \* Correlation Analysis:

$$\textcircled{1} \text{ All-confidence} = \frac{P(A \cap B)}{P(A) \cdot P(B)} \text{ or } \frac{\text{Support}(A \cap B)}{\max(\text{Support}(A), \text{Support}(B))}$$

$$\textcircled{2} \text{ Lift} = \frac{P(A \cap B)}{P(A) \cdot P(B)}$$

Note:

$$\textcircled{3} \text{ Cosine} = \frac{P(A \cap B)}{\sqrt{P(A) \cdot P(B)}}$$

$$\textcircled{4} \chi^2 = \sum \frac{(O - E)^2}{E} \quad \text{where, } O = \text{observed value}, \quad E = \text{Expected value}$$

$\chi^2$	Lift
tve < 1 -ve > 1 $A \cap B \uparrow$	cosine & All-Confidence

tve < 1 $A \cap B \uparrow$	-ve > 1 $A \cap B \downarrow$
-----------------------------------	-------------------------------------

Lift
------

tve > 1, $A \cap B \uparrow$	-ve < 1 $A \cap B \downarrow$
------------------------------------	-------------------------------------

correlation

e.g:- Suppose we are interested in analysing transactions at XYZ shopping center w.r.t the purchase of computer games and videos. Let game refer to the transaction containing computer games and video refers to the transaction containing videos of the 10,000 transaction analyzed, the data shows that 6,000 of the transactions included computer games, 7500 included videos & 4000 included both computer games & videos.

Let  $\text{Support}_{\min} = 30\%$ ,  $\text{Confidence}_{\min} = 60\%$ .

• buys(X, games)  $\Rightarrow$  buys(X, videos) predicate?

=So Total transactions = 10,000

computer game Transaction = 6,000

videos Transaction = 7,500

Both computer games & videos = 4,000

- ① Support of buys (X, games) =  $\frac{6000}{10000} \times 100\% = 60\%$
- ② Support of buys (X, videos) =  $\frac{7500}{10000} \times 100\% = 75\%$
- ③ Support of buys (X, games & videos) =  $\frac{4000}{10000} \times 100\% = 40\%$

now, confidence of buys (X, games)  $\rightarrow$  buys (X, videos)

$$= \frac{\text{A} \cup \text{B}}{\text{occurrence of A}} \times 100\%$$

$$= \frac{4000}{6000} \times 100\%$$

$$= 66.67\% \text{ which is greater than min. Confidence.}$$

So, This is strong rule i.e. it is signifying positive correlation

Confusion Matrix for  $\chi^2$ -test :

	Game	Game	$\Sigma$ rows
video	4000 $(E=4500)$	$(7500 - 4000) = 3500$ $(3000 - E)$	7500
Video	$(6000 - 4000) = 2000$ $(E=1500)$	$\frac{\text{total}}{(10000 - 4000 - 3500 - 2000)} = 500$ $(E=500)$	2500
$\Sigma$ Columns	6000	4000	10,000

now, expected values:

$$\textcircled{1} \text{ For Video-game} = \frac{6000 \times 7500}{10000}$$
$$= 4500$$

$$\textcircled{2} \text{ For Video-game} = \frac{4000 \times 7500}{10000}$$
$$= 3000$$

$$\textcircled{3} \text{ For Video-game} = \frac{6000 \times 2500}{10000}$$
$$= 1500$$

$$\textcircled{4} \text{ For Video-game} = \frac{2500 \times 4000}{10000}$$
$$= 1000$$

now,

$$\textcircled{1} \quad x^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

$$= \frac{(4000 - 4500)^2}{4500} + \frac{(3500 - 3000)^2}{3000} + \frac{(2000 - 1500)^2}{1500}$$
$$+ \frac{(-500 - 1000)^2}{1000}$$
$$= 55.56 + 83.33 + 166.67 + 250$$
$$= 555.56 > 1$$

This signifies that there is a negative correlation.

$$\text{② Lift } (\text{frame}, \text{video}) = \frac{P(A \cup B)}{P(A) \cdot P(B)}$$

$$\text{i.e., } P(A \cup B) = \frac{4000}{10000} = 0.4$$

$$P(A) = \frac{6000}{10000} = 0.6$$

$$P(B) = \frac{7500}{10000} = 0.75$$

$$\text{i.e., Lift} = \frac{0.4}{0.6 \times 0.75}$$

$$= 0.89 < 1$$

which signifies negative correlation.

$$\text{③ Cosine } (A, B) = \frac{P(A \cup B)}{\sqrt{P(A) \cdot P(B)}}$$

$$= \frac{0.4}{\sqrt{0.6 \times 0.75}}$$

$$= 0.596 < 1$$

which signifies negative correlation.

$$\text{④ All confidence} = \frac{AB}{AB + \max(\bar{A}B, \bar{B}A)} \quad \text{or}$$

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

$$= \frac{40}{75}$$

$$= 0.53 < 1$$

which signifies negative correlation.

\* Consider the following Transaction with respective data for the videos & games?

*Dataset*

Dataset	G1	$\bar{G}_1$	$G_2\bar{V}$	$\bar{G}_2$	All confidence	Cosine	Lift	$\chi^2$
D0	4000	3500	2000	0	0.53	0.6	0.84	14778
D1	4000	3500	2000	500	0.53	0.6	0.89	555.6
D2	4000	3500	2000	10000	0.53	0.60	1.73	2013

= This are already evaluated values for:

All confidence, cosine, Lift &  $\chi^2$  but we need to calculate them using formulae in initial page.

Also for the data in D2 we have calculated the values of all 4-attributes in previous example.

for D2 let:

confusion matrix:

	G1	$\bar{G}_1$	$\Sigma_{row}$
$\Sigma_{column}$	4000	3500	7500
V	2000	10000	12000
	6000	13500	19500

$$\textcircled{1} \text{ Lift} = \frac{p(q|uv)}{p(q) \cdot p(v)}$$

$$\text{So, } p(q|uv) = \frac{4000}{19500} = 0.205$$

$$p(q) = \frac{6000}{19500} = 0.31$$

$$p(v) = \frac{13500}{19500} = 0.38$$

$$\text{Lift} = \frac{0.205}{0.31 \times 0.38}$$

= 1.73 > 1 which implies positive correlation.

$$\text{Cosine} = \frac{P(A \cup B)}{\sqrt{P(A) \cdot P(B)}}$$

$$= 0.60 < 1 \quad (\text{Negative correlation})$$

$$\textcircled{2} \quad \text{All-Correlation} = \frac{A \cdot B}{AB + \max(\bar{A}\bar{B}, \bar{A}B)}$$

$$= \frac{4000}{4000 + \max(2000, 3500)}$$

$$= \frac{4000}{4000 + 3500}$$

$$= 0.53 < 1 \quad (\text{Negative correlation})$$

Note:

- i) If there exists null value in given data set, Lift &  $\chi^2$  is not efficient measurement Attribute (result varies largely)
- ii) If the given data set is unbalance (not evenly distributed)
  - All-confidence is good measurement Attributes,
  - Cosine is not good measurement Attributes.

- ① Multi-level mining
- ② Multi-dimensional mining
- ③ Constrained-based mining
  - Monotonic
  - Antimonotonic
  - succinct
  - convertible

## ① Multi-level Mining:

① Uniform & non-uniform

② Concept hierarchy

Specialization:

Continent

Country

States

Streets

Generalization

{ support = 50%  
confidence = 60% }

{ support = 40%  
confidence = 50% }

{ support = 30%  
confidence = 40% }

- Uniform means equal distribution of support & confidence in all levels.
- Non-uniform means unequal distribution of support & confidence in all levels.

e.g.: Avg-sale of Sony TV in Asia (Generalized)

let  $x$  \$

Avg sale of Sony TV in Nepal (Specialized)

## ② Multi-dimensional?

Dimension (Attributes) - Two or more than two

e.g.: Age ( $x$ ,  $30 \dots 40$ )  $\wedge$  Salary ( $x$ ,  $50K \dots 80K$ )  $\Rightarrow$  buys ( $x$ ,  $car$ )  $\vee$  sales( $x$ , bike)

constrained

constrained

Implementing constraints as :-  $30 \dots 40$

$\text{avg}(E\text{-Age}) \leq 25$

$\min()$

$\max()$  etc.