

Subject: O.W.M

Sem: V

Q. Explain different OLAP operations, with suitable example.

Ans) Online Analytical Processing (OLAP) is a category software that allows users to analyze information from multiple database systems at the same time. It is a technology that enables analysts to extract and view business data from different point of view.

i) The OLAP cube consists of numeric facts called measures, which are categorized by dimensions. OLAP cube is also called the hyper cube.

Four types of analytical OLAP operations are : 1. Roll-up 2. Drill down 3. Slice and dice 4. Pivot.

iv) Let us understand these operations through example of college database having dimensions 1. Course 2. Student 3. Time and fact (measure) : Aggregate marks.

v) Hyper-Cube:-

		DIV A				DIV B			
		Student		Course		Time			
		S ₁	S ₂	C ₁	C ₂	Q ₁	Q ₂	Q ₃	Q ₄
		120	140	70	90	100	100	120	120
		83	83	100	150	83	83	83	83
		Q ₁	Q ₂	Q ₃	Q ₄	Q ₁	Q ₂	Q ₃	Q ₄
		150	150	200	200	100	100	120	120
		Time	(Quarters)	C ₁	C ₂	C ₃	C ₄	C ₁	C ₂
		Q ₁	Q ₂	Q ₃	Q ₄	Q ₁	Q ₂	Q ₃	Q ₄

← Block - A

In Block A we have dimension Student at level individual student, Time in Quarter and courses at department level. To columns and rows i.e each block fill with aggregate marks or NA values (Empty blocks).

1 Roll-up:- Roll-up is also known as "Consolidation" or "aggregation". The Roll-up operation can be performed in 2 ways:

- (a) Reducing dimensions (b) Climbing up concept hierarchy
- e.g:- Get the year wise marks for all students is roll-up as Semester wise marks are added to get year wise marks & it's one level up hierarchy of dimension:- Time.

Rollup		(Students)			
On-Time	(Year wise)	S ₁	S ₂	S ₃	S ₄
Time	(in Year)	Y ₁	300	250	230
	Y ₂		500		
	Y ₃			200	
	Y ₄				475
		C ₁	C ₂	C ₃	C ₄
		course			

A. Block → You → 100 → 500 → 475 (extra)

2)

Drill-down:- In drill-down data is fragmented into smaller parts. It is the opposite of the rollup process. The drill-down operation can be performed in two ways:-

- Moving down the concept hierarchy.
- Increasing dimension.

e.g:- Drill down on Courses to get unit test marks (Hatif-Syllabus Performance)

Drill down on Course-1		S ₁	20	22	
		S ₂	23	23	
		S ₃	15	17	
Q ₁			25	28	
Q ₂					
Q ₃					
Q ₄					
		U.T-1	U.T-2		
		Course-1			

If we drill-down on each course, dimensions will be double of original dimension of So, we dice only Course-1.

3) Slice:- One dimension is selected, and a new sub-cube is created.

e.g:- For Single student get the marks for all courses and all semester.

<u>Slice on</u>	O_1	C_1	C_2	C_3	C_4
Student	O_2	C_1	C_2	C_3	C_4
s_1	O_3	C_1	C_2	C_3	C_4
	O_4	C_1	C_2	C_3	C_4

c_2 is empty because sub didn't allocate for it.

3) ⑥ Dice :- Two or more dimension is selected, resulting in new sub-cube.

e.g.: Select Student S₂ & S₃ in Quarter Q₃ & Q₄ and Course C₂ & C₄.

Dice $(S_2 \text{ or } S_3) \& (O_3 \text{ or } O_4) \& (C_2 \& C_4)$:-

	S_2		
	S_3	150	100
Θ_3	200	120	
Θ_4	200	120	
	C_2	C_4	

4) Pivot: - In Pivot, you rotate the data axes to provide a substitute presentation of data. Analogous to matrix transpose.

e.g:- Pivoting result obtained in slice .

C_1	120	120	107	-109
C_2	No	No	No	No
C_3	70	70		
C_4		90	90	

Q.» Consider a data warehouse for a hospital where there are 3-dimensions namely
 (a) Doctor (b) Patient (c) Time and 2 measures (i) Count (ii) Charge where, charge is the fee that the doctor charges a patient for a visit.

i) Draw Star and Snow-flake Schema.

Ans There are Four tables:- 3-Dimensions tables + 1-Fact table with 2 measures.

Dimension Table:-

1. Doctor (Doctor_ID, name, phone, location, Pin, Specialization)

2. Patient (P.ID, name, phone, state, city, location, Pin)

3. Time (T.ID, day, month, week, Quarter, Year)

Fact-Table:-

4. Fact-Table (Doctor_ID, P.ID, T.ID, Count, Charge)

Fact-TableHospitalDoctorsPatient

Patient_ID

Name

Phone

State

City

Location

Pin

Doctor_ID

Patient_ID

Time_ID

Count

Charge

Doctor_ID

Name

Phone

Location

Pin

Specialization

Time

Time_ID

Day

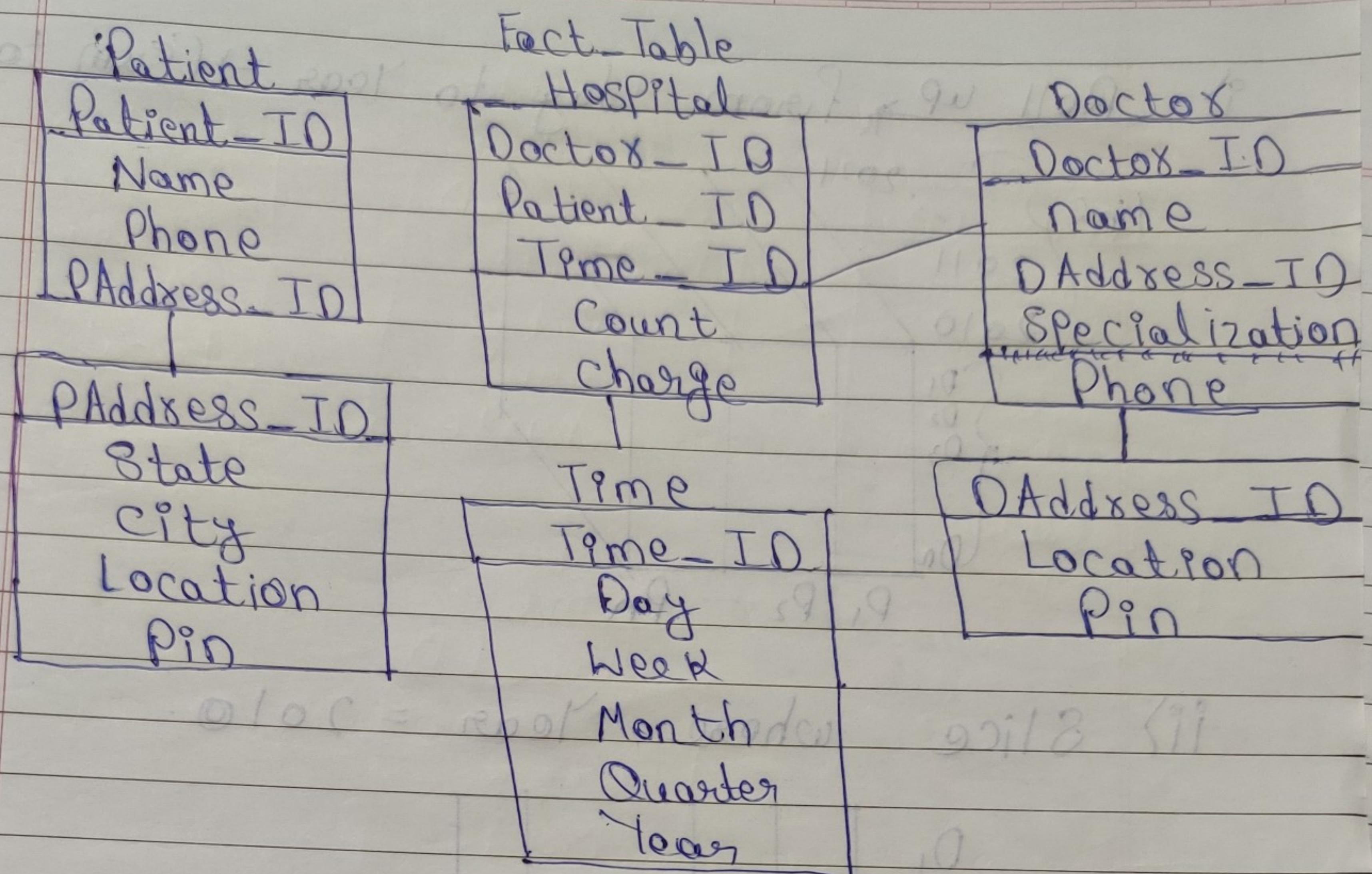
Week

Month

Quarter

Year

Star-Schema



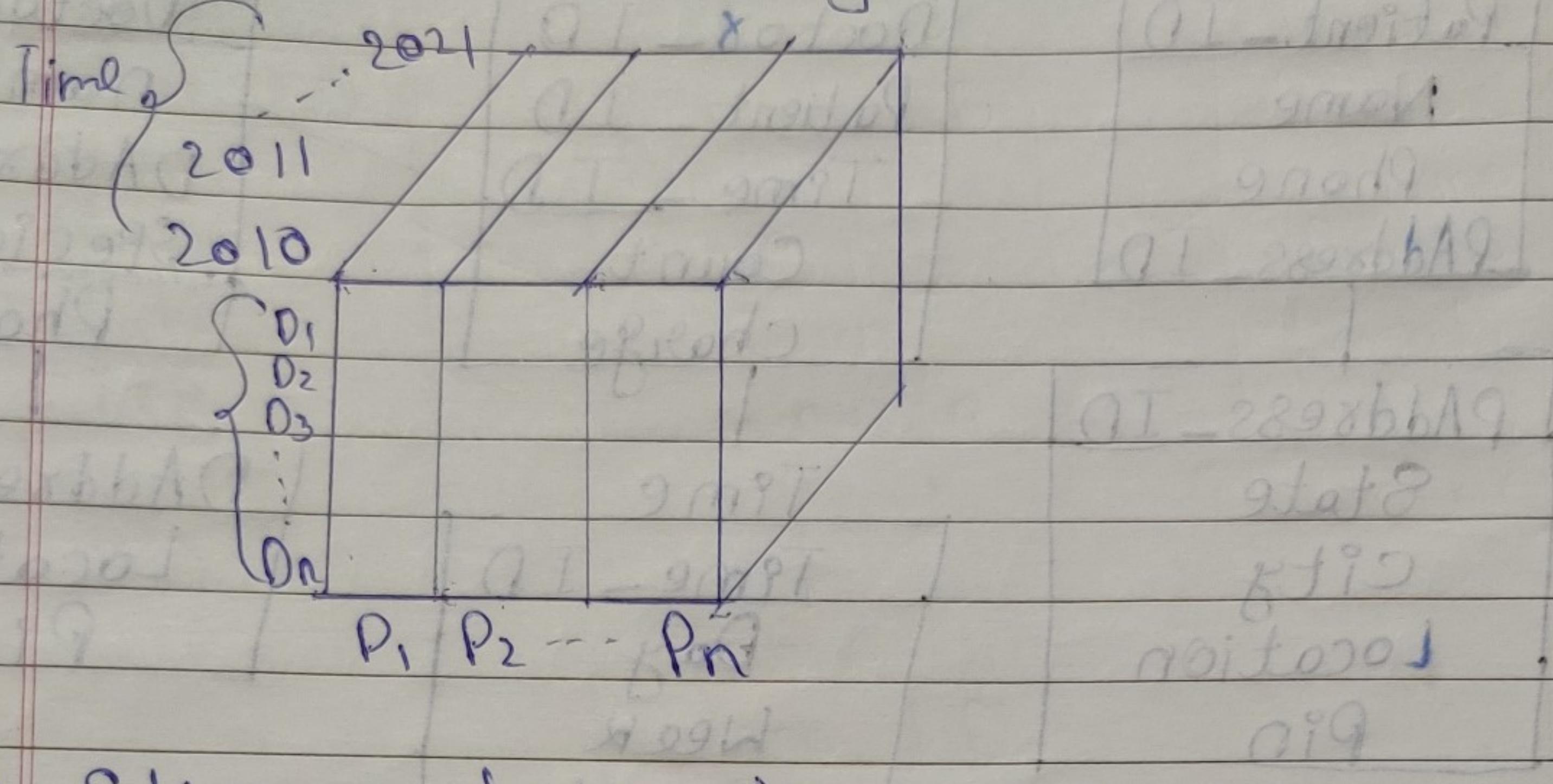
Snow-flake Schema.

ii) Starting with base cuboid [Day, Doctor, Patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2010.

Ans) Base Cuboid:-

		Time		
		T ₁	T ₂	T ₃
Doctor	D ₁	300	500	600
		400	600	900
Doctor	D ₂	350		
		400	1200	350
Doctor	D ₃	500		
		700	900	1200
Doctor	D ₄			

Roll up \rightarrow from Day to Year \rightarrow (sum to year)



ii) Slice where Year = 2010.

Roll-up on patient group by doctor.

D_1	000	000	000	cT	?	61000
D_2	000	000	000	cT	permit	12800 mid
D_3	020	020	000	cT	?	12800 mid
:	000	000	000	10	?	
D_n	020	020	000	50	?	12800 mid

Sum(Patient_Charge)

Result 003 10
— 005) 009 005 00

iii) To obtain the same list write an SQL query to the data are stored in R.D.B with the schema Fee. [Day, Month, Year, Doctor, Hospital, Patient, Count, Charge]

Ans) $\text{SELECT Doctor, SUM(Charge)} \text{ FROM Fee}$
 $\text{WHERE Year = 2010 GROUP BY Doctor.}$

Q.) Suppose that the data for analysis includes the attribute Age, the Age values for the data tuples are (in increasing order)
 $[13, 15, 16, 16, 19, 20, 25, 29, 33, 41, 44, 53, 62,$
 $69, 72]$ use min. max normalization to transform the value 45 for charge Age onto the range $[0.0, 1.0]$

$$\text{Ans) } V' = \frac{(V - \text{minage})}{\text{maxage} - \text{minage}} [\text{new_maxage} - \text{new_minage}]$$

$$V = 45, \text{minage} = 13, \text{maxage} = 72$$

$$\text{new_maxage} = 1.0, \text{new_minage} = 0.0$$

$$V' = \frac{(45 - 13)}{72 - 13} [1.0 - 0.0] + 0.0$$

$$\therefore \frac{16}{31} = 0.51613$$

Q.1) Describe the steps involved in data mining when viewed as a process of knowledge discovery.

Ans Steps involved in data mining when viewed as process of knowledge discovery:-

1.) Data cleaning:- Data from multiple source (Database, flat-files, and different form) are extracted and clean to remove noise and inconsistent data.

2.) Data integration:- Since data is extracted from multiple source they may be combined to provide useful information. Domain knowledge is necessary for proper integration. Generally, stored in R.D.B.

3.) Data Selection:- After cleaning & integration data relevant to the analysis task are derived from the data warehouse. e.g:- Feature Sub-set Selection.

4.) Data transformation:- Data are transformed and consolidated into forms appropriate for mining by performing Data reduction, different transformation techniques & aggregation.

5) Data-mining:- It is an essential process where Intelligent methods are applied, algorithms from various different fields such as Machine learning, Image processing and Natural language processing are used to extract novel data patterns.

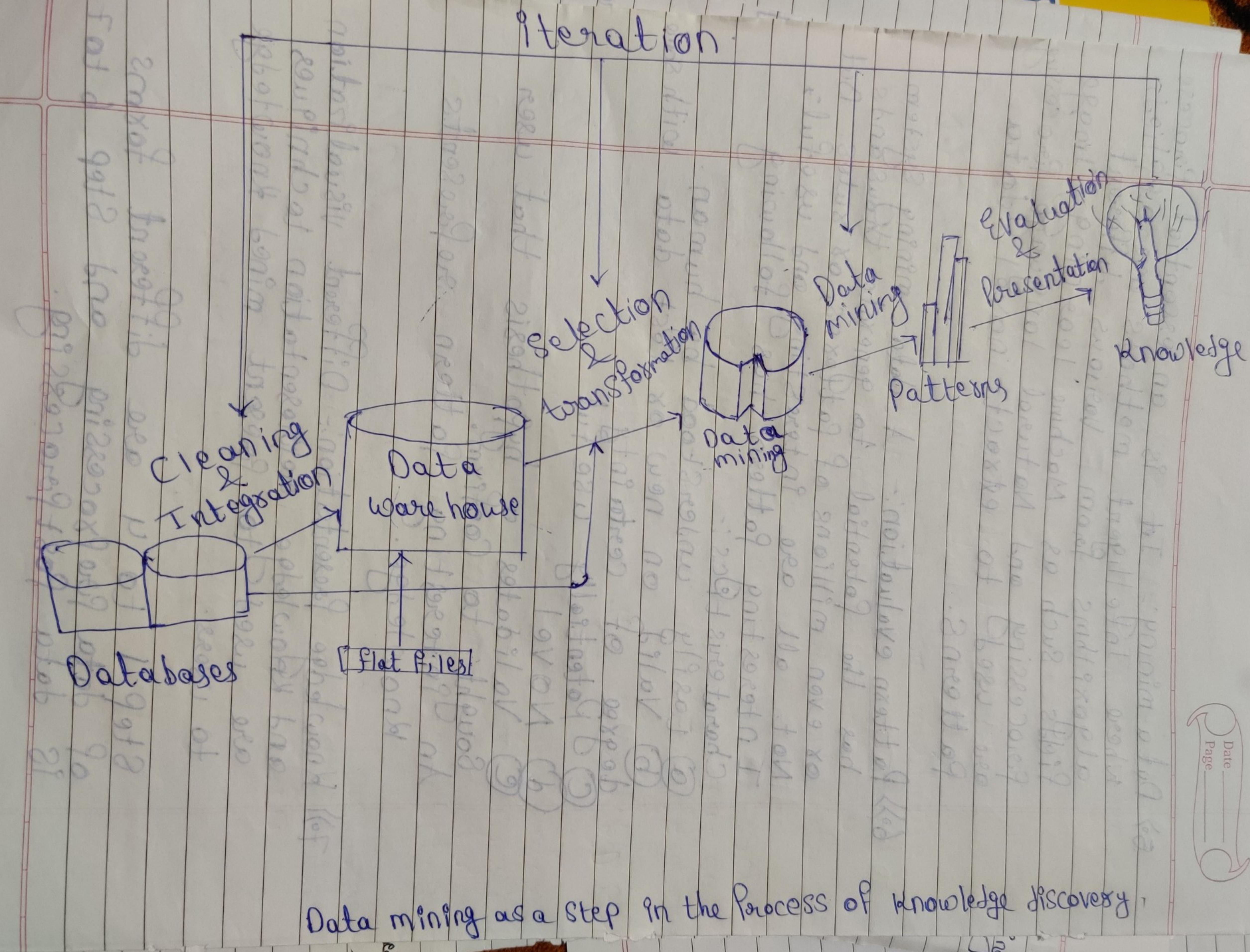
6) Pattern evaluation:- A data mining system has the potential to generate thousands or even millions of patterns, or rules. But not all are interesting and useful; Interesting pattern has following characteristics:-

- (a) Easily understood by human.
- (b) Valid on new or test data with some degree of certainty.
- (c) Potentially useful.
- (d) Novel.
- (e) Validates a hypothesis that user sought to confirm.

An interesting pattern represents knowledge.

7) Knowledge presentation:- Different visualization and knowledge representation techniques are used to present mined knowledge to users.

Step 1 to 4 are different forms of data preprocessing and Step 6 to 7 is data postprocessing.



<<Q.) Why is tree Pruning is useful in decision tree induction? what is a drawback of using a separate set of tuples to evaluate pruning. Given a decision tree you have the option of a) Converting D.T to results & then pruning the resulting rules. b) Pruning the D.T & then converting pruned tree to rules, what advantage does (a) have over (b)?

Ans)

- ① The decision tree built may overfit the training data. There could be too many branches some of which may reflect anomalies in training data due to noise or outliers.
- ② Tree Pruning addresses this issue of overfitting the data by removing the least reliable branches (using statistical measures).
- ③ This generally results in more compact and reliable decision tree that is faster and more accurate in its classification of data.
- ④ The drawback of using separate set of tuples to evaluate pruning is that it may not be representative of training tuples used to create the original decision tree.
- ⑤ If the separate set of tuples are skewed then using them to evaluate the pruned tree would not be a good indicator of the pruned tree's classification accuracy.
- ⑥ Furthermore, using separate set of tuples to evaluate pruning means there are less tuples to use for creation and testing of the tree.

(ii)

once a decision tree has been constructed, it is a simple matter to convert it into an equivalent set of rules.

Converting a decision tree to rules before has three main advantages [Advantages of (a) over (b)]:

1. Converting to rules allows distinguishing among the different contexts in which a decision node is used.

(a) since each distinct path through the decision tree/node produces a distinct rule, the pruning decision regarding that attribute test can be made differently for each path.

(b) In contrast, if the tree itself were pruned, then only two choices would be: ① Remove the decision node completely or ② retain it in its original form.

2. Converting to rule removes the distinction between attribute tests that occur near the root of the tree & those that occur near the leaves. We thus avoid messy book keeping issues such as how to reorganize the tree if the root node is pruned while retaining part of subtree below this test.

3. Converting to rules improves readability. Rules are often easier for human to understand.

To generate rules, trace each path in the decision tree, from root node to leaf nodes, recording

the test outcomes as antecedents and the leaf-node classification as the consequent.

Q) A database has 5 transactions, let
min. Support = 60% & min. Confidence = 80%

TID
T₁₀₀
T₂₀₀
T₃₀₀
T₄₀₀
T₅₀₀

Item Bought
 $\{M, O, N, K, E, Y\}$
 $\{D, O, N, K, E, Y\}$
 $\{M, A, K, E\}$
 $\{M, U, C, K, Y\}$
 $\{C, O, O, K, I, E\}$

(a) Find all frequent item set using Apriori algorithm.

(b) List all the strong association rules matching the following, where x_i is a variable representing customers & item i denotes variable representing items

$\forall x \in \text{transaction}, \text{buys}(x, i \text{ item})$
 $\wedge \text{buys}(x, j \text{ item}) \Rightarrow \text{buy}(x, i \text{ item})$

Sol: Frequent item set:

A → 1	→ 1 5	< 60% X
C → 2	→ 2 5	< 60% X
D → 1	→ 1 5	< 60% X
E → 4	→ 4 5	>= 60% ✓
K → 5	→ 5 5	>= 60% ✓
I → 1	→ 1 5	< 60% X
M → 3	→ 3 5	>= 60% ✓
N → 2	→ 2 5	< 60% X
O → 4	→ 4 5	>= 60% ✓
U → 1	→ 1 5	< 60% X

$$Y \rightarrow 3 \rightarrow 3/5 > 60\% \checkmark$$

1st frequent item-set = $\{E, K, M, O, Y\}$

2nd frequent item-set:

$$\{E, K\} \rightarrow 4 \rightarrow 4/5 \checkmark$$

$$\{E, M\} \rightarrow 2 \rightarrow 2/5 \times$$

$$\{E, O\} \rightarrow 3 \rightarrow 3/5 \checkmark$$

$$\{E, Y\} \rightarrow 2 \rightarrow 2/5 \times$$

$$\{K, M\} \rightarrow 3 \rightarrow 3/5 \checkmark$$

$$\{K, O\} \rightarrow 3 \rightarrow 3/5 \checkmark$$

$$\{K, Y\} \rightarrow 3 \rightarrow 3/5 \checkmark$$

$$\{M, O\} \rightarrow 1 \rightarrow 1/5 \times$$

$$\{M, Y\} \rightarrow 2 \rightarrow 2/5 \times$$

$$\{O, Y\} \rightarrow 2 \rightarrow 2/5 \times$$

2nd frequent item-set $\Rightarrow \{\{E, K\}, \{E, O\}, \{K, M\}, \{K, O\}, \{K, Y\}, \{O, Y\}\}$

3rd Frequent item-set:

$\{E, K, O\} \rightarrow 3 \rightarrow 3/5 \checkmark$

$\{E, K, M\} \rightarrow 1 \rightarrow 1/5 \times$

$\{E, K, Y\} \rightarrow 2 \rightarrow 2/5 \times$

$\{K, M, O\} \rightarrow 0 \times$

$\{K, M, Y\} \rightarrow 2 \rightarrow 2/5 \times$

$\{K, O, Y\} \rightarrow 2 \rightarrow 2/5 \times$

3rd - Frequent item set = $\{E, K, O\}$

Association Rule	Support	Confidence
------------------	---------	------------

$\{E, K\} \rightarrow \{O\}$	3	$3/4 \Rightarrow 75\% \times$
$\{O\} \rightarrow \{E, K\}$	3	$3/4 \times$
$\{E, O\} \rightarrow \{K\}$	3	$3/3 = 100\% \checkmark$
$\{K\} \rightarrow \{E, O\}$	3	$3/5 \Rightarrow 60\% \times$
$\{O, K\} \rightarrow \{E\}$	3	$3/3 = 100\% \checkmark$
$\{E\} \rightarrow \{O, K\}$	3	$3/4 \Rightarrow 75\% \times$

b) List of association rules that satisfies min. confidence \Rightarrow

$\langle E \wedge O \rangle \rightarrow K$
$\langle O \wedge K \rangle \rightarrow E$