

Assignment No. 01

Date:

--	--	--

Q.1

What is meant by metadata in the context of a data warehouse?
 Explain the different types of metadata stored in the data warehouse
 Illustrate with a suitable example.

- a) Metadata is simply defined as data about data. The data that is used to represent other data, is known as metadata.
- b) Metadata are created for the data names & definitions of given warehouse.
- c) Additional metadata are created & captured for timestamping any extracted data, the source of the extracted data and missing fields that have been added by data cleaning or integration process.
- d) Metadata act as a directory. The directory helps the decision support system to locate the contents of a data warehouse.
- e) Metadata is the road map to a data warehouse.
- f) Metadata can be broadly categorized into three types :

I. Business Metadata :-

It is the data ownership information metadata are the data that remains warehouse objects.

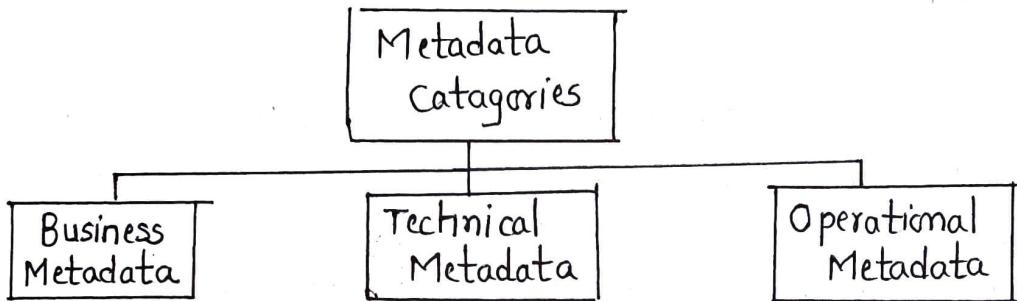
II. Technical Metadata :-

It includes database system names - table & column names and sizes, data types and allowed values. It also includes structural information, such as primary & foreign key attributes, and indices.

III. Operational Metadata :-

It includes currency of data & data and data lineage. Currency of data means whenever the data is active

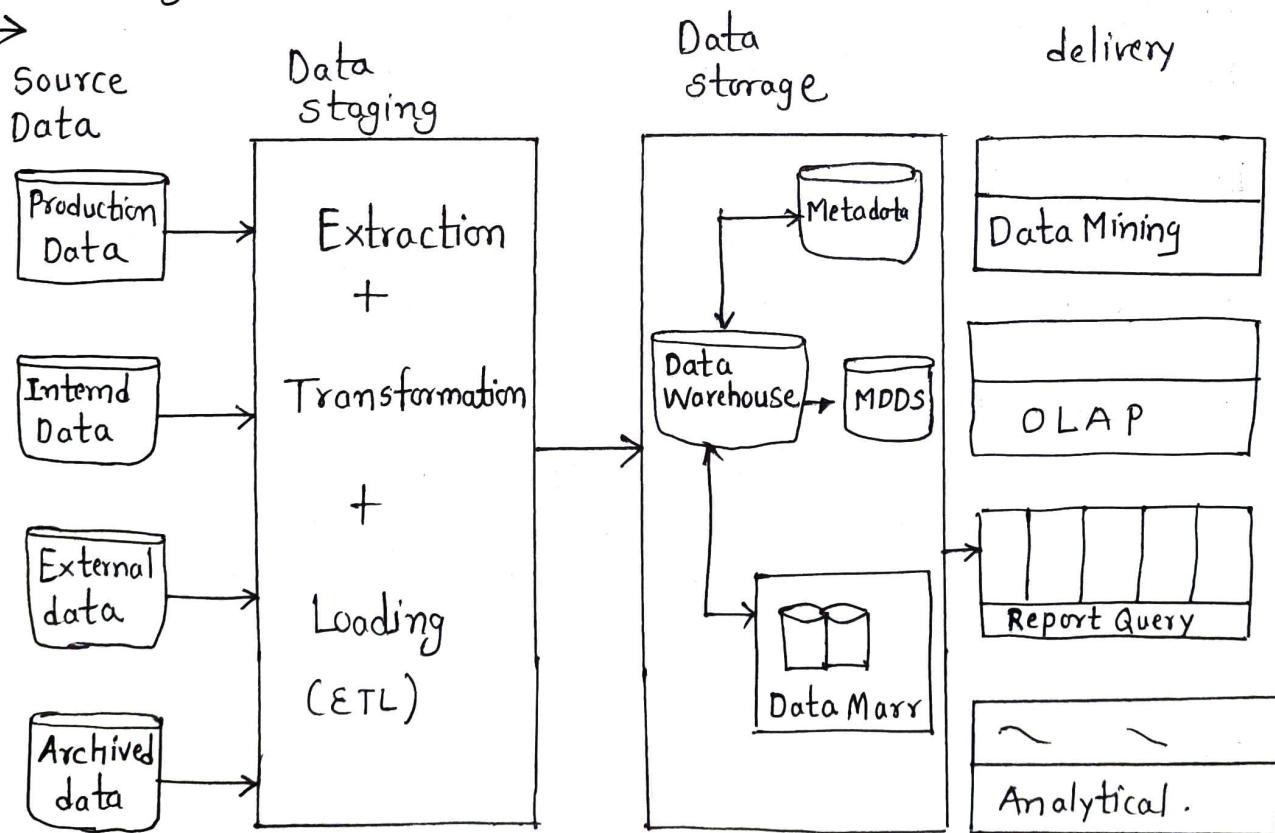
archived or purged. Lineage of data means the history of data migrated & transformation.



Example:-

The index of book service as a metadata for the contents in the book. Metadata is the summarized data that leads to detailed data.

Q.2 Explain the architecture of Data warehouse with a suitable block diagram.



Date :

i) Data Source layer :-

The layer representing various data sources that feed data into the datawarehouse. The data into the datawarehouse. The data into the datawarehouse. The data can be in any of These formats : Plain text file, relational database, Excel File & other types of data can act as data sources.

- a. Production source - Represents sales data, HR data & Product data.
- b. Internal data - Represents data of a department of an organisation such as employ data
- c. External data :- Represents data from outside the organisation or third party such as survey data of demographic.
- d. Achieved data :- Represents logs of the web server along with the user's browsing data.

ii) Data staging layer :-

The storage area for data processing where data comes before being transformed into the data, that is entered in a data that entered in a datawarehouse.

- a. Extraction - The process of extracting data from different source system and validating it against certain quality.
- b. Data warehouse - It is maintained by organisations as central warehouse of data that can be equally accessed by all business experts & end users.
- c. Loading :- The process of loading the data either from datawarehouse data mart.

(iii) Data Storage layer :-

The Layer in which the transformed data & cleaned data is stored.

- a. Data Warehouse :- It is maintained by organisations as central warehouse of data that can be equally accessed by all business experts & end users.
- b. Data Mart :- When data warehouse is created at the department level is known as data mart.
- c. Metadata :- Details about the data is known as metadata. In other words, it is a catalog of data warehouse.
- d. MDDS - It is multidimensional database that allows data to be molded & viewed in multiple dimensions. It is defined by dimensions & facts.

iv) Information Delivery:

It provides the information that reached to end users. The information that can be in any form such as tables, chart, graphs or histograms.

- a. Data Mining :- The process of finding relevant & useful information , large amount of data.
- b. OLAP :- Allows the navigation of data of different levels abstraction . such as down, rollup, slice, dice, & so on .

Q.3

We would like to view sales data of a company w.r.t. 3D's namely location, Item time. Represents sales data in the form of a 3D data cube for the above and perform Roll-up, Drill down, Slice, Dice & pivot OLAP operations on the above data cube & illustrate.



Dimensions → Location, Time
Facts → Sales Data.

		MUM	280	400	567	875	
		Kolhapur	986	85	987	800	
		Rat	786	987	504	987	
		Sangli	788	987	500	765	
	↑	Q1	788	987	500	765	
Time		Q2	678	654	987	540	
(quarter)		Q3	899	875	480	190	
		Q4	787	969	908	1000	
			PC	mase	sho	book	

Items



i) Roll-up

Maharashtra

2840
2459

(Roll-up Location)

Q1	1576	1974	1000	1530
Q2	678	654	987	540
Q3	899	875	480	190
Q4				
	PC	Mouse	Shoe	Book

ii) Roll-down

Mum
Kohi
Rat

Sangli	167	235	125	190
Jan	167	235	125	190
Feb	205	245	125	180
Mar	200	240	125	210
Apr	216	267	125	185
May				
June				
July				
Aug				
Sep				
OCT				
NOV				
DEC				

Date :

iii) Slice

Mum				
Kohl				
Rat				
Sangli	788	987	500	765

PC Mouse shoe Book

iv) Dice

Ratnagiri	786	987	
Sangli	788	987	
Q1:	788	987	
Q2:	678	697	

PC Mouse

v) Divot

				788	PC
				987	Mouse
				500	shoe
				765	Book

Mum Kohl Ratn Sangli

Q.4 Consider a data warehouse for a hospital where there are three dimensions.

- a. Doctor
- b. Patient
- c. Time

Consider two measures

i) Count

ii) charge where charge is the same fee that the doctor charges a patient for a visit for the above example, create a cube & illustrate the OLAP operations.

→ There are four tables out of 3 dimensions tables & 1 fact tables

Dimension tables:-

i) Doctor (DID, name, phone, number, Location, pin, specialization)

ii) Patient (PID, name, phoneno, state, city, Location, pin)

iii) Time (TID, Day, month, quarter, year)

Fact table :- Fact table (D-ID, PID, TID, count, charge)

Doctor				
		D2	02	03
		D1	04	01
(Time quarter)		Q1	04	01
		Q2	07	10
		Q3	02	03
		Q4	01	02

Fact -
Count (No. of times
patient visited
doctor)

Date :

doctor	D ₂	500	1000	100		
	D ₁	1000	200	100		
	Q ₁	1000	200	100		
time	Q ₂	2000	4000	100		
(quarter)	Q ₃	500	500	400		
	Q ₄	100	500	300		
		P ₁	P ₂	P ₃		

fact - sum or fees paid by patient.

- Roll-up

Doctors	1500	1200	200		
Q ₁	2000	400	200		
Q ₂	2000				
Q ₃	560				
Q ₄	100				
	P ₁	P ₂	P ₃		

Roll Up for fees.

patient

- Drill down :-

" "

Doctors		D ₂	500	1000	100
D ₁		200	200	00	
Jan	200	100	00		
Feb	400	0	00		
Mar	400	100	100		
Apr					
May					
June					
July					
Aug					
Sept					
Oct					
Nov					
Dec					

Drill Down
for fees

- Slice

D ₂	500	1000	100
D ₁	1000	200	100
	P ₁	P ₂	P ₃

Dice

D ₂	500	4000
D ₁	100	200
Q ₁	1000	200
Q ₂	2000	4000
	P ₁	P ₂

- Pivot :-

D ₂	500	1000	100
D ₁	1000	200	100



1000	500	P ₁
200	1000	P ₂
100	100	P ₃

Date :

Q.5

The college wants to record the marks for the course completed by students using the dimensions

a) Course b) Students c) Time & measure Aggregate marks

Create a cube & apply the OLAP operations on the cube for illustration.

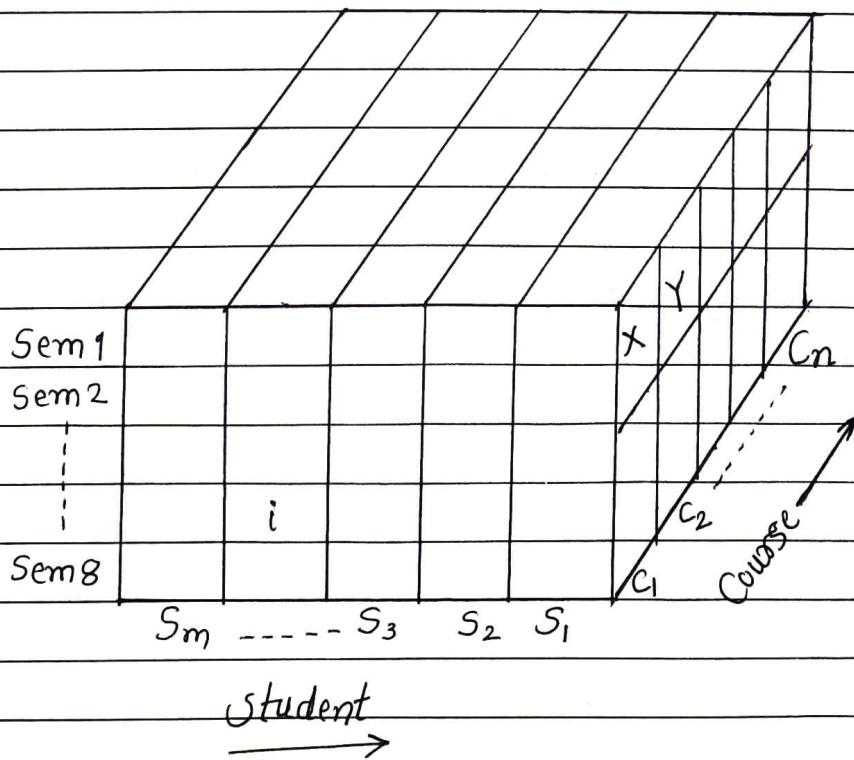


Marks obtained \Rightarrow fact

Dimensions - course, student & time

Let the current level of grain is each cell gives aggregate marks of each student for an individual cluster of course for each semester.

"X" is the total aggregate marks obtained by students 'S' in semester 1 for all the subjects that belongs to cluster c_j .

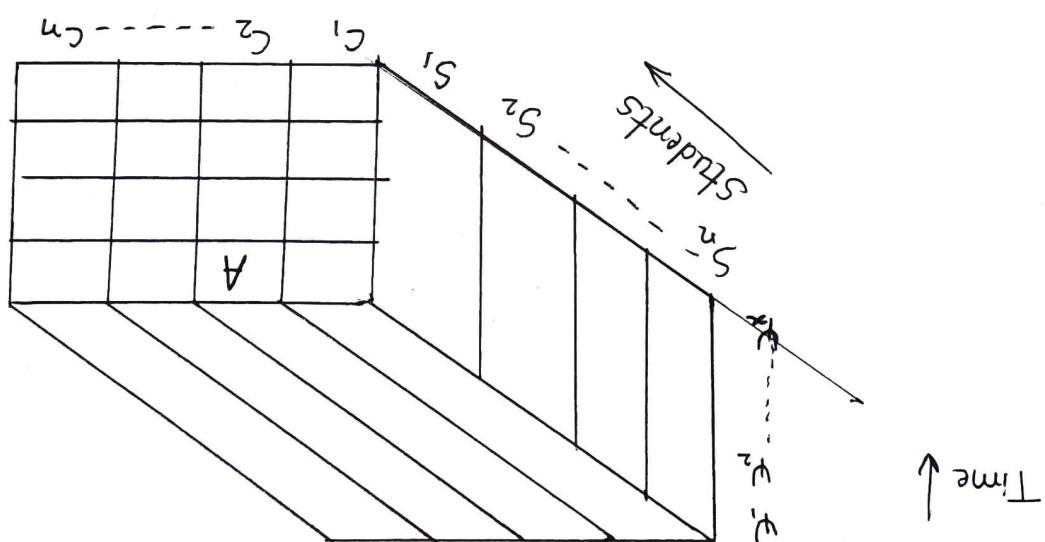


Roll-up :-

Obtaining aggregated/generalised level of information from

the current level is called roll-up.

Here we obtain year wise aggregate marks from semesterswise aggregate marks.



aggregate marks.

Here we obtain year wise aggregate marks from semesterswise

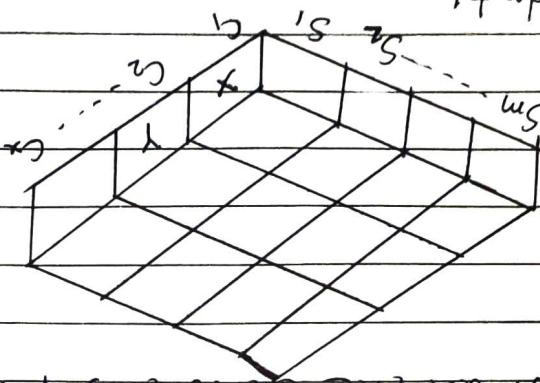
Roll-down :-

Obtaining detailed level of information from the

current level is called down.

For example, getting weekly profit from monthly profit. It is done by decomposing the cells along one or more dimensions. In this example, obtaining aggregate marks for individual courses from cluster level.

Student 1

Time ↑
Sem 1

In this example, obtaining one specific part of cube for detailed investigation is called slice & dice. for example getting profit for all products for all the years for a specific region.

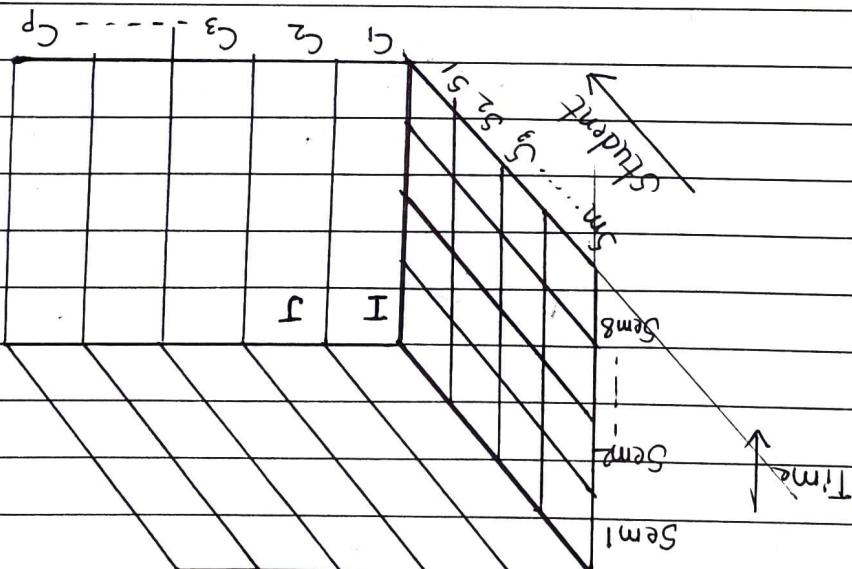
In this example aggregating all the marks for all the students for all the subjects for sem 1.

Slice & Dice :-

Here, "I" is the aggregate marks obtained by student S1 in semester for course that belongs to cluster C1.

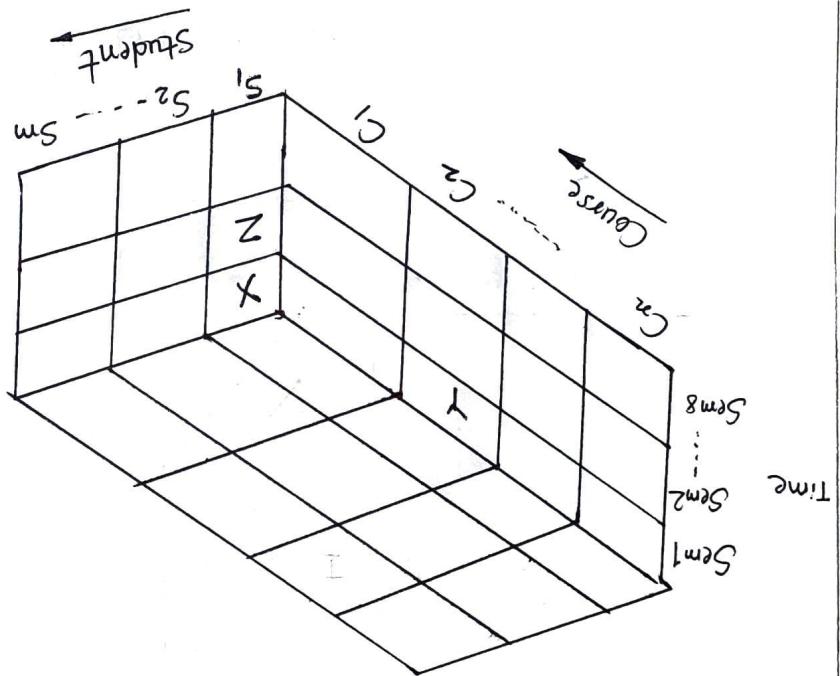
←

Course



Date :

--	--	--



If it is rotation of the cube, it is used to analyze the same data from different perspective, for example, analyzing profit margin wise & then product wise.

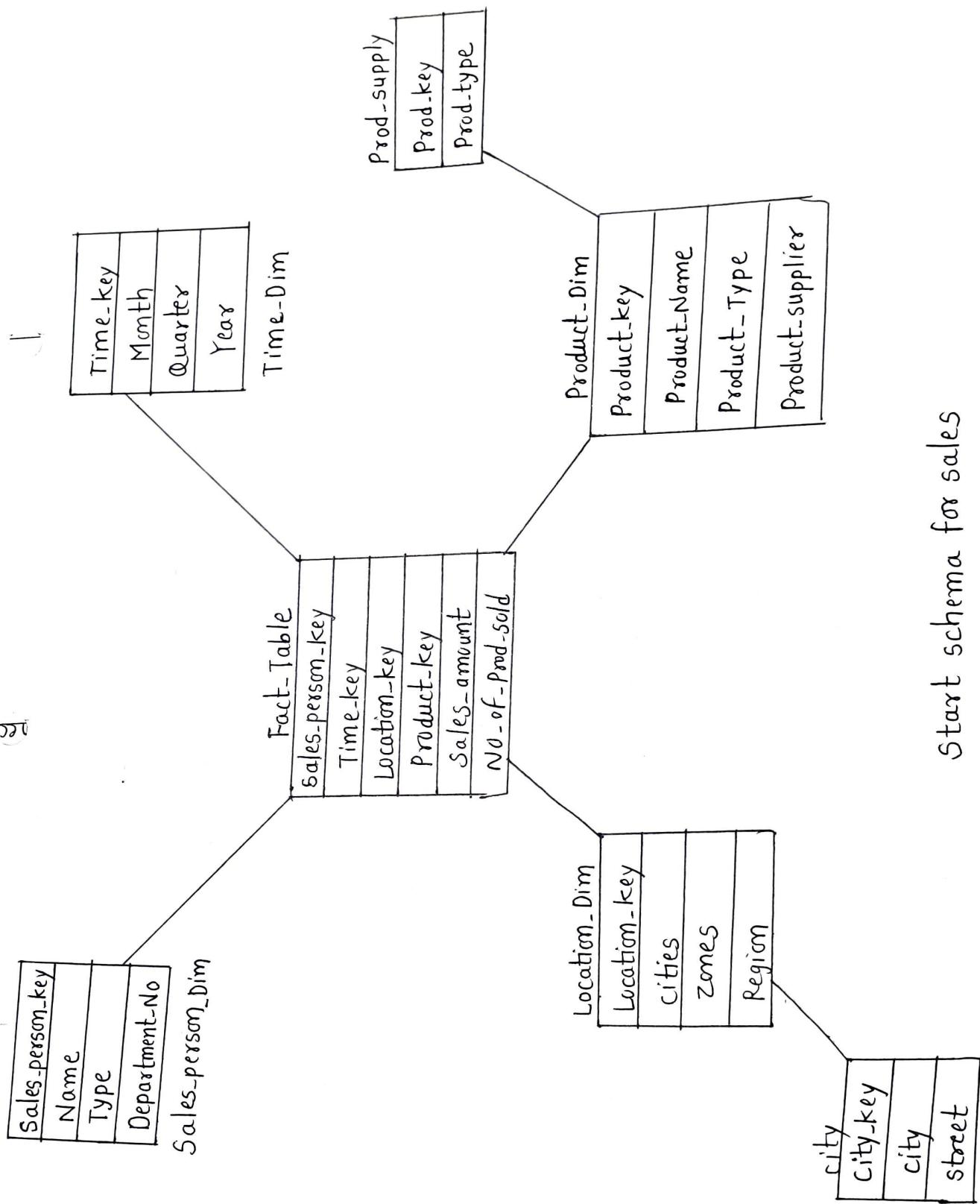
In this example, analyzing aggregate marks obtained by students character wise & then student wise.

Pivot :-

A manufacturing company has a large huge sales network. to centralize the sales, it is divided into regions, each region has multiple zones. Each zone has different cities, each sales person is allowed to different cities. The objective is to track sales figures at to different cities. The objective is to track sales figures at different granularity levels for region, sales person at time. convert the star schema to snowflake schema.

Each zone has different cities. Such sales person is allocated different cities. The objective is to track sales fig. at different cities. The objective is to count no. of products sold. granularity levels of region & to count no. of products sold. Design a star schema by considering granularity levels for region, sales person at time convert the star schema to snowflake schema.

(Late:



Star schema vs Snowflake schema.	Star Schema	Snowflake Schema
1. In star schema, The fact tables as well as tables as the dimension tables are contained.	1. In star schema, The fact tables dimensions, tables as well as sub dimensions tables are contained.	If uses more space
2. Uses less space	It uses less space.	It is bottom-up model.
3. Star schema uses top - down model.	It is bottom-up model.	4. It takes less time for execution of queries.
4. It takes less time for execution of queries.	While it takes more time than star schema for execution of queries.	5. In star schema, Normalization is not used.
5. In star schema, Normalization and denormalization are used.		

Write down differences between the following.

7.8

Date:

Data Mart	Data Warehouse
While it is a decentralised system.	Data Warehouse is a centralised system.
While it is a denormalized mart, highly denormalized takes place.	In data warehouse, highly denormalization takes place.
While in data mart, highly denormalized takes place.	Data-Warehouse is top-down model.
While it is a bottom up Model.	To build a warehouse is difficult
While to built Mart is easy.	In Data Warehouse, Fact
In data mart, star schema & snowflake schema are used.	comstellation schema is used.

b. Data Warehouse vs Data Mart.

OLTP	OLAP	OLTP & OLAP.
OLTP consists only from various databases.	operational current data.	C. OLTP & OLAP.
It is subject oriented used for Data mining, Analytics, Decision making etc.	Used for business tasks.	3. It provides a multidimensional view of different business tasks.
If reveals a snapshot of present	Business tasks.	4. The size of the data is relatively small as the historical data is archived.
Large amount of data is stored typically in TB PB.	Few e.g. MB, GB.	

Date: _____