

Module 1

Q1) ER modelling and dimensional modelling

\Rightarrow	ER modelling	Dimensional modelling
①	Transaction oriented	① Subject oriented
②	Its core components are entities and relationships	② Its core components are fact tables and dimension tables
③	It eliminates redundancy	③ It plans for redundancy
④	Highly volatile data	④ Non-volatile data
⑤	Physical and logical model	⑤ Physical model
⑥	Normalization is suggested	⑥ De-normalization is suggested.
⑦	OLTP application	⑦ OLAP application
⑧	Type of information used is real-time information	⑧ Type of information used is historical information
⑨	High transaction volume	⑨ Low transaction volume.
⑩	Application is used for buying products	⑩ Application is used to analyze buying patterns

(Q2) OLTP and OLAP

	OLTP	OLAP
①	Online transaction processing	Online analytical processing
②	It makes use of standard DBMS.	It makes use of data warehouse
③	Tables are normalized	Tables are not normalized.
④	The data is used to perform day-to-day operations.	The data is used in planning, problem-solving and decision making.
⑤	The size of the data is relatively small	The size of the data is relatively large.
⑥	Queries are very fast	Queries are relatively slow
⑦	Both read and write operations	Only read and rarely write operations
⑧	It serves the purpose to Insert, update, and delete the information from the database	It extracts the information for analysis and decision-making.

(3) Star Schema & Snowflake schema

⇒ Star Schema Snowflake schema

- | | |
|--|---|
| (1) Top-down model | (1) Bottom up model |
| (2) Uses more space | (2) Uses less space |
| (3) Takes less time | (3) Takes more time |
| (4) Normalization is not used | (4) Normalization and De-normalization are used |
| (5) Design is very simple | (5) Design is complex |
| (6) Query complexity is low | (6) Query complexity is higher |
| (7) Less number of foreign keys | (7) More number of foreign keys |
| (8) High data redundancy | (8) Low data redundancy |
| (9) The fact tables and dimension tables are contained | (9) The fact tables, dimension tables as well as sub dimension tables are contained |

(q4) Top down approach & bottom up approach

⇒ Top down approach Bottom up approach

① Centralized data ① Data marts are built first.

② Development time is longer ② Development time is shorter.

③ High complexity ③ Initially low complexity

④ High scalability ④ Moderate scalability

⑤ High risk

⑤ Low risk

⑥ Ensures consistency ⑥ Does not ensure consistency.

⑦ High investment required.

⑦ Incremental investment required as data marts grow (low)

(Q4) Data warehouse design strategies

⇒ (1) A data warehouse is a single data repository where a record from multiple data sources is integrated for OLAP.

(2) It collects, organizes and stores data from different places to help you make better decisions.

(3) When building a data warehouse, there are two main ways to do it:

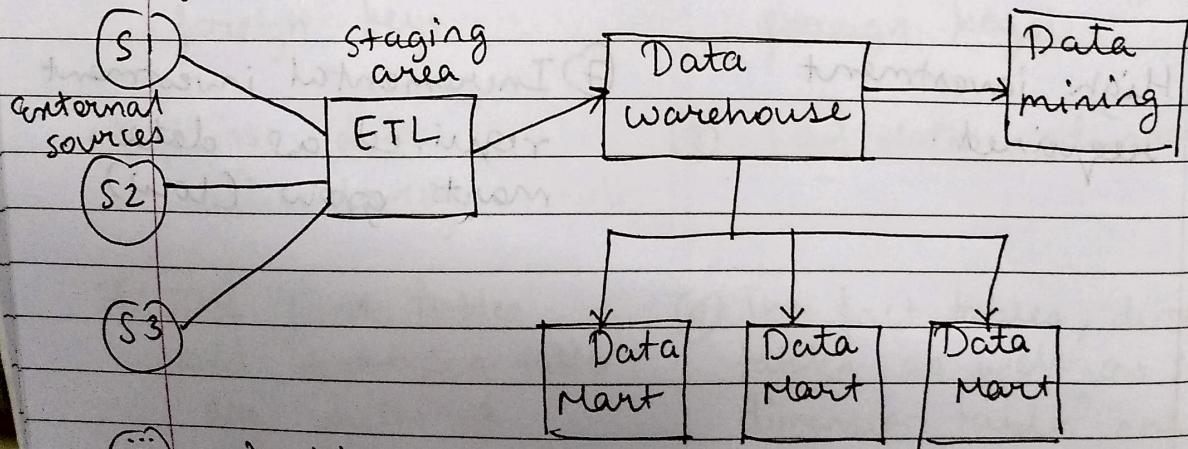
i) Top down approach

ii) Bottom up approach

(4) Top down approach

i) This method starts by building a big, central data warehouse for the entire organization.

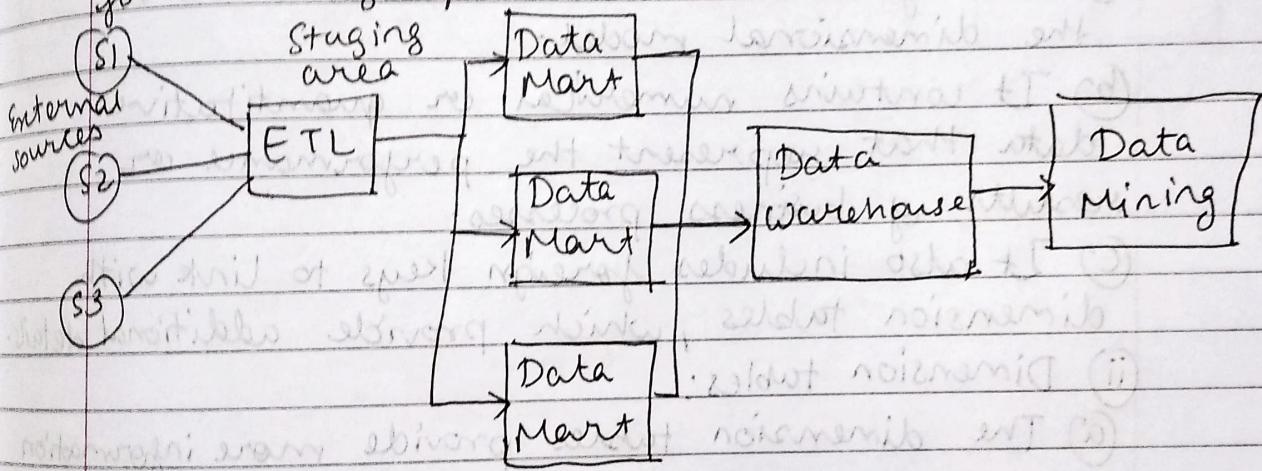
ii) Then, smaller, focused sections called data marts are created from the main warehouse for specific departments.



iii) Working

- Data is collected from various sources
- Data is processed & validated using ETL tools
 - Extract: Data is collected from sources
 - Transform: Data is standardized & cleansed
 - Load: Data is loaded in the data warehouse

- (e) After cleaning, the data is stored in single, organized location, data warehouse.
- (d) Smaller data marts are created for each department to focus only on their needs.
- (c) The company uses reporting tools to get insights, trends and predictions from the data.
- (b) Bottom up approach
- (i) The method focuses on creating small, separate data marts first for specific departments.
- (ii) These data marts are then connected later to form a larger, centralized data warehouse.



- (iii) Working:
- (a) Data is collected from various external sources.
 - (b) Data is processed using ETL tools, similar to the top-down approach.
 - (c) Individual data marts are built for specific needs.
 - (d) These data marts are then combined to form a bigger, central data warehouse.
 - (e) Each department can use the data warehouse for analyzing the data.

(P5) Dimensional modelling.

⇒ ① Dimensional modelling is a way to organize data in a database to make it easy to analyze and understand.

② It is mostly used in data warehouses for decision-making and organizing data into a cube for better viewing and analysis.

③ Dimensional modelling involves two types of tables:

i) Fact table

a) The fact table is a central table in the dimensional model.

b) It contains numerical or quantitative data that represent the performance or results of business processes.

c) It also includes foreign keys to link with dimension tables, which provide additional details.

ii) Dimension tables:

a) The dimension tables provide more information for the facts in the fact table.

b) It contains descriptive or categorical attributes that helps to explain the measures in the fact table.

④ Example :

A city and state can view a store summary in a fact table. Item summary, customer information can be viewed in dimension table Sales (storeID, ItemID, CustID, qty, price)

StoreID (storeID, city, state)

ItemID (ItemID, category, brand, color, size)

CustID (CustID, name, address).

Fact table

StoreID	ItemID	CustID	qty	price
---------	--------	--------	-----	-------

StoreID	ItemID	CustID

Dimension
tables

(Q6) Basic building blocks of data warehouse
OR

Components of data warehouse

OR

Data warehouse architecture.

- ⇒ (1) A data warehouse is a system that stores data from various sources in one place, making it easier for businesses to analyze and make decisions.
- (2) Its components or building blocks work together to collect, organize, store and present data.
- (3) The components are:
 - i) Data sources
 - (a) These are the initial points from where data is collected and generated.
 - (b) Its types:
 - (i) Operational data : Data from day - to - day business systems.
 - (ii) Internal data : Private data such as reports
 - (iii) Archived data : Older data stored for historical analysis.
 - (iv) External data : Data from outside sources
 - ii) ETL Process
 - (a) Extract : Extracts raw data from various sources using queries.
 - (b) Transform : Data is standardized and cleaned to remove errors
 - (c) Load : Save the prepared data in the data warehouse.
 - iii) Data storage
 - (a) The central repository where the cleaned,

processed data is stored.

(b) Components :

(1) Staging area : Temporary storage where raw data is held before processing.

(2) Warehouse storage : The main storage area.

(3) Data marts : Subsections of warehouse.

(iv) Metadata

(a) Metadata describes the structure, content and the meaning of the data stored in warehouse.

(v) Information delivery

(a) This component ensures data is accessible and usable for decision-making.

(b) Eg: A dashboard showing monthly revenue trends

(vi) Data modelling

(a) The process of organizing the data in the warehouse to support efficient querying and analysis.

(v) Common models :

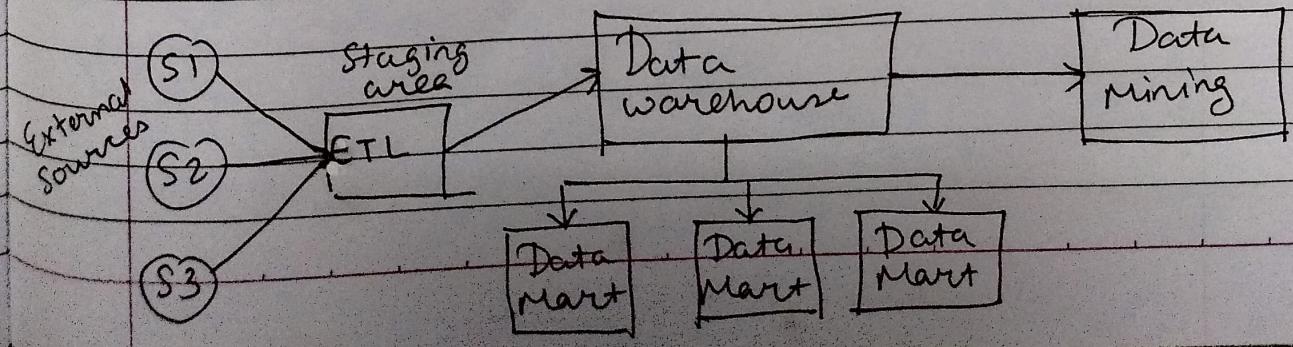
(1) Star schema

(2) Snowflake schema

(vii) Management and control

(a) This includes the systems and processes that maintain the warehouse's functionality and security.

(b) It schedules ETL tasks, monitors performance, ensures data integrity, security and backups



(Q7) Design the data warehouse dimensional model for a wholesale furniture company.

=>

Furniture

id_furniture
type
category
material

Customer

id_customer
age
gender
city
region
state

Sales

id_sales
id_furniture
id_customer
id_time
Quantity
income
discount

Time

id_time
day
month
year

Star Schema

for
wholesale furniture company

(8) Design a star schema for company sales with three dimensions: location, time & item.

Location

location_id

country

state

city

Sales

sales_id

location_id

item_id

time_id

quantity
amount

Item

item_id

name

category

price

manufacturer

5000

hi res

smash

crush

300

Time

time_id

date

day

month

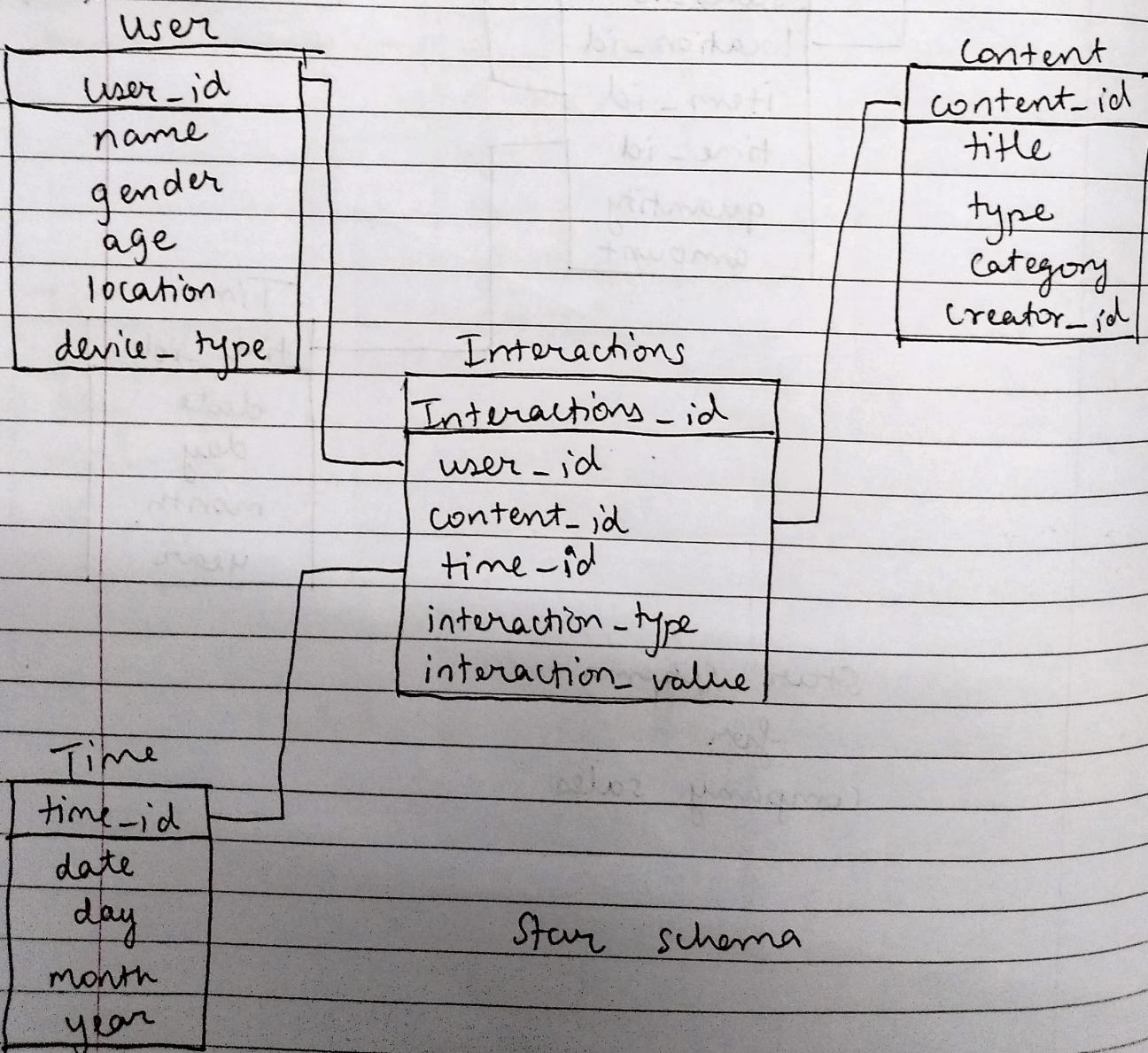
year

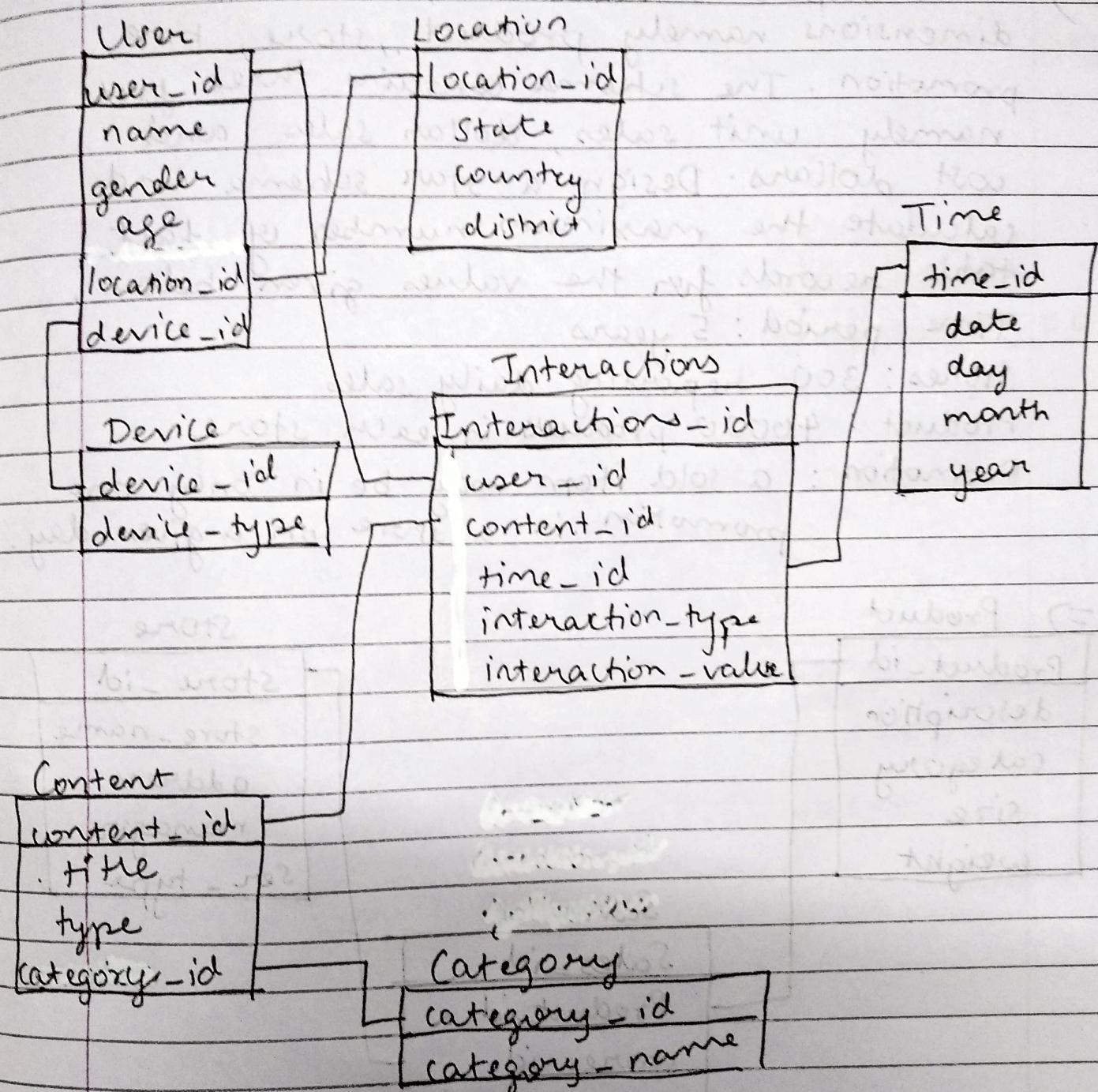
Star Schema

for
company sales

(Q9) A social media platform wants to analyze user engagement data to improve content recommendations and user experience. Design a star schema and snowflake schema for the interactions fact table, which contain details about user interactions, including interaction details, user information, content details and time periods.

⇒





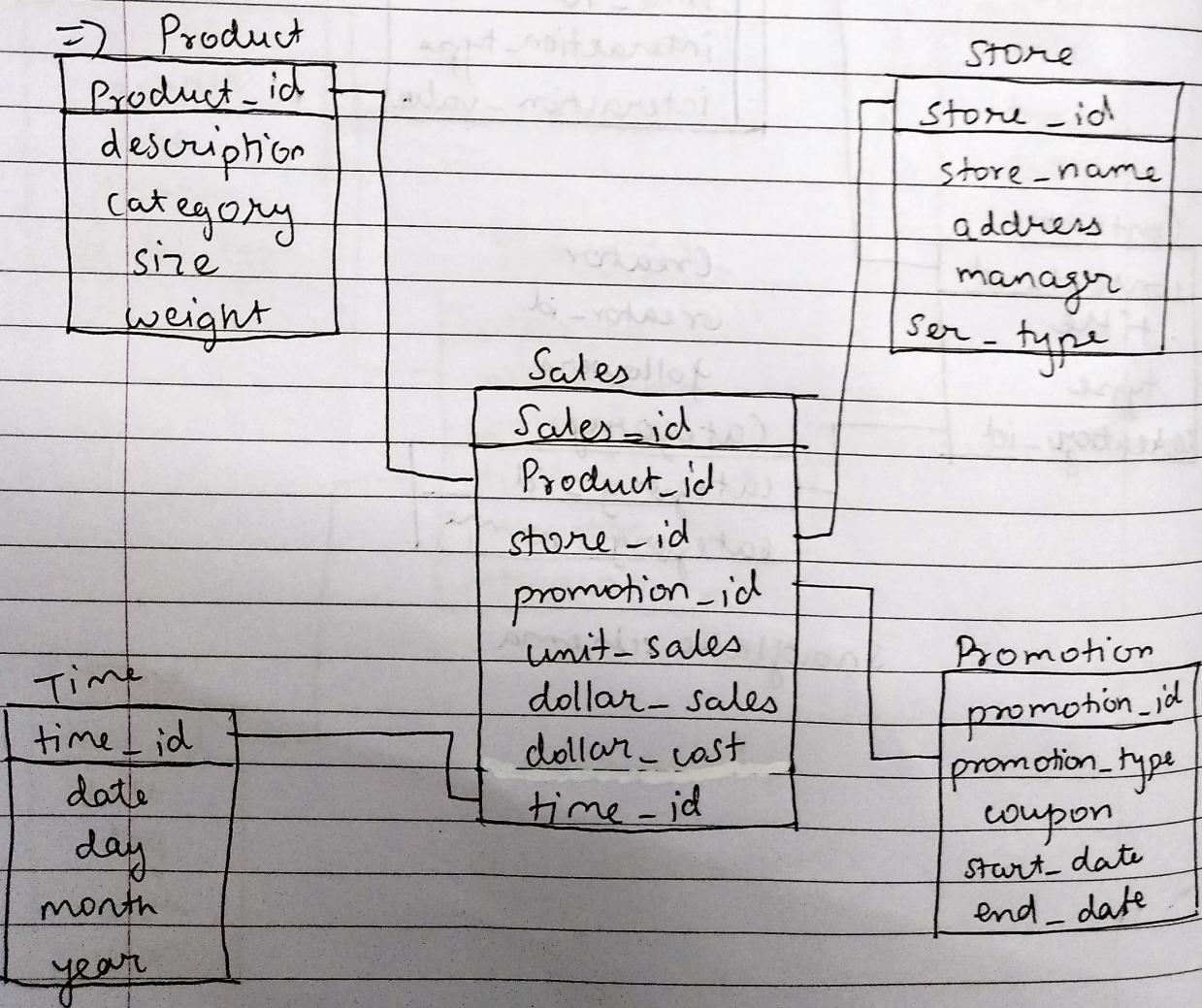
(Q10) For a supermarket chain, consider the dimensions namely product, store, time, promotion. The schema contains three facts namely unit sales, dollar sales and cost dollars. Design a star schema and calculate the maximum number of fact table records for the values given below:

Time period : 5 years

Stores : 300 reporting daily sales

Product : 40000 products in each store

Promotion : A sold item may be in only one promotion in a store on a given day.



$$\text{Time period} = 5 \text{ years} \times 365 \text{ days} = 1825 \text{ days}$$

There are 300 stores

Each store daily sale 4000

Promotion = 1

Maximum number of fact table records

$$= 1825 \times 300 \times 4000 \times 1 = 2,19,00,00,000$$