

# Database Design & MySQL

What is a data warehouse?

→ A company has 3 departments → Marketing, Sales, Finance

Each department maintains a separate database

This could lead to a situation wherein each department has its own version of the facts, for a query as

"What is the total revenue of the last quarter?"

Every department might have a diff ans.

∴ each department draws info from a separate database

Hence the data warehouse can prove to be useful.

A data warehouse would thus be the central repo of the data of the entire enterprise.

Definition : A data warehouse is a centralized repo where large amount of structured data from multiple sources are stored and managed to support data analysis, reporting etc.

A data warehouse is

- ① Subject-oriented
- ② Integrated - repository
- ③ Non Volatile
- ④ Time Variant

Collection of data in support of management decisions.

So the Data warehouse is a collection of Data.  
It has following properties:

⇒ **Subject Oriented**: A data warehouse should contain information about a few well-defined subjects rather than enterprise.

⇒ **Integrated**: A data warehouse is the integrated repository of the data, it contains information from various systems within organization.

⇒ **Non-Volatile**: the data values in database cannot be changed without a valid reason.

⇒ **Time Variant**: A data warehouse contains historical data for analysis.

### Structure of Data Warehouse

One of the primary methods of designing a Data warehouse is dimensional modelling.

The two key elements of dimensional modelling are facts & dimensions.

In short "Dimensional modelling is a popular way to design data warehouses, aiming to make data easy to retrieve and analyze."

Facts → Numerical Data

Dimensions → Metadata (data explaining some other data)

Dimensional Modelling Organizes the data  
into two main parts:

1) Facts : facts are the numbers or measurements  
you want to analyze, they usually  
represent something measurable like  
Sales amount, quantity sold or total revenue.

→ facts are stored in fact table where each  
row represents an event or transaction

e.g. In sales database, "Total Sales Amount",  
"Quantity Sold" are facts because they provide  
measurable data about each sale.

2) Dimensions : dimensions are the descriptive  
details about the facts, they provide the  
context so we can understand the numbers  
better and answers questions like "Who", "What",  
"Where" & "When"

→ Dimensions are stored in dimensions table,  
where each row describes an attribute (like  
product or time period)

e.g. In the same sales database

→ Product Dimension : describe details  
about each product, like prod  
name, category, brand.

→ Customer Dimension : describe customers,  
including customer name, location,  
age

## Putting It all together

A fact table links to dimension table so we can analyze facts with context.

for eg you might have a fact table that records each sale with fields like

→ Sales Amount (fact)

→ Quantity Sold (fact)

→ Links to dimensions like ProductID, CustomerID, DateID

Using these links, we could answer questions like

→ What was the total sales amount for Product X in 2023.

e.g In BANK

Facts {  
    → Withdrawl amount  
    → Account balance after withdrawl  
    → Transaction charge amount  
    → Customer ID  
    → ATM ID  
    → Date of withdrawl

\* A schema diagram in a data warehouse is a visual representation of structure and organization of data within a warehouse. It shows the relationships b/w different tables, columns, keys and other database objects.

There are commonly 3 types of Schema models

In Data warehouse:

① Star Schema: It is a simple and widely used Schema, a central fact table links to multiple dimension tables, resembling a Star shape.

e.g

### E-commerce

Date of Purchase	Product	Customer	Quantity
05-01-2017	Pizza	John	2
06-01-2017	Coke	John	5
-	-	-	-
-	-	-	-
-	-	-	-



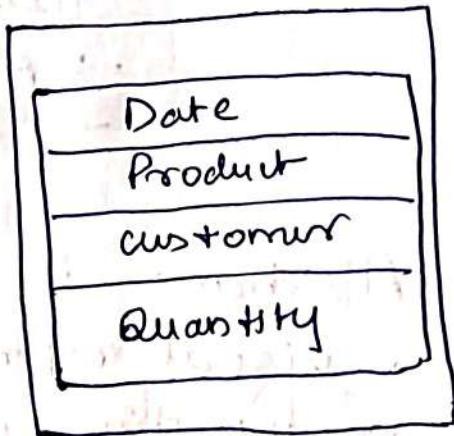
### Schema Diagram

Now suppose we have few questions?

① What do people from Bangalore purchase in general?

② In general what's the soft drink sales?

These info are not there



So star schema is useful for augmenting our main data with the additional information useful for analysis.

Dimension Table

"Date"
Day
Holiday
-
-

Dimension Table

"customers"
Name
city
category
-
-

"Date"
"Product"
"customers"
quantity
-

Main Schema  
(fact Table)

Dimension Table

"Product"
"Name"
"category"
-
-

Now we can use every dimensional table for analysis.

Benefits of Star Schema

- Queries use very simple joins while retrieving the data and thus by query performance increases.
- It is simple to retrieve data for reporting at any point of time for any period.

## Disadvantages of Star Schema

- ⇒ If there are many changes in the requirement, the existing star schema is not recommended.
- ⇒ Data Redundancy is more as tables are not hierarchically divided.

## Snowflake Schema

Snowflake Schema organizes data into a multi-dimensional model, unlike the simple star schema whose dimension table are typically denormalized

(ie dimension table consist redundant data i.e. duplications of data across tables or records. When data is redundant, the same piece of info is stored in multiple places), the snowflake schema normalize these tables, creating a more complex Branched structure

Dimension Table

Time Key
Day
Month
Years

DIT

Branches key
Branch name
Branch type

Fact Table

Time-key
item-key
branch-key
location-key
Dollars sold

Sales

item-key

item-name
Brand
Type
Supplier-key

DT

Supplier-key
Supplier-type

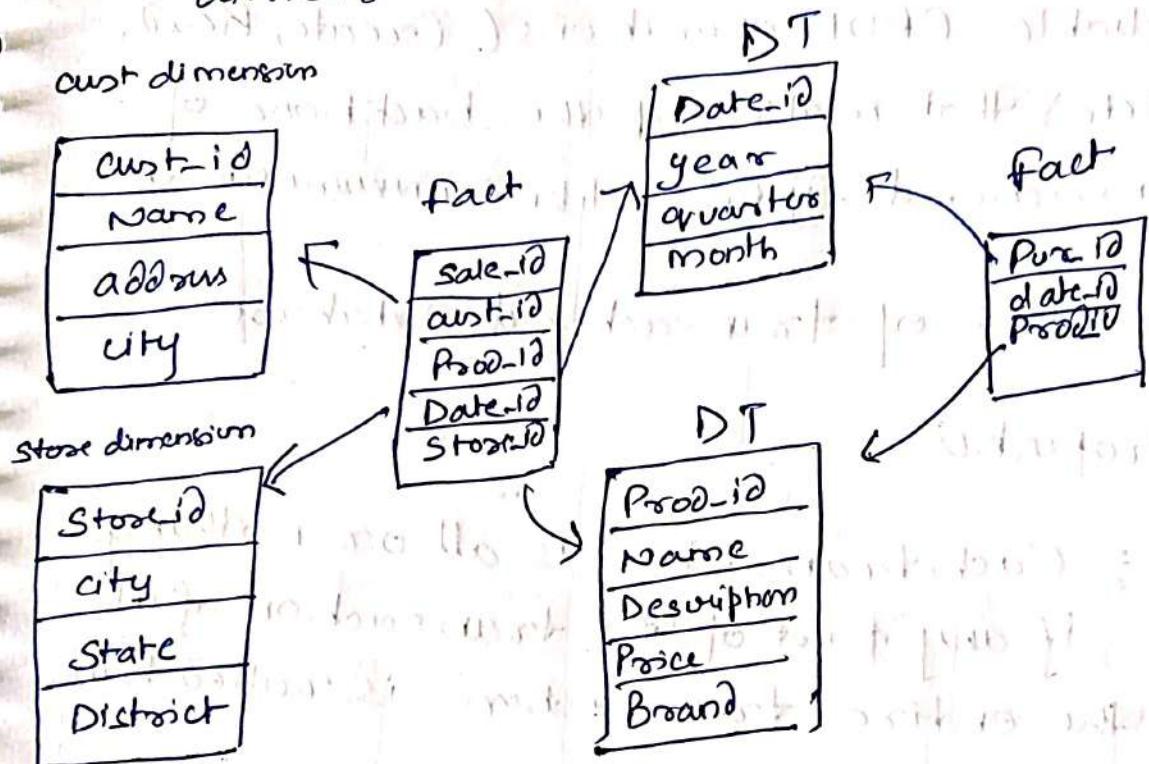
DT

Location key
state
city-key

city-key
city
state

DT → Dimension Table

 Galaxy Schema : Also called as fact constellation.  
It is used when multiple fact tables share dimension tables.



## OLAP vs OLTP

Now we will learn what features of Data warehouses differentiate it from a regular database.

**The Wrong Question:**  
Diff b/w Datawarehouse  
and Database

Database is the collection of data, also the data warehouse is a collector of data.

**The Right Question:**  
Diff b/w DataWarehouse &  
transaction Database

Transactional Database helps users carry out day to day business transactions.

## Transactional Database

It is a type of database designed to handle real time transactions efficiently. It's optimized for performing fast and reliable CRUD operations (Create, Read, Update, Delete) that make's up the backbone of online transactional systems like e-commerce sites etc.

### Key characteristics of transactional databases

#### ① ACID Properties

Atomicity : Each transaction is all or nothing. If any part of the transaction fails, the entire transaction is rolled back.

Consistency : transactions move the database from one valid state to other maintaining predefined rules such as constraints.

Isolation : transactions are isolated from one other preventing concurrent transactions from interfering and leading to incorrect data.

Durability : Once a transaction is committed, the data is permanently saved, even in the event of system crash.

#### Eg's of Transactional Databases

- Relational Databases : MySQL, Postgres, Oracle
- NoSQL Databases : MongoDB

# Diff b/w Transactional Database & Data Warehouse

Level	<u>Transactional DB</u>	Data Warehouse
Purpose	1. to <del>transform</del> perform business transactions.	1. To conduct analysis and decision making.
Users	2. End users (like us) for business transactions	2. Used by decision makers like Managers, CEOs etc
Technology	3. NO ETL	3. ETL process, Extract transform and load process.
Front end Tech	1) Web servers 2) Browsers for interaction 3) Mobile Phones 4) Integration with Payment gateway 5) NO BI	<p style="text-align: center;">Data Warehouse</p> <pre> graph TD     DW[Data Warehouse] --&gt; Integrated[Integrated]     DW --&gt; MultipleSources[Multiple sources]     Integrated --&gt; MultipleDatabases[Multiple Databases]     MultipleSources --&gt; MD[Multiple databases]   </pre> <p>1. BI tools Business Intelligence</p>
Database Design	1) Based on Integrity rules and normalisation 2) No star schema 3) Nature of Access : CRUD (Create, Retrive, Update, Delete)	1) Based on dimensional modelling 2) Creation of Star schema 3) Nature of Access : Retrievive
Also known as	OLTP : Online transaction Processing System	OLAP : Online analytical processing System

## ETL Process Data Warehouse

↳ used to ingest data onto schema and perform ops.

Select : Identify the data relevant for analysis.  
It involves determining which datasets, tables or fields are necessary to achieve specific goals of the analysis.

Extract : Connecting to the particular data source and pulling out the data.

Transform : Modifying the extracted data to standardise it.

Load : Pushing the data into Data warehouse.

## Entity Constraints

Constraints are the rules that are used in MySQL to restrict the values that can be stored in the columns of a database.

This ensures data integrity, which is nothing but the accuracy and consistency of the data stored in the database.

### ① Not Null constraint

- ensures that the column cannot contain a null value.
- used when a value is required for every row in a column (eg Primary key fields, crucial fields like email in users table)

g Create table students ( student\_id int NOT NULL, name varchar(50) NOT NULL );

## 2. UNIQUE constraint

- ensures all values in a column are unique across the table.
- commonly used for columns like email or username where duplicate values aren't allowed.

eg Create table employees ( emp\_id int unique, email varchar(100) unique );

Note: the UNIQUE constraint does allow null values, it won't violate the uniqueness rule because null values are treated differently from regular data.

## 3. Primary Key constraint

- combination of not null and unique, it uniquely identifies each row in a table.
- A table can have only one primary key which can be a single column or a combination of columns (Composite primary key)

eg Create table orders (

order\_id int Primary Key, order\_date DATE NOT NULL

);

#### 4. Foreign Key Constraint

- enforces a link b/w two tables by referencing a primary key in another table.
- Maintains referential Integrity by ensuring that values in a column or set of cols. match values in the referenced table.

eg Create table Orders (

```
    Order-ID int Primary Key,  
    Customer-ID int,  
    Foreign Key (Customer-ID) References Customers  
        (Customer-ID)
```

#### 5. CHECK Constraint

- ensures that all the values in col meet a specified condition.
- useful for setting range limits or enforcing rules on dataformats.

eg Create table Product (

```
    Product-ID int Primary Key,  
    Quantity int check (Quantity >= 0)
```

);

## 6) Default constraint

→ Sets a default value for a col, if no value is provided

Ex: Create table users (

```
user_id int Primary key,
registration_date DATE DEFAULT
current_date
```

);

## Referential constraints

entity type constraints pertain to the value in a single table.

the second type of constraint is called referential constraint, they are used to restrict the values that are taken by a column in one table based on the values that exist in another table.

We have 2 tables

customers

cust_id	firstname	lastname	Age
1	John	Wick	52
2	Sheldon	Cooper	41
3	Charlie	Harper	45

orders

order_id	order_num	cust_id
1	1234	2
2	3142	2
3	1246	3
4	4567	1

Cust-ID is a foreign key in the "Orders" table, because it is used to reference the customers table.

We can easily determine which customer placed that particular order and find out all the details of that customer by looking up the corresponding cust-ID in customers table.

referential Integrity constraint

The diagram illustrates a referential integrity constraint between two tables. On the left, the 'Employees' table has columns: emp-id, name, and dept-ID. It contains two rows: E01 (name N001, dept-ID D04) and E02 (name N002, dept-ID D01). An arrow points from the dept-ID column of the Employees table to the dept-ID column of the 'Departments' table on the right. The 'Departments' table has columns: Dept-ID and Dept-name. It contains two rows: D01 (Dept-name Finance) and D04 (Dept-name HR).

emp-id	name	dept-ID	Dept-ID	Dept-name
E01	N001	D04	D01	Finance
E02	N002	D01	D04	HR

According to the referential constraint b/w two tables the value that appears as a foreign key in a table is valid only if it also appears as a primary key in the table to which it appears.

Note: foreign key need not be unique, many employees can belong to the same department.

A table can have any number of foreign key.

### Semantic Constraints

Adds rules based on the specific meaning or business rules of the data, for ex a constraint ensuring an "age" col has values greater than zero, or Salary field is within a range.

# In Summary

## ① Entity Constraint

→ Primary Key constraint : Ensures each row in a table is unique and not null, making it a unique identifier for rows.

→ Unique Constraint : Guarantees that all values in a column are unique.

→ NOT NULL Constraint : Prevents null values in a specific col, ensuring that a value is always provided.

## ② Referential Constraint

→ Foreign key constraint : Enforces a link b/w tables by requiring that the values in a foreign key col correspond to the values in the primary key col of another table.

→ On Delete Cascade : If the referenced row is deleted, any rows referencing it are automatically deleted, maintaining referential integrity.

e.g. 2 tables : Department & Employee

Department table

↳ Dept-ID (primary key)

Dept-name

Employee table

↳ emp-ID (Primary key)

emp-name

Dept-ID (Foreign key)

Now if we set On Delete cascade on Employee. Department ID of means when a row in Dept table is deleted, all rows in employee table

that reference Dept-ID will also be deleted.

eg Create table order-details

order-id int primary key,

customer\_id int References customer-details

on delete cascade,

Order- Date Date

);

→ On update cascade

⇒ If primary key is updated, all related foreign keys are automatically updated to match the value.

### 3. Semantic constraints

⇒ Check constraint: Enforces a specific condition on the data within a col

eg age > zero

eg Create table employee (

eid int(11) NOT NULL,

name varchar(40) Default Null

Date-of-joining DATE NOT NULL CHECK

(Date-of-joining > '2019-02-01')

Primary key (eid)

);

## ERDs

An entity-relationship diagram or ERD can be thought of as a map of the database schema.

We can visualise the structure of the entire schema and answer the following questions just by looking at it.

- What are the tables that it contains?
- What are the columns that each table contains?
- What are the datatypes and constraints for each column?
- What are the relationships b/w various tables?

**Entity:** Small piece of information that we want to track on a database.

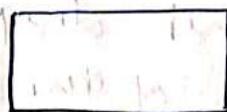
**Entity type:** Collective name given to the bunch of entities

## Why use ER Diagrams in DBMS

- ER Diagram represents the ER model in the database making them easy to convert onto relations (tables)
- ER diag requires no technical knowledge and no hardware support
- It gives standard way of visualizing the data logically.

# Symbols used In ER models

Rectangle



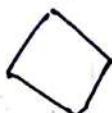
Entities In  
ER Model

Ellipse



Attributes In  
ER model

Diamond



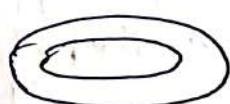
Relationship  
among entities

Line



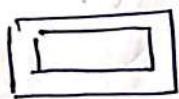
Attributes to entities  
and entity sets with  
other relation type

Double  
ellipse



Multi-valued  
Attributes

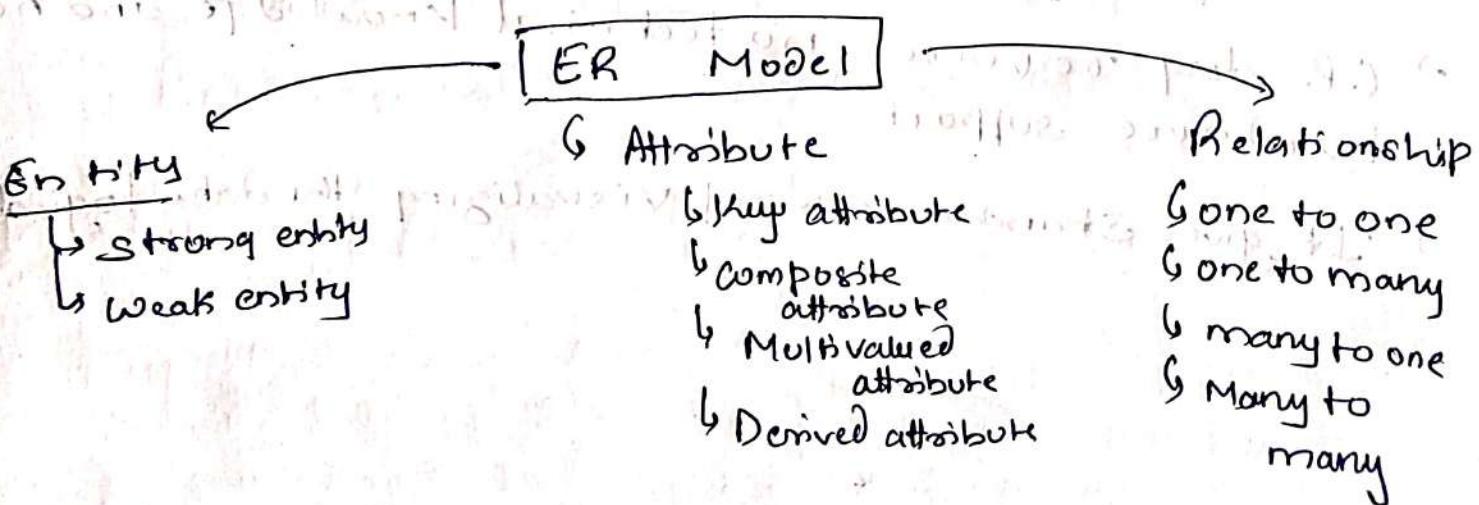
Double  
Rectangle



Weak  
entity

## Components of ER Diagrams

ER model consist of Entities, Attributes, Relationships,  
among entities In a Database system.



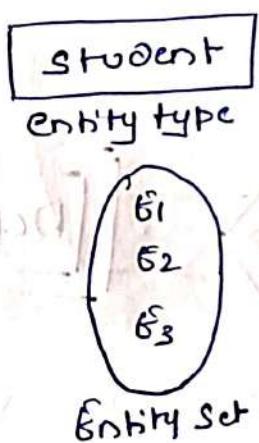
## What is Entity?

An entity may be an object with a physical existence - a particular person, car, house or employee  
or it may be an object with a conceptual existence, a company, a job or a university course.

## What is entity-set?

An entity is an object of entity type and the set of all the entities is called as entity set.

In ER diagrams, entity type is represented as

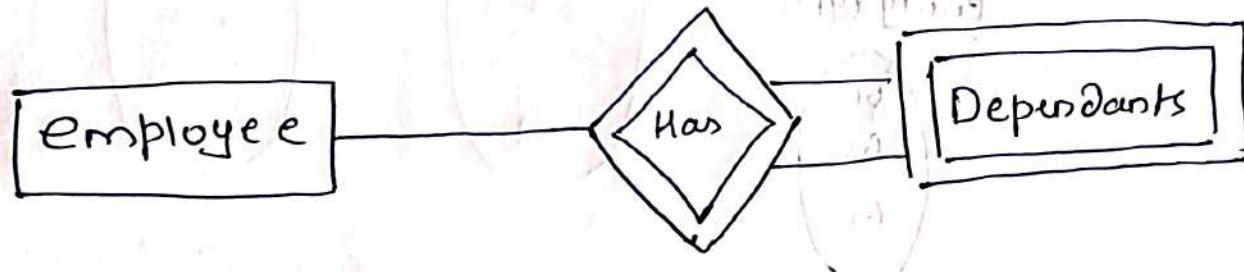


We have 2 types of entities

- ① A strong entity is a type of entity that has a key attribute. Strong Entity does not depend on other entity in the schema. If has a primary key, that helps in identifying it uniquely, and it is represented by a rectangle, they are called strong entity types.

A weak entity type is an entity that has a key attribute that uniquely identifies each key in the entity set, but some entity type exists for which key attributes can be defined.

eg A company stores the information of dependents (Parent, children, spouse) of an employee, but the dependent's can't exist without an employee. So the dependent will be a weak entity type and Employee will be identifying entity entity type for Dependent which means employee is a strong entity.



### Attributes

These are the properties that defines the entity type.  
eg Roll-no, name, DOB are the attributes that define entity student

Represented as : Attribute



### Types of Attributes

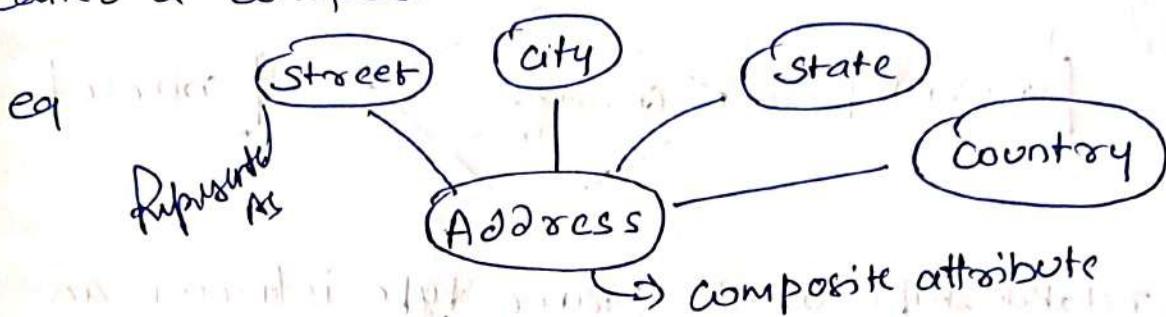
- 1) Key attribute : The attribute which uniquely identifies each entity in the entity set is called key attribute, for eg roll-no will be unique for each student

Represented as



## 2. Composite attribute

An attribute composed of many other attributes is called a composite attribute.



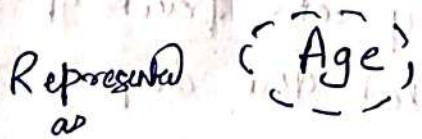
## 3. Multivalued Attribute

An attribute consisting of more than one value for an entity. Eg Phone-no (People contains more than 1)

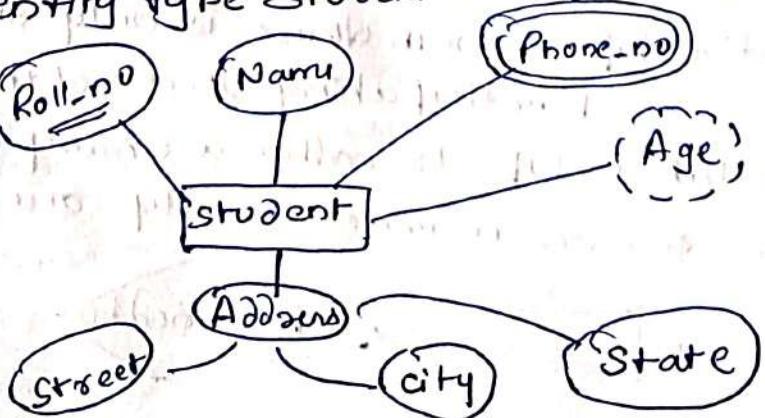


## 4. Derived Attribute

An attribute that can be derived from other attributes of the entity type is known as derived attribute.



Complete entity type student and its attributes



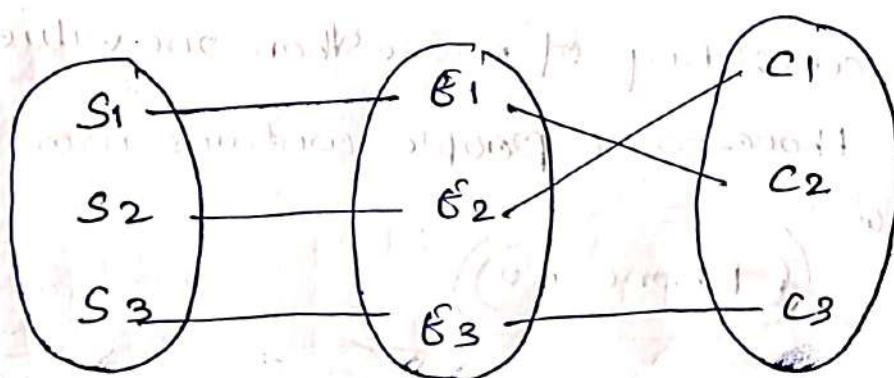
## Relationship Type & Relationship Set

A relationship type represents the association b/w entity types. For eg "enrolled in" is a relationship type that exists b/w entity type student and course.

Represented As



A set of relationships of the same type is known as a relationship set.

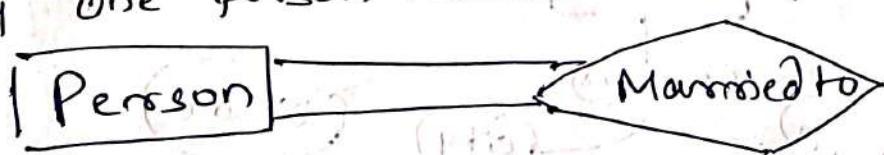


S1 enrolled in C2, S2 enrolled in C1, S3 enrolled in C3

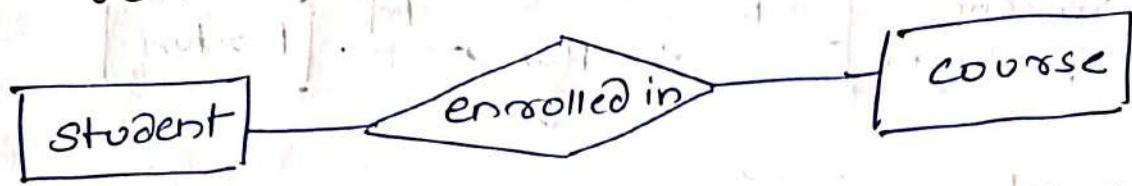
### Degree of a relationship set

↪ The nos of different entity sets participating in a relationship set is called the degree of a relationship.

- ① Unary Relationship : When there is only one entity set participating in a relation, the relationship is called a unary relationship.  
eg One person married to only one person



2. Binary Relationship : When there are two entity sets participating in a relationship, the relationship is called as Binary Relationship.



3. Ternary Relationship : Where there are 3 entity sets participating in a relationship.

4. Many relationship : When there are n entity set participating in a relationship.

### What is Cardinality

The number of times an entity of an entity set participating in a relationship is known as cardinality.

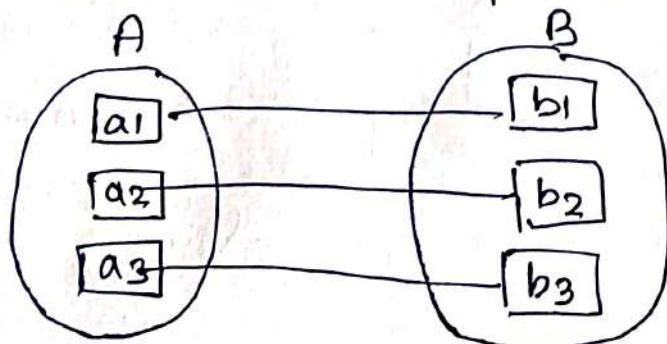
Different types of cardinality are

① One-to-One : When each entity in each entity set can take part only once in the relationship (one to one)

e.g. Male can marry one female and vice versa

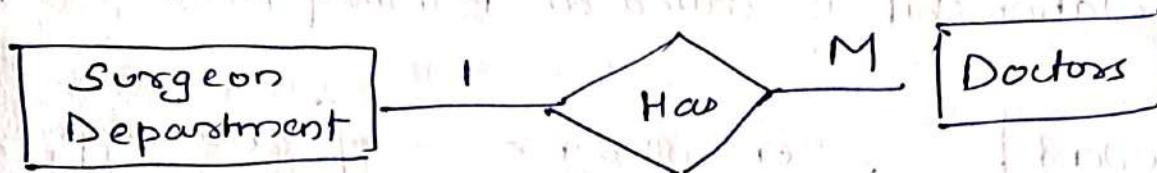


Using sets can be represented as

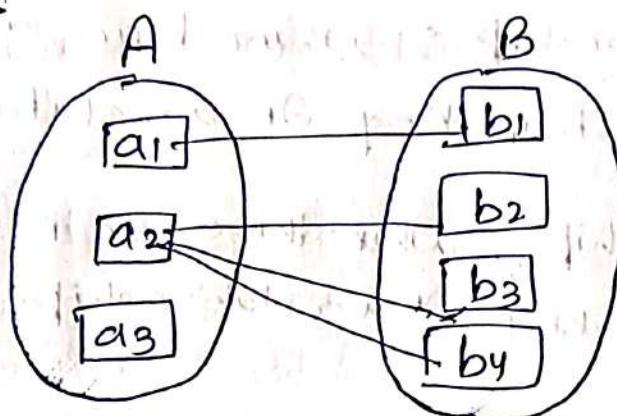


2. One to Many : each entity can be related to more than one entity.

e.g.



Using sets

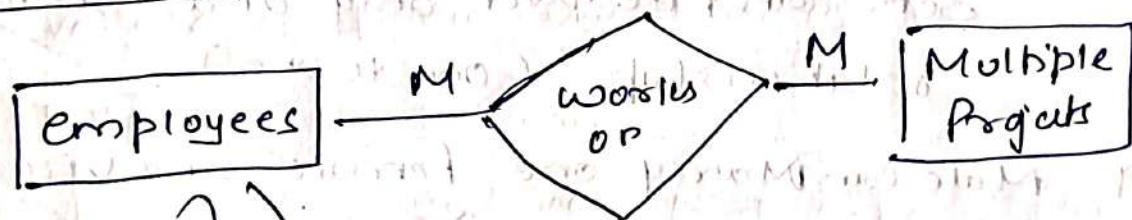


3. Many to one :

e.g.



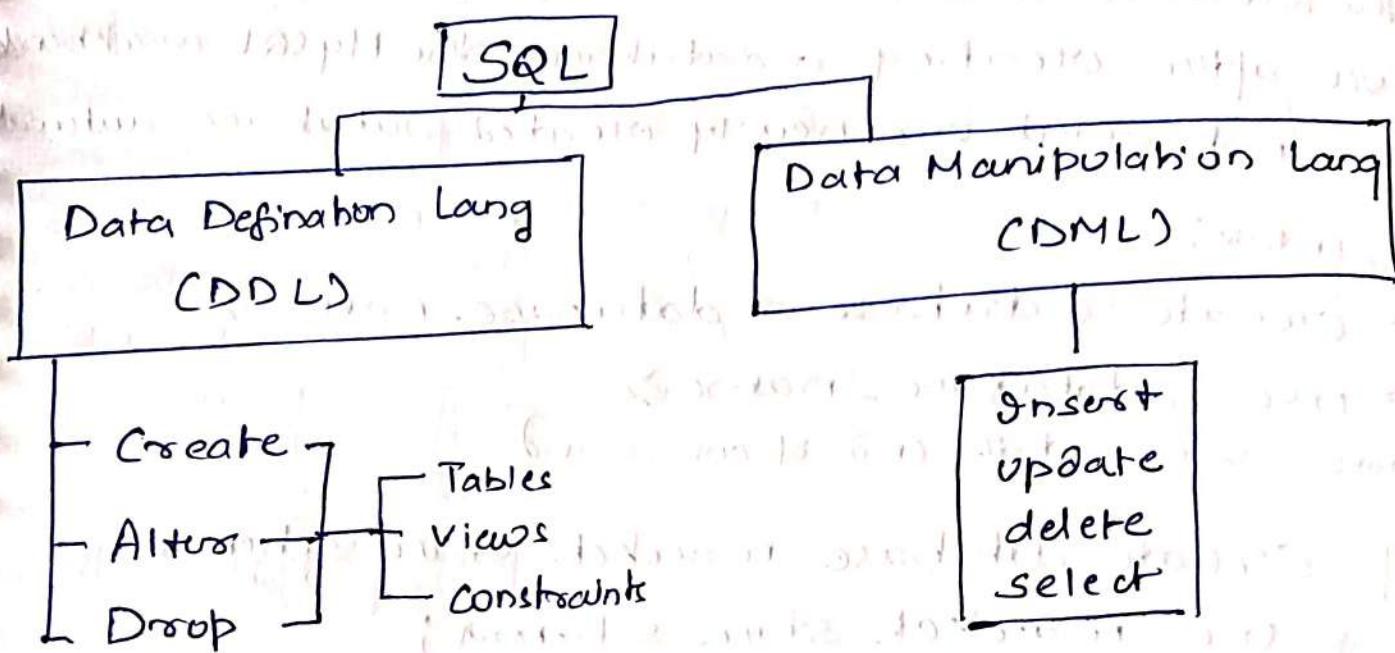
4. Many to Many :



(~~the door password~~)  
MySQL Workbench

# Introduction to DDL and DML Statements

SQL = Structured Query Language



So the commands available in SQL can be broadly categorised as follows:

⇒ Data Definition Language (DDL)

⇒ Data Manipulation Language (DML)

DDL as the name suggests is used to create a new schema as well as modifying an existing schema. The typical commands available in DDL include Create, Alter, Drop.

As a data scientist mostly we will be working with DML commands to generate insights.

## DDL Statements : Demonstration

Now we will learn how to create and use a database (also called schema), we need to specify the database that we need to use, this is because even after creating a database, the MySQL workbench does not select our newly created database automatically.

Syntax:

- Create database < database-name >
- use < database-name >

Note: use ; at the end of command

Eg → Create database market\_star\_schema;  
→ use market\_star\_schema;

# Creating tables and developing star-schema

Create table shipping\_mode\_dimen (

ship-mode varchar(25),

Vehicle-company varchar(25),

Toll Required boolean

); # Varchar : Data type in MySQL

which corresponds to string attribute.

Variable character

String attribute.

length → Here 25 refers to the maximum length it can take.

If we want to store fixed length → char(9) → Max 9 characters

# Alter is a command useful for modifying any schema object.

# We will add a primary key constraint to an existing table.

> alter table shipping-mode-dimensions add constraint Primary Key (ship-mode);

eg create table india (

matches\_played int,

matches\_won int,

matches\_lost int,

net-run-rate Decimal (4,3)

);

Note : In DECIMAL (4,3) → 4 is the total nos of digits and 3 is the total nos of digits after decimal point

## DML Statements : Demonstration

# Insert Command

insert into shipping-mode-dimension

values

('Delivery Truck', 'Ashok Leyland', false),

('Regular Air', 'Air India', false));

↳ Here we are not specifying the order of the attribute values.

Other format to insert the value, better way would be

Insert into shipping-mode-dimensions (Ship-mode)  
Vehicle-company, Toll-required)

Values

('Delivery Truck', 'Ashok Leyland', false),

('Regular Air', 'Air India', false);

Here we can mention the order of the attributes

Now we will learn about "update-set-where"

& "delete-from-where" commands

# update command is used to modify data in the database.

# In the above inserts we set toll-required as false for Delivery Truck we need to update it as true

update shipping-mode-dimensions set Toll-required = true  
where Ship-mode = 'Delivery Truck';

# Delete AirIndia entry

Delete from shipping-mode-dimensions where  
vehicle-company = 'Air India';

4 delete from tablename where condition;  
4 update tablename set col1 = val1 where condition;  
(syntax)

Notes: A value that is a primary key in a table cannot be deleted, if another table has the same value as a foreign key.

So it's necessary to sequence the delete commands.

## Modifying Columns

# Alter command to add a new column

```
alter table shipping-mode-dimen add vehicle-num  
varchar(20);
```

ie alter table <table-name> add <col-name>  
<datatype>;

Now this new col will have null values for all  
the existing rows.

```
> update shipping-mode-dimen set vehicle-number  
= "MH-05-R1234";
```

# Dropping a column

```
Alter table <tablename> drop column <col-name>
```

```
Alter table shipping-mode-dimen drop column  
Vehicle-Numbers;
```

↳ # this col is gone forever.

# Change the col name "Toll-Required" to "Toll-Amount" also change datatype to int

> alter table shipping-mode-dimen

change TollRequired Toll-Amount int;

C keyword

# for dropping the table

" drop table <tablename>"

> drop table shipping-mode-dimen;

Note : ① Create, alter, drop, rename, truncate are DDL commands

② Select, Insert, update, delete are DML.

Note : Drop removes the table completely from the database, we can't retrieve it back, all the Integrity constraints are removed.

Truncate only drops the table rows, it doesn't delete the schema from the database.

## Querying in MySQL

SQL statements and Operators

→ In Database  
Market\_star\_schema

# to print the entire data of all the customers

Select \* from cust-dimen;

↳ tablename

# list names of all the customers

> Select customer\_name from cust\_dimen;

# Print name of all customers along with their city and state

> Select customer\_name, city, state from cust\_dimen;

# Note : In general the data is shown in the order in which it was inserted into the table.

# we have a set of functions that we can use in SQL like for eg Aggregate functions, these functions cumilate the data in multiple ways.

eq # Print the total nos of customers

> Select count(\*) as total\_customers from cust\_dimen;

total-customers  
1883

( Note : Here we used "Count(\*) as total\_customers"

This is called as, Aliasing

$\equiv$

eq for Aliasing

Select customer\_name AS "Customer Full Name",

city AS "Customer City" from cust\_dimen;

Customer Full Name Customer City

-

-

-

Note : the alias is limited to that specific query only.

# we can also apply filter conditions

Ex # Show details of the customers belonging to West Bengal.

> Select \* from cust\_dimen where state = "West Bengal";

# Similarly, we can write multiple conditions

> Select \* from cust\_dimen where state = "West Bengal" and city = "Kolkata";

Note: filtering is always evaluated first.

Ex 2 How many customers are from West Bengal

> Select count(\*) as Bengal\_customers from cust\_dimen where state = "West Bengal";

Ex 3: Print the names of all customers who are either corporate or belong to Mumbai

> Select customer\_name, city, customer\_segment from cust\_dimen where city = 'Mumbai' or customer\_segment = 'Corporate';

# List the details of all the customers from southern India, namely Tamil Nadu, Karnataka, Telangana, Kerala

Select \* from cust\_dimen where state in ('Tamil Nadu', 'Karnataka', 'Telangana', 'Kerala');

# Point the details of all non-small business customers

Select \* from cust-dimen where customer\_segment != "Small Business";

G = Negation operator (Not equal to)

# Dealing with numeric operators

e.g# list the order\_ids of all those orders which caused losses.

Select order\_id, Profit from market\_fact-full where Profit < 0;

e.g# List the orders with '-5' in their order\_ids and shipping costs BETWEEN 10 and 15. (Like)

Select order\_id, shipping\_cost from market\_fact-full where order\_id like '%\\_-5%' and shipping\_cost between 10 and 15;

e.g# Write a query to display cities in the cust-dimen table that begins with letter K

Select city from cust-dimen where city like 'K%';

Note K% → cities starts with K

%K → returns cities ending with K.

# Between operators e.g

Select \* from market\_fact-full where shipping\_cost between 10 and 15;

The following wildcards are supported in SQL

→ Both % and \_ (underscore)

% is a multi-character wildcard, whereas  
\_ is a single character wildcard.

e.g. "\_%%" → returns words with second character  
as % followed by zero or more  
characters at the end.

Note: In Between both are Inclusive

## Aggregate functions

SQL Provides several aggregate func, each  
designed for specific data analysis tasks,  
here are the most commonly used aggregate func  
in SQL.

- ① **SUM()** → calculates the total sum of a numeric column.
- ② **Avg()** → calculates the avg (mean) of a numeric col.
- ③ **Count()** → count the nos of rows or non null values in a column
- ④ **Min()** → find the minimum value in a column
- ⑤ **Max()** → find the maximum value in a column
- ⑥ **Group\_concat() or string\_agg()**, depending on SQL version → Concatenates values in a col into a single strng.

## Example of Aggregate func with Group By

eq table : employees with following columns

- employee\_id
- department\_id
- salary
- job\_title

### 1. SUM() - calculate total

to find the total salary by department

> Select department\_id, SUM(salary) as total\_salary  
from employees group by department\_id;

### 2. AVG() - calculate Average

to find Avg salary by job-title

> Select job\_title, AVG(salary) as average\_salary  
from employees group by job\_title;

### 3. COUNT() → Count Rows or Values

eq count the total nos. of employees in each department

> Select department\_id, COUNT(\*) as employee\_count  
from employees group by department\_id

Eq 2 of count()

To count the non-null values in a specific col

Select department\_id, count(salary) as non-null-salaries from employees group by department\_id;

This counts only the rows where salary is not null for each department.

4. MIN() → to find minimum value

To find the lowest salary in each dept

Select department\_id, min(salary) as minimum\_sal from employees group by department\_id;

5) MAX() → to find the maximum value

To find the max salary in each department

Select department\_id, max(salary) as max\_sal from employees group by department\_id;

Using GROUP BY with Multiple Columns

Eq. to calculate the avg salary by both department and job title

Select department\_id, job\_title, avg(salary) as avg\_salary from employees group by department\_id, job\_title;

## Group- Concat()

Suppose we have a table employees with 2 cols department-id, employee-name. We want to list all the employee names in each department separated by commas.

### 1) Group-Concat()

The group-concat() in MySQL are used to combine values from multiple rows within each group into a single concatenated string.

> Select department-id, group\_concat(employee-name) as employee-names from employees group by department-id

group department-id	employee-names
1	Alice, Bob, Charlie
2	Dave, Emma, Frank

## GroupBy and Having

The groupby and Having clauses are commonly used together for grouping data and filtering grouped results.

① Group By : This clause groups rows that have the same values in a specified col, Aggregate func (like SUM, AVG, COUNT etc)

② Having : this clause is used to filter groups based on aggregate values.

While the "where" clause filters rows before grouping, "Having" filters after the grouping and aggregation have completed.

Ex table : Sales

Assume we have a table Sales with following cols

→ Prod-ID , → category , → Sales-amount , → Sales-date

# Suppose you want to find the total sales of each Product

we will use Groupby with sum func

> Select prod-ID , SUM(sales-amount) as total-sales from sales groupby prod-ID

# If you only want to see products with total sales above 10000

Select prod-ID , SUM(sales-amount) As total-sales from sales group by prod-ID

Having SUM(sales-amount) > 10000;

## Ordering

In SQL the `order by` clause is used to sort the result set based on one or more columns; we can sort the data in ascending / descending order.

Syntax:

```
Select col1, col2, ...
  from Table-name
  order by col1 [ASC|DESC], col2 [ASC|DESC], ...
```

By default `order by` sorts data in ascending order

e.g. table  $\rightarrow$  employees with cols as

$\rightarrow$  employee\_id

$\rightarrow$  name

$\rightarrow$  salary

$\rightarrow$  department

e.g. # list of employees sorted by salary in descending

```
> Select name, salary from employees order by
  salary Desc;
```

# for ascending

```
> Select name, salary from employees order by
  salary;
```

## Order by Multiple Columns

We can sort by multiple columns, for eg we want to sort by department in ascending order and then by salary in descending order within each department

> Select name, department, salary from employees  
order by department Asc, salary desc;

## Order by Aggregate functions

When using aggregate func and groupby we can use order by to sort the result based on the aggregated values

> Select department, sum(salary) as total\_salary  
from employees group by department  
order by total\_salary Desc;

## Order By with Aliases

If we use an alias for a column, we can use the alias in the order by clause

> Select name, salary <sup>AS</sup> Emp\_sal from employee  
order by emp\_sal desc;

## Ordering Null Values

SQL standards vary by database on how null values are ordered.

→ Some databases put null values at the beginning of ascending order and the end of descending order.

→ Others allow explicit control over the placement of null values e.g. NULLS FIRST or NULLS LAST  
e.g. In PostgreSQL

Select name, salary from employees order by salary desc, NULLS Last;

e.g. we have Sales table

→ Sale-ID → Product-ID → amount → sale-date

We want to retrieve the total sales amount for each product that has a total sales amt > than \$10000, sort the result by total\_sales in descending order and limit the O/P to top 5 prod

→ Select Product-ID, sum(amount) as total\_sales  
from Sales  
group by Product-ID  
having sum(amount) > 10000 order by  
total\_sales desc limit 5;

# String and Date - Time functions

We have many built-in string functions that help with manipulating text data

1. `concat()` → concatenate strings

→ Combines two or more strings into one.

→ Syntax: `CONCAT(string1, string2, ...)`

eg `Select CONCAT(firstname, ' ', lastname)`  
as full-name from employees;

2. `upper()` → Convert to Uppercase

→ Converts a string to uppercase letters

→ Syntax: `UPPER(string)`

eg `Select upper(name) from products;`

3. `lower()` → Convert to Lowercase

→ Converts a string to lowercase letters

→ Syntax: `LOWER(string)`

eg `Select LOWER(name) from products;`

4. `SUBSTRING()` → Extract Part of a String

→ Extract a portion of a string from a specified position for a specified length.

→ Syntax: `substring(string, start, length)`

eg `Select substring(name, 1, 3) from employees;`  
→ returns 1st 3 letters of the name

5. Length () → Returns the length of a string  
→ Returns the number of characters in a string  
Syntax: LENGTH (string)  
eg Select Length (name) from employees;

6. TRIM () → Removes leading and trailing spaces  
→ Removes white space from the beginning  
and end of the string.  
→ Syntax: TRIM (string)  
eg Select TRIM (name) from employees;

7. Replace () → Replace the Part of String  
→ Replace the occurrence of a substring within  
a string with another substring.  
→ Syntax: Replace (string, old-substring, new-substring)  
eg Select Replace (name, ' ', '-') from products;  
(replaces space with hyphens)

8. POSITION () → Find position of substring  
→ finds the position of the first occurrence of a substring  
within a string

Syntax: Position (SubString IN string)

eg Select Position ('app' IN name) from products

9. REVERSE () → Reverse the string (REVERSE (string))  
eg Select REVERSE (name) from products;

# Date Time func in SQL

1. CURRENT\_DATE : Returns the current Date
2. CURRENT\_TIME : Returns the current time
3. CURRENT\_TIMESTAMP : Returns the current date and time  
eq

Select CURRENT\_DATE;

Select CURRENT\_TIME;

Select CURRENT\_TIMESTAMP;

## 2. GETDATE()

Returns the current date and time (similar to current\_timestamp)

→ Select GETDATE();

## 3. DATEADD (datepart, number, date)

→ Add a specific time interval to date

→ datepart → Part of date to increment  
eq year, month, day

eq

Select DATEADD(day, 5, '2023-01-01');  
#Add 5 days to 2023-01-01

## 4. DATEDIFF (datepart, startdate, enddate)

↳ Returns the diff b/w two dates

```
> Select DATEDIFF(day, '2023-01-01',  
'2023-01-10'),
```

↳ Returns 9 days

5. YEARS(), MONTHS(), DAYS()

↳ Returns specific parts of a date

eq Select YEARC('2023-01-01') j → 2023

SELECT MONTH('2023-01-01') ; → Returns 1

SELECT DAYC('2023-01-01') ; → Returns 1

~~4:~~ SYSDATETIME()

SYSDATE TIME ()  
↳ Returns date and time including fractional seconds

→ SELECT SYSDATETIME()

## Q Employees table

↳ first-name, last-name, Job-title, №

Write a query to retrieve full names of all the employees along with employee no.

eg employee number full name

1002 Diane Murphy

→ Select employeeNumber, CONCAT(firstname, ' ',  
lastname) from employees;

# REGEX → Regular expression

Used for pattern matching or string matching

[abc] → a, b or c

[^abc] → any character except a, b, c

[a-zA-Z] → a to Z

[A-Z] → A to Z

[a-zA-Z] → a to z, A to Z

[0-9] → 0 to 9

→ Quantifiers In Regex

Using quantifiers we can tell the computer about the repetitions.

⇒ [ ]? if ? is there it means  
↳ whatever the rule is present b/w brackets it can occur 0 or 1 times

⇒ [ ]+ If + sign is there it means  
↳ whatever rule is present in b/w brackets it can occur 1 or more times.  
(ie atleast 1 time it has to occur)

⇒ [ ]\* → occurs 0 or more no's of times

⇒ [ ]{n} → occurs n times

[ ] {n,} occurs n or more times

[ ] {y, z} occurs atleast y times but less than z times

## Regex Metacharacters

Meta characters are basically short forms

Instead of writing [0-9] we can write as \d

for [^0-9] → \D

for [a-zA-Z-\d] → \w

and \W is for [^\w]

Note: \ → backwood slash tells computer to treat following characters as search characters, because + sign means different → this means diff, so in order to search

\+ also known as escape character

e.g. mobile number → starts with 8 or 9 and total digits = 10

[89][\underbrace{0-9}] {9}

↓ this has to repeat 9 times

1st digit can be 8 or 9

Q2 first char should be uppercase; it should contain lower case alphabets, only one digit is allowed in b/w.

$[A-Z][a-z]^* [0-9] [a-z]^*$

occurs 0 or more nos of times

1 digit

0 or more nos of times

Q3 Email ID

Q sadya123@gmail.com

1st Part    2nd    3rd

(alphabets +, -, -) these will come with a dot

$[a-zA-Z0-9_ \cdot \backslash \cdot ] + [\text{@}] [a-zA-Z]^+ [\cdot \cdot ] [a-zA-Z]$

(for allowing hyphens we need to use escape character \ )

Small a to z, A to Z, 0 to 9, -, ., @, \_ at least one

others [ @ ] for @ others [a-zA-Z]^+ for

gmail others [ \cdot ] for . others [a-zA-Z]^+ for

for com or in

In SQL we use as

select \* from table\_name where column\_name  
REGEXP 'Pattern';

## Nested Queries

While generating insights from the data we may need to refer to multipletables in a query, there are two ways to deal with such type of queries:

- ① Joins
- ② Nested Queries / Subqueries

e.g. Print the order numbers of the most valuable order (eq maximum sale) by sales

→ Select Good-id, Sales from market-fact-full where  
Sales = (select max(sales) from market-fact-  
full);

Nested Query

Nested Queries are better ∵ we don't need to hardcode anything here creating a generic form of query.

Note: When is a query we are composing it with NULL we are not supposed to use = instead we should use 'is null' or 'is not null.'

for eq :

Select prod-id from market-fact-full where  
prod-base-margin = null ; X

Correct way : Select prod-id from market-fact-full  
where prod-base-margin is null ; ↗

Eq Select \* from prod-dimensions where prod-id  
in ( select prod-id from market-fact-full,  
where prod-base-margin is null );

Note : the subquery is executed first then the main query.

SQL based Nested Queries also known as Subqueries or inner queries are SQL queries that are embedded within another SQL query, they are often used to perform operations where the result of one query depends on the outcome of another query.

### Types of Nested Queries

① Single row subquery : these query return only 1 row or 1 col as output

Select name from employees where salary =

(Select MAX(salary) from employees);

2) Multi-row Subqueries: These return multiple rows and are typically used with operators like 'in', 'any' or 'all'.

e.g. Select name from employees where dept\_id is

(Select dept\_id from departments where loc = "New York");

3) Multi-column Subqueries: These subqueries return more than 1 column.

e.g. Select name, department from employees where

(department\_id, job\_id) in (Select department\_id, job\_id from Job-assignments where active = 1);



## COMMON Table Expressions

These are temporary result sets that you can reference within a select, insert, update or delete statement, they help to break down the complex queries into more readable parts.

Syntax of CTE CTE is defined using "with" clause

With cte-name AS (

select col1, col2

from table-name where condition)

Select \* from cte-name;

eg with AvgSalary AS ( )  
Select dept\_id, Avg(salary) as avg-sal  
from employees GROUP BY department\_id

Select e.name, e.salary from employees e  
Join AvgSalary a ON e.department\_id =  
a.department\_id where e.salary > a.avg-salary

## \* Derived table

Derived table in SQL are subqueries that are used in the From clause of main query they act as temporary tables that exist only during the execution of the query.

Syntax :  
Select dt.col1, dt.col2  
From ( Select col1, col2  
From table-name  
Where condition  
) As dt Where dt.col1 > 100;

Views  $\Rightarrow$  these are the temporary tables which can be used

Definition : Can be thought as a saved query that can be used as a table, however it doesn't hold any physical data

Create view view-name as Select col1, col2 from table where condition;

# Joins & Set Operations

## Set theory

A set is a collection of distinct values. It is denoted by comma-separated values enclosed by curly braces.

e.g. Set A → even nos, Set B → Prime nos

$$A = \{2, 4, 6, 8\}, B = \{2, 3, 5, 7\}$$

the union of A and B ( $A \cup B$ ) will contain all the values that are present in either A or B.

$$\therefore A \cup B = \{2, 3, 4, 5, 6, 7, 8\}$$

the intersection of A and B ( $A \cap B$ ) will contain all the values present in both A and B

$$\therefore A \cap B = \{2\}$$

the set difference of A and B ( $A - B$ ) will contain the values that are present in A but not in B

$$A - B = \{4, 6, 8\}$$

Q Total students in class = 50

$$\begin{array}{l} E + H = 10 \\ E = 32 \end{array} \quad \left| \begin{array}{l} \text{need to find only Hindi} \\ A \cup B = A + B - A \cap B \\ 50 = A + 32 - 10 \end{array} \right.$$

$$50 - 22 = A = 28$$

$28 \Rightarrow$  total students enrolled in Hindi

$$28 - 10 = 18 \quad (\because H = (E + H))$$

Q  $A \cap (B \cup A) = ?$

$$(A \cap B) \cup (A \cap A) = (A \cap B) \cup A$$

$$= (A \cup A) \cap (B \cup A)$$

$$= A \cap (B \cup A)$$

$$A = \{1, 2, 3, 6\}, \quad B = \{4, 5, 6\}$$

$$B \cup A = \{1, 2, 3, 4, 5, 6\}$$

$$\Rightarrow A \cap (B \cup A) = A$$

## Types of Joins

In SQL joins are used to combine rows from two or more tables based on a related column b/w them. Joins allow us to query data across multiple tables in a relational database by creating

relationships b/w tables.

## Types of Joins in SQL

### ① inner join:

- ↳ Returns records that have matching values on both tables
- ↳ used when we only need rows where there is a match in both tables.

eg Select \* from table1 inner join table2 on  
table1.id = table2.id;

### ② Left Join or (left\_outer join)

- ↳ Returns all the records from the left table and the matched records from the right table, if no match is found nulls are returned for columns from the right table.

↳ used when we want all data from primary (left) table even if there are no matches found, Nulls are returned for cols from right table (if no match)

eg Select \* from table1 left join table2 on  
table1.id = table2.id;

### ③ Right Join:

- ↳ Returns all the records from the right table and matched records from left table

eg Select \* from table1 right join table2 on  
table1.id = table2.id;

#### 4. full join

↳ Returns all the records when there is a match in either left or right table. If there is no match nulls are returned for non matching rows from either table.

↳ Useful for combining data when you need to see all possible records regardless of matches.

eg Select \* from table1 full join table2 on table1.id = table2.id;

#### 5. Cross join :

↳ Produces the Cartesian product of two tables, meaning it joins every row of the first table with every row of the second table.

→ Useful when you need every combination of rows b/w two tables

eg Select \* from table1 cross join table2;

#### 6) Self join :

↳ A join where a table is joined with itself.

↳ Useful when comparing rows within same tables.

eg Select a.column\_name, b.column\_name  
from table1 a, table1 b where  
a.id = b.Parent\_id;

## Multijoin

A multijoin refers to a SQL operation where more than two tables are joined together within a single query. This type of join is useful when we need to combine data from multi-related tables.

e.g scenario: Imagine a database for an e-commerce platform with the following tables

- Customers → contains customers details
- Orders → Contains order info
- Products → Contains Product details
- OrderDetails → contains orders line items linking orders and products

e.g Select

customers. CustomerName,  
Orders. OrderDate,  
Products. ProductName,  
OrderDetails. Quantity

From

customers

Join

orders on customers.cust\_id = orders.cust\_id

Join

OrderDetails on Orders. order\_id = OrderDetails.order\_id

Join

Products ON OrderDetails.prod\_id = Products.prod\_id

## ON vs USING

The diff b/w ON and Using is SQL pertains to how conditions are specified. In Join clauses, both are used to define how tables are matched when performing joins, but there are subtle distinctions b/w them.

ON → used to specify join conditions b/w cols of the tables being joined.

It's more flexible than Using because it allows conditions that are not simple col comparisons, you can use any boolean expression.

Syntax: Select \* from tableA Join tableB ON  
tableA.col1 = TableB.col2;

USING : It is used when both tables have a col with the same name and you want to join on that column.

Syntax: Select \* from TableA Join TableB  
Using (column-name);

\* Views : Views in SQL are virtual tables created by a query which allows you to save complex queries for easier access and use later.

A view doesn't store data or logic or retrieves data from the underlying tables every time it is queried.

### Creating View with Joins

e.g Employees table

emp_id	name	dept_id
1	John	101
2	Jane	102
3	Alice	101

Department Table

dept_id	department_name
101	Sales
102	Marketing
103	IT

# Now we want a view that shows each employee's name and department name

> Create view EmployeeDepartmentView AS  
Select emp.emp\_id, emp.name, dept.dept\_name  
from employees emp innerjoin department dept  
on emp.dept\_id = dept.dept\_id;

# To see the view

Select \* from EmployeeDepartmentView;

# To Drop a view => Drop view EmployeeDepartmentView;

# SET Operations with SQL

Set operations in SQL allow you to combine the results of two or more select queries.

## ① UNION

- ↳ Combines the result of two or more select queries and returns only distinct values (remove duplicates)

Rules:

- the number and order of columns must be the same in all queries
- the datatype of columns must be compatible

e.g.

Table A

ID	NAME
1	Aliu
2	Bob

Table B

ID	NAME
2	Bob
3	Charlie

Select ID, Name, from Table A

UNION

Select ID, Name from Table B;

ID	NAME
1	Aliu
2	Bob
3	Charlie

## ② Union All

- ↳ combines the result of two or more select queries and return all rows, including duplicates

Select id, name from table A

Union all

Select id, name from table B

ID	NAME
1	Aliu
2	Bob
2	Bob
3	Charlie

# Joins Adv eq (Advanced Joins)

✳ Natural Join: Joins tables based on columns with the same name and datatype; it automatically matches columns without specifying condition.

e.g. Select columns from table 1 Natural Join table 2;

✳ Using Clause: Specifies the column to join on when columns have the same name in both tables.

e.g. Select cols from table1 Inner join table2 Using (col\_name);

✳ Join with subqueries: Allows joining with the result of a subquery.

e.g. Select main-table.columns, subquery-table.columns from main-table

Join (Select columns from table where condition) subquery-table on main-table.column = subquery-table.column;

## Multiple Joins : Joining more than 2 tables,

e.g.

Select columns from table1

Inner join table2 on table1.col-name = table2.col-name

Inner join table3 on table2.col-name = table3.col-name

## OFFSET In SQL

Offset In SQL Is used to skip a specific no. of rows in the resultset before starting to return rows , It's often used with limit or fetch to paginate results.

Syntax → Select columns from table

order by column OFFSET num-of-rows-to-skip Rows;

e.g. Select employee-id, employee-name from employees order by employee-name

offset 10 Rows;

↳ Query will skip first 10 rows and return the rows left after the skipped ones.

e.g. with limits:

Select \* from employees order by empname offset 2 limit 1;

↳ 3rd record

## Rank functions - I

When we used the `Group By` clause It reduced the no. of rows ; It also leads to the loss of individual properties of various rows.

To address the two points i.e reduction of rows and loss of individual properties of rows, SQL provides a special clause known as "`OVER`" clause

which displays group characteristics in the form of a new column.

Ex

for eg, If we want to find the top 3 records based on somebody's salary

↳ So we sort it on salary in Descending Order then apply `Limit 3`.

↳ Problems with this approach:

- It is not a standard way of doing
- It is very rigid i.e. we can only do top3, top5 kind of a query.

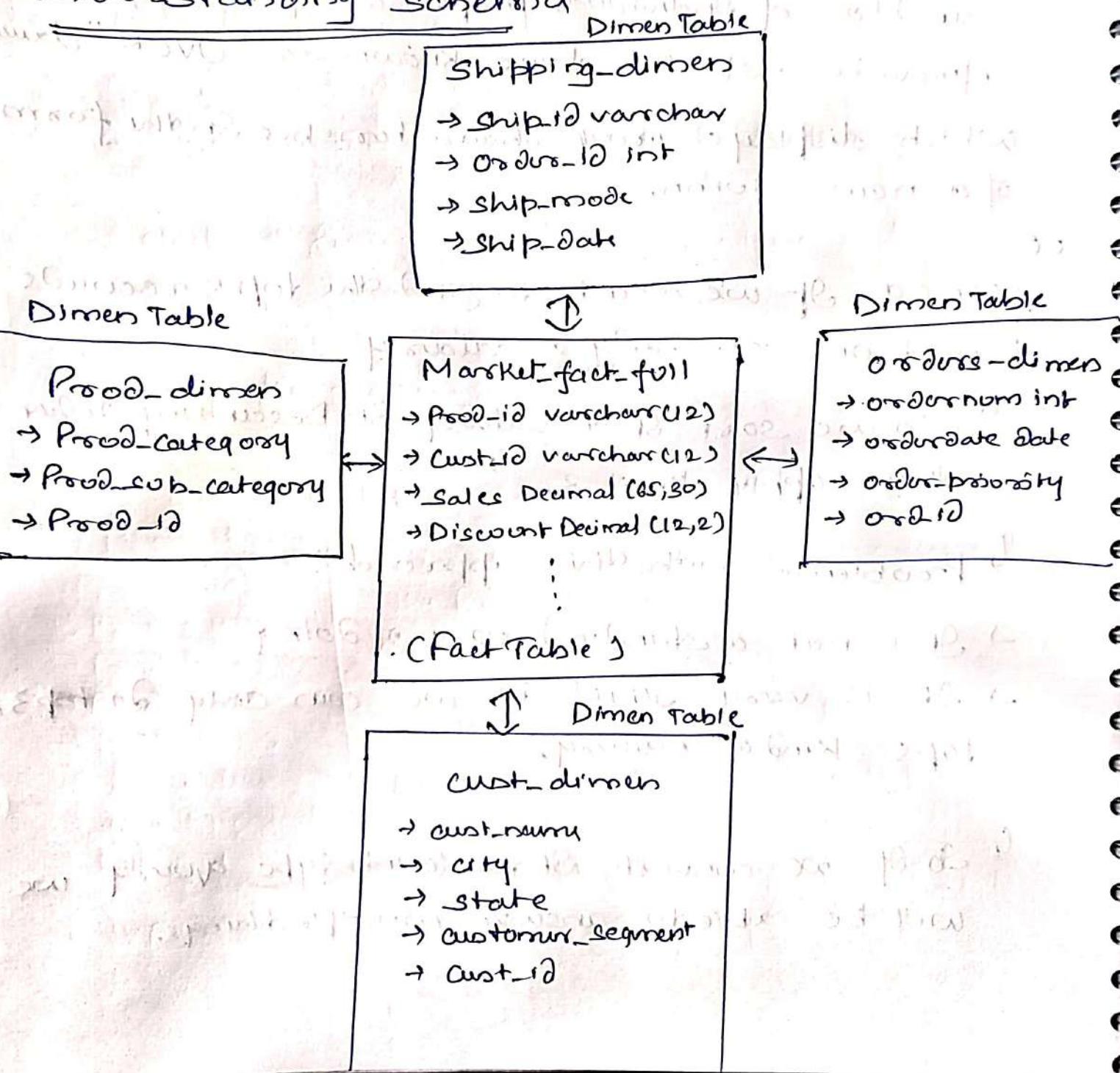
↳ So if we convert it to rank type query we will be able to answer multiple things.

# RANK Syntax

RANK () OVER (

PARTITION BY <expression> [ {, <expression>} ]  
ORDER BY <expression> [ ASC | DESC ], [ {, <expression>} ... ]

## Understanding Schema



> USE market\_star\_schema;  
 # we will try to rank all the sales based on  
 order id in descending order  
 SELECT customer\_name, order\_id,  
 ROUND(sales) AS rounded\_sales,  
 RANK() OVER (ORDER BY SALES DESC)  
 AS sales\_rank  
 FROM market\_fact\_full AS m  
 INNER JOIN  
 cust\_dimen AS c  
 ON m.cust\_id = c.cust\_id  
 WHERE customer\_name = 'RICK WILSON';

	customer_name	order_id	rounded_sales	sales_rank
1	Rick Wilson	5186	26095	1
2	Rick Wilson	3841	8496	2
3	Rick Wilson	1207	8407	3

# Now if we want to filter top 10 records  
 o for this we will use common table  
 expression o for this we use "With"  
 keyword.

With rank\_info as

(  
Select customer-name,  
order\_id, Round(sales) As rounded\_sales,  
RANK() OVER (ORDER BY sales DESC) AS  
sales\_rank

From market\_fact\_full as m

Inner join

Cust-dimensions as c

On m.cust\_id = c.cust\_id

Where customer-name = 'Rick Wilson'

)  
Select \* from rank\_info where sales\_rank  
 $\leq 10$

6 top 10 records fetched.

Eq 2 Given a table name Products, with following cols

Products

- Prod\_code
- Prod\_name
- Prod\_line
- Prod\_scale
- Prod\_vendor
- Prod\_desc
- avQuantityInStock
- buyPrice

→ Write a query to retrieve the rank of the products in decreasing order of avQuantityInStock

→ Select avQuantityInStock,  
RANK() OVER (ORDER BY  
avQuantityInStock DESC) AS  
avQuantityRank FROM products;

avQuantityInStock	avQuantityRank
9995	1
6534	2
5839	3

# We have dense-rank and percent-ranks functions.

Let's say we have 5 students in a class

↳ 3 students scored 98/100

4<sup>th</sup>, 5<sup>th</sup> student scored 95, 90

Now if we apply rank thru 1st 3 students in the rank func will get 1, 1, 1 then 2 and 3 is skipped because we already covered 3 students 4<sup>th</sup> and 5<sup>th</sup> student will get 4, 5 as ranks.

# Instead of this approach, if we want to give the rank to 4<sup>th</sup> and 5<sup>th</sup> student as 2<sup>nd</sup>, 3<sup>rd</sup> given that they got 2<sup>nd</sup> highest, 3<sup>rd</sup> highest marks → for this we use dense-rank

What does percent rank do?

↳ Converts the entire ranking into % age, whole class will get nos from 0 to 1, similar to percentile

Types of Rank:

→ RANK(): Rank of current row within its partition, with gaps.

→ Dense-rank(): Rank of current row within its partition without gaps.

→ Percent\_rank(): Percentage Rank Value

## DENSE\_RANK and PERCENT\_RANK syntax

Dense\_rank() over (

Partition by <expression> [ {, <expression>}... ]  
Order by <expression> [ASC | DESC],  
[ {, <expression>}... ]

Percent\_rank() over (

Partition by <expression> [ {, <expression>}... ]  
Order by <expression> [ASC | DESC],  
[ {, <expression>}... ]

Now let's understand how to use different rank functions in queries to solve various problems.

eg for dense\_rank

> Select ord\_id,  
discount,  
customer\_name,  
RANK() over (order by discount Desc) as  
disc\_rank,

DENSE\_RANK() over (order by discount Desc)  
As disc\_dense\_rank

From market\_fact\_full as m Innerjoin customers as  
c on m.cust\_id = c.cust\_id where customer\_name = "Rick Wilson"

ord_id	discount	cust_name	disc_rank	disc_dens_rnk
Ord-3855	0.10	Rick Wilson	1	1
Ord-5186	0.09	Rick Wilson	2	2
Ord-1209	0.09	Rick Wilson	2	2
Ord-3841	0.08	Rick Wilson	4	3

RANK() DENSE\_RANK()

We learnt about the 'RANK', 'DENSE\_RANK', 'PERCENT\_RANK' functions, In this segment we will learn about fourth type

→ Using RANK() and Dense\_Rank we can get the rows having exact same rank, So to avoid this we have a concept of Row Number, It gives a unique value for each and every row.

## Row NUMBER Syntax

Row\_NUMBER() OVER (

Partition By <expression> [ , <expression>... ]

Order By <expression> [ASC | DESC], [ , <expression>... ]

)

We can use "Row\_NUMBER()" for the following type use cases.

→ To determine the top 10 selling products out of a very large variety of products

→ To determine the top 3 winners in a car race.

Q9 We want to rank various customers based on the nos of orders placed.

> Select customer\_name,  
Count(Distinct order\_id),  
RANK() over (Order by count(Distinct order\_id)  
DESC),  
DENSE\_RANK() over (Order by count(Distinct order\_id) DESC),  
Row\_Number() over (Order by count(Distinct order\_id) DESC)  
from market\_fact\_full as m  
Inner Join  
cust\_dimen As c  
ON m.cust\_id = c.cust\_id  
Group By customer\_name;

### \* Partitioning

We know that 'rank()' and 'dense()' func are used with Over Clause, However this may not be enough if we want to rank groups of rows based on certain criteria,

for eg we can rank the top 10 Batsmen in the world using rank() function, what if we want to find out top 10 batsmen from each team

↳ for this use case we use "partition" and "over" clauses together.

-- Partitioning example --

first creating data

Select ship-mode, Month(ship-date) AS shipping-month, count(\*) as shipments

From

Shipping-dimen

Group By ship-mode,  
Month(ship-date);

ship-mode	shipping-month	shipments
Regular Air	10	467
Regular Air	11	467
Delivery truck	6	91
Regular Air	5	564
Delivery truck	9	86
		488

On top of this we have to rank

with shipping-summary AS

( Select ship-mode, MONTH(ship-date) AS  
shipping-month,

COUNT(\*) AS Shipments

From

Shipping-dimen

Group by ship-mode,

Month(ship-date)

) Select \*, RANK() OVER (Partition BY  
Ship-mode ORDER BY Shipments DESC) AS  
Shipping-Rank

From Shipping-summary;

Ship-mode	Shipping-month	Shipments	Shipping-Rank
Delivery Truck	12	103	1
Delivery Truck	9	105	2
Delivery Truck	3	105	2
Express Air	12	96	1
Express Air	7	91	2
Express Air	4	86	3

Rank got reset here

## Over function

The over func in SQL is used to perform with the window func to perform calculations across a set of rows related to the current row, without requiring grouping into a single result row.

This is particularly useful in analytical queries where you want to get results on a row by row basis with some aggregation or calculation applied.

### Common Window function with Over

- ① Row\_Number : Assigns a unique sequential integer to rows within a partition.

```
eg Select employee_id, dept_id, salary,  
Row_Number() over (Partition By  
dept_id ORDER BY salary DESC)  
AS rownum FROM employees;
```

# Here, Row\_Number() will restart the numbering for each department, ordered by salary in descending order

2. RANK : Similar to Row-NUMBER, but if two rows have the same value, they get the same rank, and the next rank is skipped.

> Select employee\_id, dept\_id, salary,  
RANK () over (Partition By dept\_id Order by  
salary DESC) As rank from employees;

3. DENSE\_RANK : Similar to RANK, but without skipping ranks if there are ties.

> Select employee\_id, dept\_id, salary,  
DENSE\_RANK () OVER (PARTITION BY dept\_id  
ORDER BY salary DESC) AS dense-rank  
from employees;

4. SUM, AVG, MIN, MAX : Aggregate func that provide a running total, average, minimum, or maximum for a specific window

> Select employee\_id, salary,  
SUM (salary) OVER (Order by employee\_id)  
AS running\_total from employees;

## Components of the OVER Clause

- Partition By : Divides the result set into partitions, similar to groups in aggregate functions.
- ORDER BY : Defines the order of rows within each partition.
- ROWS : Specifies the range of rows for each calculation.

Eg Table named → ordersdetails with following cols

Ordersdetails

- OrderNumber
- ProductCode
- quantityOrdered
- PriceEach
- OrderLine Number

expected o/p

OrderNumber	IndividualOrderAmount	TotalOrderAmount
10154	2332.13	4465.85
10154	2133.72	4465.85

Q-there are some orders that have multiple order amounts corresponding to same order , write a query to retrieve the individual and total order amounts for each order along with order numbers

> Select

OrderNumber,  
(quantityOrdered \* priceEach) As individualOrderAmount,  
sum(quantityOrdered \* priceEach) Over (Partition By  
OrderNumber) As totalOrderAmount

From OrdersDetails

Order By OrderNumber ASC, individualOrderAmount DESC

# Windows In SQL

In SQL, a window function is a feature that allows you to perform calculations across a set of table rows related to the current row, without needing to use GROUP BY.

The concept is called window because the calculations are done over a subset or "window" of rows within a resultset.

## Key Concepts:

1. Window func: functions that performs calculation across the rows in a specified window or Partition.

Common window func are:

- ROW\_NUMBER(): Assigns a unique rownumber within the partition.
- RANK(), DENSE\_RANK(): Assigns a rank to each row within a partition
- SUM(), AVG(), MIN(), MAX(): Performs aggregate calculations within a ~~weak~~ partition

2. OVER(): Specifies the window (or range of rows) to use with the window function

→ Partitioning (PARTITION BY): Divides the data into partitions or subsets for each calculation e.g. Partition by department, country etc.

→ Ordering (Order By): Determines the order in which rows within partition are arranged.

e.g. Let's say we have "Sales" table with cols sales-person-id, sales-amount, sales-date, we want to calculate running total of sales-amount for each salesperson

> Select

sales-person-id,

sales-date,

sum(sales-amount) over(partition by

sales-person-id order by sales-date desc)

As running-total

From "Sales" table must get sales' total amount for each salesperson

## Named Window

A Named window function in SQL is a feature that allows you to define a reusable window specification with a name, which can then be referenced in multiple windows functions within a query, particularly useful for simplifying complex queries.

### Syntax:

We use "window" clause followed by window-name and the window-specification, you then refer to the named window in the over clause

ee  
WINDOW window-name AS (window-spec)

Here window-spec lets us define details like Partition By, Order By etc

Ex  
We have a table "Sales" with columns Salesperson-id, sale-amount, sale-date

We want to calculate running total and moving avg for each salesperson

eq

```
    Select  
        Salesperson-ID,  
        Sale-Date,  
        Sale-amount,  
        Sum(Sale-amount) over w as total  
        Avg(Sale-amount) over w as average  
    From  
        Sales  
    Window w As ( Partition By Salesperson  
                    Order By Sale-amount )  
    Return Sale-amount, total, average;
```

## FRAMES

The frame clause is used with window function in frames in SQL. Whenever we use a window function in SQL, it creates a ~~partition~~ "window" or "partition" depending upon the column mentioned after the "partition by" clause in the 'over' clause and thus it applies the window func to each of those partitions. Inside these partitions, we can create a subset of records using the frame clause.

∴ A frame clause specifies a subset

Frames are useful when you want to calculate values that depend on the nearby rows.

## Using frames with window functions

frames work with window func in SQL

A window func like SUM() or AVG()

lets you perform calculation over a window of rows around each row. When we define a frame we tell SQL which rows around the current row should be included in the calculation.

e.g Select employee\_id,

salary,

SUM(salary) OVER (ORDER BY employee\_id)

AS running-total

from employees;

# Here sum() calculates the running total of salaries and the order by employee\_id sorts the rows by employee\_id.

## What Does a frame DO?

A frame lets you control the specific range of rows in each window

- which rows to start and end with (e.g. Start from beginning look at only the current row or look at the last two rows).
- How far forward or backward to look (e.g. include two rows before and one row after the current row).

## Defining frames with Rows and RANGE

- **Rows** : Looks at the fixed number of rows relative to the current row.
- **Range** : Looks at rows with specific value in the ordered column, allowing you to include row with the same value.

### frame Options

e.g.: Running total using a frame

to get a running total of salaries for each employee, you can include all rows from start up to current row.

Select employee\_id, salary,

\$ 0.00

SUM(salary) OVER (ORDER BY employee\_id)

100.00

ROWS BETWEEN UNBOUNDED

PRECEDING AND CURRENT ROW) AS  
running\_total FROM employees;

- UNBOUNDED PRECEDING :! Include all the rows preceding from start
- Current row → Includes the current row.  
So for each employee, this frame adds up the salaries of all employees up to that point.

Q2: Moving Avg of Salaries Over a few rows  
to calculate a moving avg of salaries for each employee based on two rows before and the current row.

? Select employee\_id, salary, Avr(salary) OVER (ORDER BY employee\_id ROWS BETWEEN 2 PRECEDING AND CURRENT ROW) AS Moving\_avg  
from employees

Q9

Sales_id	employee_id	sale_amount	sale_date
1	101	100	2024-01-01
2	102	150	2024-01-02
3	101	200	2024-01-03
4	103	300	2024-01-04
5	101	250	2024-01-05
6	102	400	2024-01-06

# to calculate the running total of sales upto each sale date

Select sales-id, sale-date, sale-amount,  
SUM(sale-amount) OVER ( ORDER BY sale-date,

ROWS BETWEEN UNBOUNDED  
PRECEDING AND CURRENT ROW) AS

Running-total  
from Sales ;

G O/P

Sales-id	sale-date	sale-amount	running-total
1	2024-01-01	100	100
2	2024-01-02	150	250
3	2024-01-03	200	450
4	2024-01-04	300	750
5	2024-01-05	250	1000
6	2024-01-06	400	1400

Q2 # calculating the 3 day moving avg of sales

Select sales-id, sale-date, sale-amount, AVG(sale-amount)  
OVER ( ORDER BY sale-date ROWS BETWEEN 2 PRECEDING  
AND current-row) AS moving-avg from sales;

C

Sales-id	sale-date	sale-amount	moving-avg
1	2024-01-01	100	100
2	2024-01-02	150	125
3	2024-01-03	200	150
4	2024-01-04	300	216.67
5	2024-01-05	250	250
6	2024-01-06	400	316

## SQL eq

Select ship-date , SUM( shipping-cost ) AS daily-total  
From market-fact-full as m Inner join shipping-dimen AS s on m.ship-id = s.ship-id  
Group By ship-date;

C

Ship-Date	Daily-total
2010-10-20	93.92
2010-11-24	40.57
-	-
0.21	0.21

-- Frames examples --

With daily\_shipping\_summary as

C Select ship-date ,  
SUM(shipping-cost) AS daily-total

From

market-fact-full as m

Inner Join

shipping-dimen as s

on m.ship-id = s.ship-id

Group By ship-date

) Select \*,

SUM(daily-total) OVER w1 as running-total,

```
AVG(daily-total) OVER w2 AS moving-avg  
FROM daily-shipping-summary  
WINDOW w1 AS (ROWS BY daily-total ROWS  
UNBOUNDED PRECEDING),  
w2 AS ORDER BY daily-total ROWS 6 PRECEDING  
;
```

## RANGE Clause

Range clause within a window function defines a specific range of values relative to the current row for calculations like averages, sums, and other aggregations. Unlike the rows clause which counts a fixed number of rows, the range clause looks at values in a specified column.

⇒ the RANGE clause is used in the OVER() part of window func

### Syntax

```
<window-func> OVER(  
    ORDER BY <column>  
    Range between <n> PRECEDING AND <m>  
        FOLLOWING )
```

Q9 Kohli In ODIs  
Table given below contains the number of runs scored by Virat Kohli over a time period

Kohli-Batting      to calculate  
                        4 Moving Avg.

Year	Runs
2008	159
2009	328
2010	995
2011	1382
2012	1028

Years	Runs	Moving Avg
2015	623	1145.4

Q Find back and forth for moving average

Select

Year, Runs,

Avg(Runs) OVER ( ORDER BY Year )

RANGE BETWEEN 4 Preceding AND

Current Row ) AS Moving-Average

From

Kohli-Batting

Order By

Year;

So, some necessary keyword in frames are  
Unbound preceding, following, between etc

(Eq)

Row1

Row2

Row3 ← current cursor

Row4

Row5

Row6

Row7

Order by x Rows 2 following

↳ Row4 and Row5

↳ current Row of Neeche of 2 rows

(Eq)

Row1

Row2

→ current cursor → Row6

Row9

Row10

Order by x Rows Between 2 Preceding  
and 3 following

↳ Row4 to Row9

(Eq)

Row1

Row2

Row3

Row4

Row5

Row6

Row7

Row8

Row9

Row10

→ current cursor → Row6

↳ Order by x Desc Rows

Between unbounded Preceding  
And 2 following

↳ Row4 to Row11

∴ the Desc keyword has been used  
the order of the rows on which the  
frame will be created will be reversed.

## Lead and Lag functions

The Lead and Lag functions are useful in SQL and data analysis for accessing data from a row before or after the current one without needing a self join or complex subqueries. They are often used in window functions, where data is organized based on specific partitions or orderings, making it easier to analyze sequential data patterns.

### LEAD Function

The LEAD function allows you to access data from a subsequent row in your dataset, based on the specific offset. For e.g. you might want to see the next row's value in each current row, useful for comparisons with future datapoints.

#### Syntax

```
LEAD (col-name, offset, default-value) OVER  
( Partition by Partition-col ORDER BY  
order-column)
```

- col-name : the column to retrieve the data from
- offset : How many rows ahead to look  
(Default is 1)
- default-value : Value to return if the requested row doesn't exist (Default is NULL)

Eg . We have table Sales with columns "day"  
& "revenue"

day	revenue
Mon	100
Tue	150
Wed	200
Thurs	250
Fri	300

# we want to see each day revenue along with revenue of the next day

? Select day, revenue, LEAD (revenue, 1, 0)  
over (order by day) as next-day-revenue from Sales;

O/P

day	revenue	next-day-revenue
Mon	100	150
Tue	150	200
Wed	200	250
Thurs	250	300
Fri	300	0

Note: For Friday, next-day-revenue is 0 because there is no following day's revenue.

## LAG func

The LAG func works similarly but access data from a preceding row, this is useful when we want to compare each row with previous value

Syntax:

LAG ( col-name, offset, default-value ) OVER  
( Partition By partition-column Order By  
order-column )

Eg # If we have same sales table, we want to compare each day's revenue with prev day's revenue

Select day, revenue, LAG(revenue, 1, 0)  
OVER (Order by day) as prev-day-revenue

	day	revenue	prev-day-revenue
	Mon	100	0
	Tue	150	100
	Wed	200	150
	Thur	250	200
	Fri	300	250

## Applications

→ LEAD and LAG are excellent for tracking trends over time, such as changes in sales, stock prices etc.



## CASE Statements

We encountered if else and case statements in python, java, similarly we have it in SQL too.

### CASE Statement Syntax

CASE

when condition1 then result1

when condition2 then result2

when conditionN then result N

ELSE result

END as col-name

Eg Demo # In the market\_fact full table we will be fetching market\_fact\_id, profit columns

Select market\_fact\_id, profit from  
market\_fact\_full;

Market_fact_id	Profit
1	-3051
2	4.56
3	1148.56
4	729.0

Now we want to categorise the profit into 4 buckets (Huge loss, bearable loss, decent profit, great profit)

Case when example --

/\* Profit < -500 → Huge loss

Profit -500 to 0 → Bearable loss

Profit 0 to 500 → decent profit

Profit > 500 → great profit \*/

> Select market\_fact\_id,  
Profit,

CASE

when profit < -500 Then 'HugeLoss'

when profit Between -500 AND 0  
Then 'Bearable Loss'

when profit Between 0 and 500  
Then 'Decent Profit'

ELSE 'Great Profit'

END AS Profit-type

From

market\_fact\_full;

eg Below is the table for Income tax slabs

### Tax Slabs

Taxable Income Slabs	Tax Slabs
Upto Rs 2.5 lakh.	A
Rs 250001 to Rs 500000	B
Rs 500001 to Rs 1000000	C
Rs 1000000 and above	D

# Add another col to output the tax slab along with given data.

Select Name, Salary,

CASE

When Salary <= 2.5 Then "A",

When Salary Between 2.5 and 5

Then 'B'

When Salary Between 5.1 And 10 Then 'C'

End AS Tax-slab

from salaries;

Name	Salary (in lpa)
Sundar Pichai	6.8
Jeff Bezos	8.7
Bill Gates	9
Anil Ambani	14

### Common table expression CTE's

CTE is a temporary named result set in SQL  
often used to simplify complex queries and  
improve readability

→ CTE's are created using "with" clause and  
can be referenced within the main query or  
other CTE's.

CTEs are useful when you need to use the same set of data multiple times within a query.

### Syntax

With cte-name as (

```
    select col1, col2  
    From table-name  
    where condition
```

)

```
Select * from cte-name)
```

} # Defining the array for CTE.

Ex

with ProductRevenue AS (

```
    Select product_id, sum(quantity * price) As  
        total_revenue from sales  
    group by product_id
```

)  
Select product\_id, total\_revenue from  
ProductRevenue  
where total\_revenue > 5000;

Ex Tax Slabs 2

### Tax Rates

Taxable Income Slab

upto rs 2.5 lakh

rs 250001 to 500000

rs 500001 to 1000000

rs 100001 and above

Income Tax rates

No.1

5% of (Total Income - 250000)

12500 + 20% (Total Inc - 500000)

112500 + 30% (Total Inc -

1000000)

Select name, salary,

CASE

when salary <= 2.5 then 0

when salary between 2.5 and 5 then

ROUND(0.05 \* (salary - 2.5) \* 100000)

when salary between 5 and 10 then

ROUND(12500 + 0.2 \* (salary - 5) \* 100000)

else ROUND(112500 + 0.3 \* (salary - 10) \* 1000000)

End as Tax-amount

From salaries;

## UDFs (User Defined functions )

You already know how to deploy various inbuilt MySQL functions, such as sum(), avg() and concat().

Now it's possible to find the sum of two numbers in MySQL without using sum() func, Now we can use arithmetic + operator.

However there are operations that you may want to repeat multiple times in a piece of code because they do not have Inbuilt func → for this we use user-defined functions (UDFs) or stored func

## Syntax

\* Create func is a DDL statement

Syntax :

```
Create Function function-name (func-parameter1,  
                                func-parameter2,...)  
    RETURN datatype [characteristics]  
    /* func-body */  
    Begin  
        <MySQL statements>  
        Return expression;  
    End
```

Note: the func body must always contain one return statement.

-- Example --

```
DELIMITER $$  
CREATE FUNCTION profitType (Profit int)  
RETURNS VARCHAR(30) DETERMINISTIC  
BEGIN  
    DECLARE message VARCHAR(30);  
    IF Profit < -500 THEN  
        SET message = 'Huge Loss'  
    ELSE IF Profit BETWEEN -500 and 0, THEN  
        SET message = 'Bearable Loss'  
    ELSE If Profit BETWEEN 0 and 500 then  
        Set message = 'Decent Profit'
```

```

ELSE
    SET message = 'Great Profit'
END IF;
RETURN message;
END;
$$
DELIMITER ;

```

Note: A stored func returns one and only one value.

Deterministic keyword is used to ensure that the output is the same for the same input values

calling func

```

Select ProfitType(10);
Select ProfitType(-20);

```

## Stored Procedures

We want to write certain snippets which keeps on getting used again and again.

### Syntax

DELIMITER \$\$

CREATE PROCEDURE PROCEDURE-name [<Parameter list>]

BEGIN

SQL statement of stored Procedure { Can be  
N nos  
of SQL  
Statement }

END \$\$

DELIMITER ;

# Call Statement

CALL Procedure-name;

CQ

DELIMITER \$\$

Create procedure get\_sales\_customers (sales\_input  
int)

BEGIN

Select DISTINCT cust\_id, ROUND(sales) AS  
sales\_amount

From

market\_fact\_full

where ROUND(sales) > sales\_input

order by sales;

END \$\$

DELIMITER ;

CALL get\_sales\_customers (300)

Note : ## For Dropping the Procedure

DROP PROCEDURE <Procedure\_name>

# Difference b/w UDF and Stored Procedure

UDF

Stored Procedures

- It supports only one input
- It can be called using select statement
- It must return a value
- Only select operation is allowed
- It supports input/output and input-output Params
- It can be called using CALL func
- It need not return a value
- All database operations allowed

## Best Practices

① Write appropriate comments

In MySQL for single line comments

↳ -- Single line comment

for Multiline Comments

/\* comment section

with multiple lines

\*/

② Always use table aliases when SQL statement involves more than one source

↳ In Joins we use AS

③ Assign simple descriptive names for columns

④ Write all keywords in CAPS and variable names in small case

⑤ To use proper Indentation

⑥ use new line for different sections of the query

⑦ Use SQL formatter (Many available online) for beautifying the query

## Indexing

One of the ways in which you can greatly reduce the runtime of your queries is by using indexing.

Indexing is the process of referring to only the required values directly, instead of going through entire table. This prevents the query engine from looking up the values one by one, instead it returns the exact value that you want right away.

Indexing is necessary for querying extremely large datasets, a primary key is an index because it helps in identifying each record in a table uniquely.

Note: Although you can not view indexes whenever you execute queries, they occupy some storage space, indexes are internal to database engine.

When we are writing our select query

Select <columns>

from \_\_\_\_\_

where Salary < 10000

or DeptId = 5;

the most inefficient way to fetch the record line by line and making decision on a row by row basis  
But notice the where clause here

```
Select <columns>  
from employee  
where salary < 10000  
or DEPT_ID = 5;
```

gt Attributes where clause of use here  
Ex 3rd ex indexing over E |

MySQL has a specific DDL command called as "CREATE INDEX" → used for specifying the attribute on which you want to create an index.

## INDEX SYNTAX

→ Index is a DDL statement

Creating an Index

```
CREATE INDEX index-name ON table-name  
    (column1, column2, ...);
```

# we can also use Alter command

```
Alter table table-name ADD INDEX index-name  
    (column1, column2, ...);
```

# Dropping an Index

ALTER TABLE table-name DROP INDEX index-name;

-- Index Demo --

# first creating a temporary table from an existing table.

> CREATE TABLE market-fact-temp AS  
SELECT \* FROM market-fact-full;

# Create Index (using combination of cust-id, ship-id, prod-id)

> CREATE INDEX filter-index ON market-fact-temp  
(cust-id, ship-id, prod-id);

# for dropping the index

> ALTER TABLE market-fact-temp DROP INDEX  
filter-index;

# the "Where" command in SQL works efficiently when we are using indices because the index will enable us to get the exact rows that you need without processing entire table row by row.

# CLUSTERED VS NON CLUSTERED

## INDEXATION

### Clustered Index

- this is mostly the primary key of the table
- It is present within the table
- there is no requirement for a separate mapping as this index is part of the table.
- relatively faster

### Non-clustered Index

- It is a combination of one or more columns of the table.
- the unique list of keys will be outside the table
- the external table will point to different sections of the main table
- relatively slower

## Orders of Execution of the Query

the order in which the various SQL statements appear in the query is as follows

- 1) Select
- 2) From
- 3) [Join]
- 4) Where
- 5) Group By
- 6) HAVING
- 7) Window
- 8) Order By

} the order in which various statements are executed by the database engine is not same.

However the order in which the various statements are executed by the database engine is not the same

From (including joins)

↳ Where

↳ Group By

↳ Having

↳ Window Func

↳ select

↳ Distinct

↳ Order By

↳ Limit

↳ Offset

Some of the tips points

that you should keep in mind  
while writing a query

→ Use inner joins wherever possible  
to avoid having unnecessary rows  
in resultant table.

→ Avoid using Distinct while using group by clause as it slows down query processing

### Joins vs Nested Queries

EMP

Eid	Name	DID

DEPT

DID	DNAME
1	HR
2	IT

Now we want to get the name of the employee whose department name is HR

one way using subquery

6 Select name from Emp  
where DID in (

Select DID  
from Dept  
where DName = 'HR');

Using Joins

Select Name from  
Emp innerjoin  
Dept on Emp.DID  
= Dept.DID  
where DName =  
'HR');



## Profitability Analysis

Problem Statement: Growth team wants to understand sustainable (profitable) product categories.

- Sustainability can be achieved when we make better profits or at least positive profits
- We can look at the profits per product category.
- We can look at profits per product subcategory.
- We can check Avg profit per order
- Also consider Avg profit % per order.

Tables that we will use:

- ⇒ market\_fact\_full
- ⇒ Prod\_dimen
- ⇒ Orders\_dimen

# Schema that we will use :

① shipping-dimen

- ship-id VARCHAR(12)
- Order-id int(11)
- Ship-mode varchar(25)
- Ship-date DATE

② market-fact-full

- Prod-id VARCHAR(12)
- cust-id VARCHAR(12)
- sales Decimal
- Discount Decimal
- OrderQuantity
- Profit
- Shipping-cost
- Product\_Base\_Margin
- Market-fact-id
- ord-id
- ship-id

③ Prod-dimen

- Product-category
- Product-subcategory
- Prod-id

④ Order-dimen

- Ordernumber
- Orderdate
- Order-priority
- Ord-id

⑤ cust-dimen

- customer-name
- city
- state
- customer-segment
- cust-id

Problem statement: (growth team) sustainable product categories  
→ Sustainability can be achieved when we make better profits or at least positive profits.

- ↳ we can look at the profits per category
- ↳ we can look at the profits per product subcategory
- ↳ we can check avg profit per order
- ↳ we can check Profit% per order

# first get the data first

```
Select P. Product_category,  
      P. Product_SubCategory,  
      SUM(m.Profit) AS Profits
```

FROM

market\_fact\_full AS m

INNER JOIN

Prod\_dimen AS P

ON m.prod\_id = P.prod\_id

GROUP BY

P. Product\_category

P. Product\_SubCategory,

Order By SUM(m.Profit);

→ O/P

Prod_cat	Subcateq.	Prod
—	—	—
—	—	—
—	—	—
—	—	—

We are getting profits made at a Product\_category and Product\_subcategory level but we aren't able to get any insights out of this.

↳ +ve Profits are at top.

# for Profits Per category

```
Select P. Product_category, SUM(m.Profit) AS Profits  
From market_fact_full AS m Inner Join Prod_dimen  
AS P On m.prod_id = P.prod_id Group By P. Product_  
category Order By SUM(m.Profit);
```

O/P	Product category	Profits
	furniture	117433.03
	Office supplies	518021.43
	Technology	886313.52

All product categories report good profits, hence they do not contribute to the solution, but we can not identify the reasons for any loss the company is facing.

∴ Profit per product category is still an important subcategory metric that can be used for creating business reports, Profit per product subcategory is a useful metric as it shows subcategory that report heavy loss.

## # Profit per Product Subcategory

Select p. Product\_Category,

p. Product\_Subcategory,

Sum(m.profit) As Profit

From

Market\_fact\_full as M

Inner Join

Prod\_dimen as P

ON m.prod\_id = P.Prod\_id

Group By p.product\_category

p.product\_subcategory

Order by sum(m.profit)

# Exploring order\_id, order-number from orders table

> Select order\_id, order-number

from

orders-dimen

group by

order\_id,

order-number;

# From this we are not able to understand whether order id to order no is a unique mapping.

# Let's go with another approach

Select

Count(\*) as rec-count,

Count(Distinct order\_id) as order\_id-count,

Count(Distinct order-number) as order-number-count

from

orders-dimen

o/p

rec-count

order\_id-count

order-number-count

5506

5506

5496

→ Multiple order\_ids can have same order num