

# Database Management & Big Query

1) Why prefer database over excel?

1) Limit to 1 million rows.

2) Data in excel might not be secured.

3) Don't support concurrency.

4) Data integrity.

5) Performance.

## 2) Data Types in SQL

a) Int

b) Char(n) → fixed size string of length n.

c) Varchar(n) → variable size string of length upto n.

d) float(n,y) → n #s before decimal point & y #s after.

e) Date → yyyy-mm-dd

f) Time → hh-mm-ss

g) Datetime → yyyy-mm-dd hh-mm-ss

h) Timestamp → Datetime as an universal time co-ordinate.

## Keys

1) Primary key → unique & not null

2) Unique key → unique & can be null

3) Foreign key → It is a primary key in some other table.

4) Candidate key → It is a combination of columns to make a primary key.

## Query Code :-

1) Select & Select rows of table.

2) where & Conditional Filtering.

between & Select & from movies

where

year between x and y;

3) select & from movies

where title like "Toy Story %";

& Special operators.

- Match with two or more substrings.
- Can be used only with like or not like.

- where is always written after the from clause.

## String Manipulation functions :-

1) concat(x, y, ..., z)

2) lower/upper(x)

3) substr(x, start\_index, length) <sup>& length of substring</sup> to extract.

Q:- why = for comparison instead of == ?

- In sql there is no assignment operator so = is used instead.

- From other languages who have variable assignment  
== is used for comparison instead.

& used to skip rows.

& offset is always used after limit.

& limit is compulsory for offset to work.

Between & Between a certain range (both inclusive).

Ex: Find book assignments for vendor\_id & b/w dates  
2019-04-03 & 2019-05-16 &

Select \* from vendor\_book\_assignments  
where vendor\_id = ? and market\_date between '2019-01-03'  
and '2019-05-16';

Syntax of the code of order of execution &

Code &

Select  
from  
where  
orderby  
limit  
offset

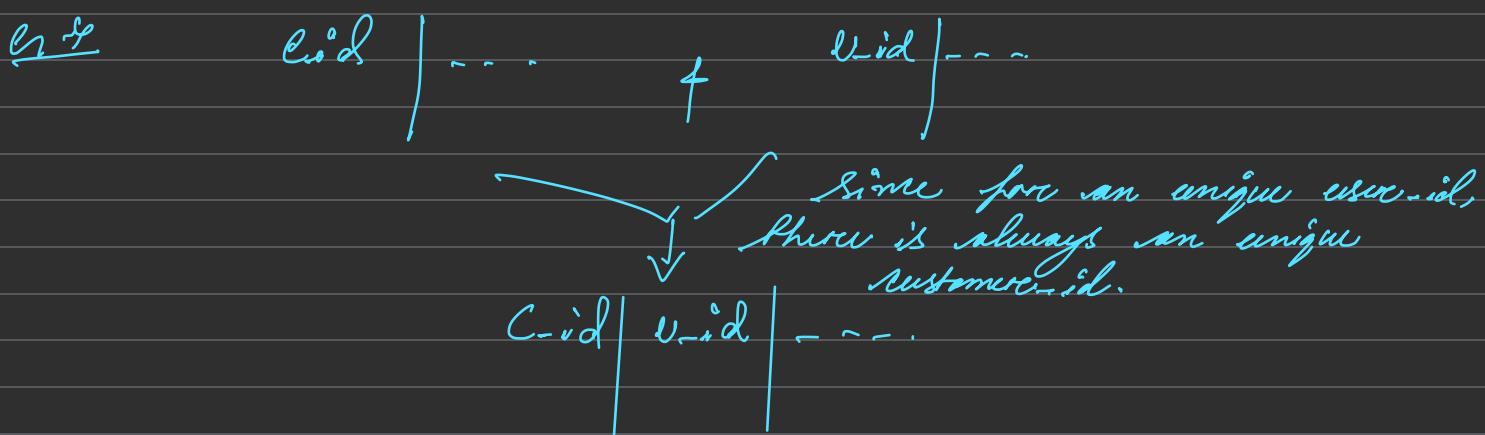
Order of Execution &

from  
where  
select  
orderby  
offset  
limit

Entity Relationships  $\rightarrow$  Relationship b/w tables in ER diagram.

1) one to one  $\rightarrow$   $++$   $\rightarrow$  one customer has one user-id.  
one user-id is assigned to one customer.

- very rarely seen in ER diagram.
- Most often, the two tables can be combined into one table.



2) one to many / many to one  $\rightarrow$   $++$   $\leftarrow$  many  
many  $\rightarrow$   $++$   $\leftarrow$  one

Ex ① A brand can make many products.

A product can be made by only one brand.  
 $\rightarrow$  brand  $++$   $\leftarrow$  products.

2) A teacher can teach many students.

A student can be taught by only one teacher in a class.  
 $++$   $\leftarrow$

3) Child of Parents  $\rightarrow$

1 child can have 1 parents.  $++$   $++$

1 parents can have many children.  $++$   $\leftarrow$   
 $\rightarrow$   $++$  (many to one)

3) Many to many ↗ ↘  
ent Book & Author ↗

| Book can have many authors.  
| authors can have many books.

→ ↘ (many to many).

↗ Languages involved in sql ↗

1) DDL -> Data definition language.

- Used for developmental work on a table.  
Drop/create/Insert/Alter table.

2) DML -> Data Manipulation Language.

- used for CRUD operations.

3) DCL / PCL -> Data f Transaction control language.

↳ Commit/Rollback  
↳ Grant/Revoke.

4) DQL -> Data query Language.

→ There are only 10 clauses in DQL.

select, from, where, orderby, limit,  
offset

# Filtering & Sub-queries &

where → filter things out.

Exact match  
 $(=, !=, in, not in)$

Partial match → used to partially  
match strings.

like

"%" → match any # of characters.  
\_ → match only one character.

Ex: select \*

from customers

where lower(first\_name) like "%ana%"

match all the customers who have  
ana in their penultimate position.

Distinct → Returns unique values in a column.

If 2 columns → Returns all the unique combinations  
of the columns.

→ Considers Null as an unique value.

ifnull ( expression/column, value to replace with ) →

if the value is null, replace it with value in 2nd parameter.

→ Subqueries & Executing queries inside queries.

↳ Subquery used in where section of outer query.  
↳ can be used to include multiple subqueries.  
↳ Subquery should return a list or a single value.  
↳ Inner subquery gets executed first.

## Window functions & (Aggregation functions)

- Always creates a new column.

- count(column), sum(), min/max, avg, lead, lag

Over → Creates a window over the specified conditions.

by      select \* , max(salary) over()  
          from employees;                          ↳ Considers the entire table as one window.

→ partition by → Divides the table into multiple windows based on the # of unique elements in the partition by column.

→ order by → Sort according to the order by columns in the partitioned window.

Groups  $\rightarrow$  slice  $\rightarrow$  Dividing the dataset into n groups of fixed sizes.

ex- $\uparrow$  select  $ntile(4)$  over(order by salary),  
salary  
from employees  
order by employees;

Approach  $\uparrow$

select  
from (select ...  $ntile(4--)$ )  
order by --.  
group by --.  
order by --.)  
 $\uparrow$  query2 (inside)  
- Table Pre-processing  
- Generated first.  
 $\uparrow$  query1 (outside)  
- Generated second.

2) lead & lag  $\uparrow$

lead  $\uparrow$  shift the columns up by n steps.

- Initiates bottom n rows to null.

lag  $\uparrow$  shifts the column down by n steps.

- Initializes the first n rows to null.

ex- $\uparrow$  month on month analysis  $\uparrow$

Compare sales of current month w.r.t sales(next month)

select  $ntile(4)$  over(order by month) as quarters,  
sum(sales) as sales  
lead(ntile(4) over(order by month), 1) as next-month  
from (select extract(month from marketdate) as  
month, sum(quantity) as sales  
from customer\_purchases  
group by 1  
order by 1  
)

group by 1  
order by 1;

→ Cumulative Sum using window function

<u>date</u>	<u>sales</u>	<u>1</u>	<u>2</u>
d <sub>1</sub>	100	100	100
d <sub>2</sub>	200	300	300
d <sub>3</sub>	100	400	400
d <sub>4</sub>	500	900	900

1) `sum(sales) over()` → Returns the total sum for each row.

2) `sum(sales) over (order by date)` → Returns the cumulative sum for each row.