

## CASE STATEMENT

- Evaluates a list of condition and returns a value when the first condition is met.

Syntax . start of logic

CASE

Result if condition  
is True

condition to be evaluated  $\leftarrow$  WHEN condition 1 THEN result 1

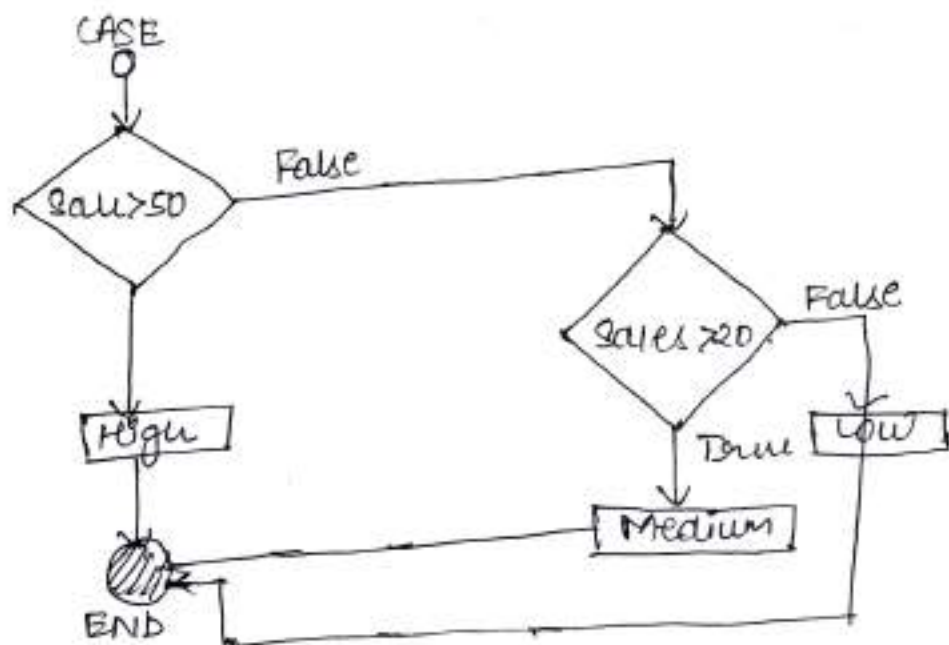
WHEN condition 2 THEN result 2

---

Default value  $\leftarrow$  ELSE Result  
(optional)

If none of the  
when conditions  
are True.

END  $\rightarrow$  End of logic



CASE

WHEN sales > 50 THEN 'HIGH'

WHEN sales > 20 THEN 'Medium'

ELSE 'LOW'

END

## # USE CASE

### 1) CATEGORIZING DATA

- Main purpose of case statement is Data Transformation.

#### # Derive New Information

- Create new column based on existing data.
- Group data into different categories based on certain conditions.

Ex: Report for Total Sales for each category:

High if sales > 50

Medium if sales > 20 AND sales > 50

Low if sales < 20

Sort result from lowest to Highest.

SELECT

category,

SUM(sales) AS Totalsales

FROM (

SELECT

ORDERID,

sales,

CASE WHEN sales > 50 THEN 'High'

WHEN sales > 20 THEN 'Medium'

ELSE 'Low'

END category

FROM sales\_order

dt

GROUP BY category

## CASE Statement :- Rules.

Data type of result must be matching.

### \* Mapping. Value

↓  
Transform the values from one form to another.

```
SELECT
```

```
Employee ID,
```

```
Firstname,
```

```
Lastname,
```

```
Gender,
```

```
CASE
```

```
  WHEN Gender = 'F' THEN 'Female'
```

```
  WHEN Gender = 'M' THEN 'Male'
```

```
  ELSE 'NOT available'
```

```
END GenderFullText
```

```
FROM sales - Employees.
```

Change M → Male

### \* CASE Statement - Quick Form

```
CASE
```

```
  WHEN Country = 'Germany' THEN 'DE'
```

```
  WHEN Country = 'India' THEN 'IN'
```

```
  WHEN Country = 'United States' THEN 'US'
```

```
  WHEN Country = 'France' THEN 'FR'
```

```
  ELSE 'n/a'
```

```
END
```

Full Form

use Country multiple times



CASE Country  $\rightarrow$  column name (only one).

WHEN 'Germany' THEN 'DE'

WHEN 'India' THEN 'IN'

WHEN 'United States' THEN 'US'

ELSE 'n/a'

END

### \* Handling NULLS

- Replace NULLS with specific value.
- NULLS can lead to inaccurate results, which can lead to wrong decision-making.

SELECT

customer\_id,

lastname,

score,

Avg CASE

WHEN score IS NULL THEN 0

ELSE score

END) OVER() AvgScore

FROM sales.customers.

### \* Conditional aggregation - use case.

Apply aggregate functions only on subsets of data that fulfill certain conditions.

SELECT

customer\_id,

SUM CASE

WHEN sales > 30 THEN 1

ELSE 0

END) TotalOrdersHighSales,

COUNT(\*) TotalOrders

FROM sales.orders

GROUP BY customer\_id

• aggregate functions.

```
SELECT  
customer_id,  
COUNT(*) AS total_order  
SUM(sales) AS Total-Sales  
AVG(sales) AS avg-sales  
MAX(sales) AS max highest-sales  
MIN(sales) AS lowest-sales  
FROM orders  
GROUP BY customer_id.
```

```
Score SELECT  
customer_id,  
COUNT(*) AS Total-n-Customer  
SUM(CASE  
    WHEN score IS NULL THEN 0  
    ELSE score  
END) AS Total-Score  
AVG(CASE  
    WHEN score IS NULL THEN 0  
    ELSE score  
END) AS avg-score  
FROM sales.customer  
GROUP BY customer_id.
```

## WINDOWS FUNCTION.

- perform calculations (e.g. aggregation) on a specific subset of data, without losing the level of details of ~~data~~ rows.

```
SELECT  
SUM(sales) Totalsales  
FROM sales.orders.
```

Total sales of order.

```
SELECT  
productID,  
SUM(sales) Total sales  
FROM sales.orders  
GROUP BY productID
```

Total sales per  
each product

```
SELECT  
orderID,  
orderDate,  
productID,  
SUM(sales) Total sales  
FROM sales.orders  
GROUP BY productID
```

# Group By Rule:

- All columns in SELECT must be included in Group By.

- Can't do aggregations and provide details at same time.

```
SELECT  
SUM(sales) OVER()  
FROM sales.orders.
```

# Result Granularity

Window function return a result per each row.

```
SELECT  
SUM(sales) OVER (PARTITION BY productID)  
FROM sales.orders
```

# Total sales by each product.

```
SELECT  
orderID,  
orderDate,  
productID,  
SUM(sales) OVER (PARTITION BY productID) AS  
Totalsales by product  
FROM sales.orders
```

Window function

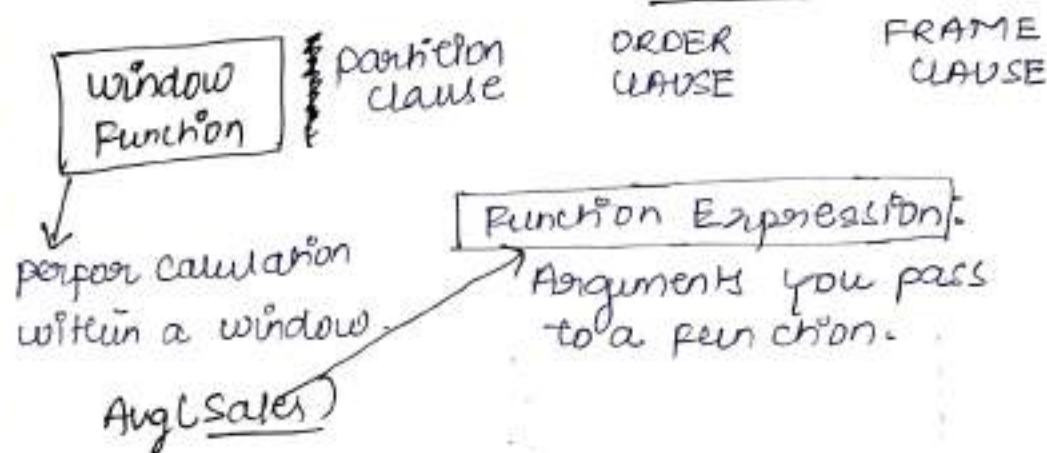
Advanced Data

Analysis (Aggregations  
+ Details)



## #THE SYNTAX.

### OVER CLAUSE



Empty ~~Rank()~~ ~~OVER~~ Rank() ~~ORDER BY~~ orderDate)

COLUMN Avg(Sales) OVER (ORDER BY orderDate)

NUMBER Ntile(2) OVER (ORDER BY orderDate)

Multiple Arguments LEAD(Sales, 2, 10) OVER (ORDER BY orderDate)

Conditional Logic SUM(CASE ) OVER (ORDER BY orderDate)

# OVER() CLAUSE : Tells SQL that the function used is a window function.

• It defines a window or subset of data.

### # PARTITION BY Category

↓  
partition clause

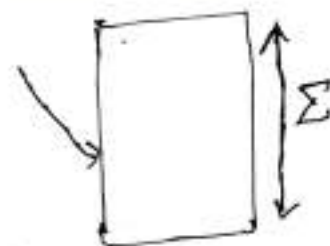
- Divides the dataset into windows (partition)
- Divides the result set into partitions (windows)

## # partition clause

divides the rows into groups, based on column/s.

$\text{SUM}(\text{sales}) \text{ OVER} ()$

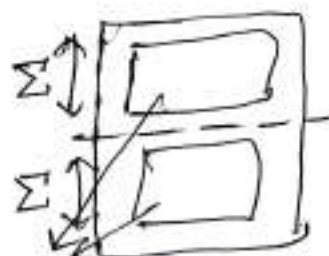
calculation done on entire dataset.



↓

```
SELECT
  orderID,
  orderDate,
  SUM(sales) OVER() Totalsales
FROM sales.orders
```

$\text{SUM}(\text{sales}) \text{ OVER} ()$   
 $\text{PARTITION BY productID}$



calculation done individually on each window.

```
SELECT
  orderID,
  orderDate,
  productID,
  SUM(sales) OVER(PARTITION BY productID)
  Totalsales Beam
FROM sales.orders
```

## Flexibility of window

Allow aggregation of data at different granularities within same query.

## # ORDER BY

Sort the data within a window

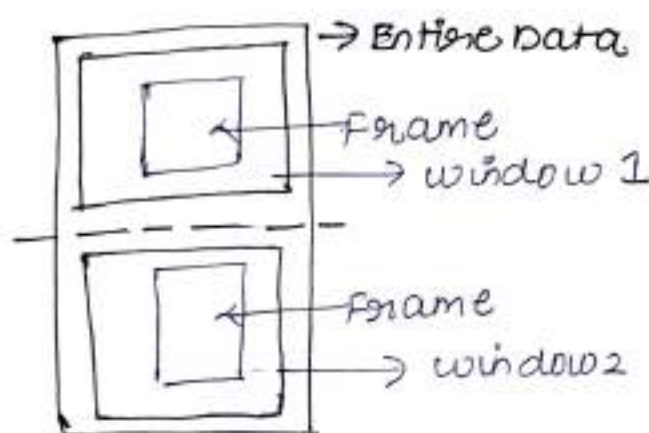
```
SELECT
  ORDERID,
  sales,
  RANK() OVER(ORDER BY sales DESC)
FROM sales.orders
```



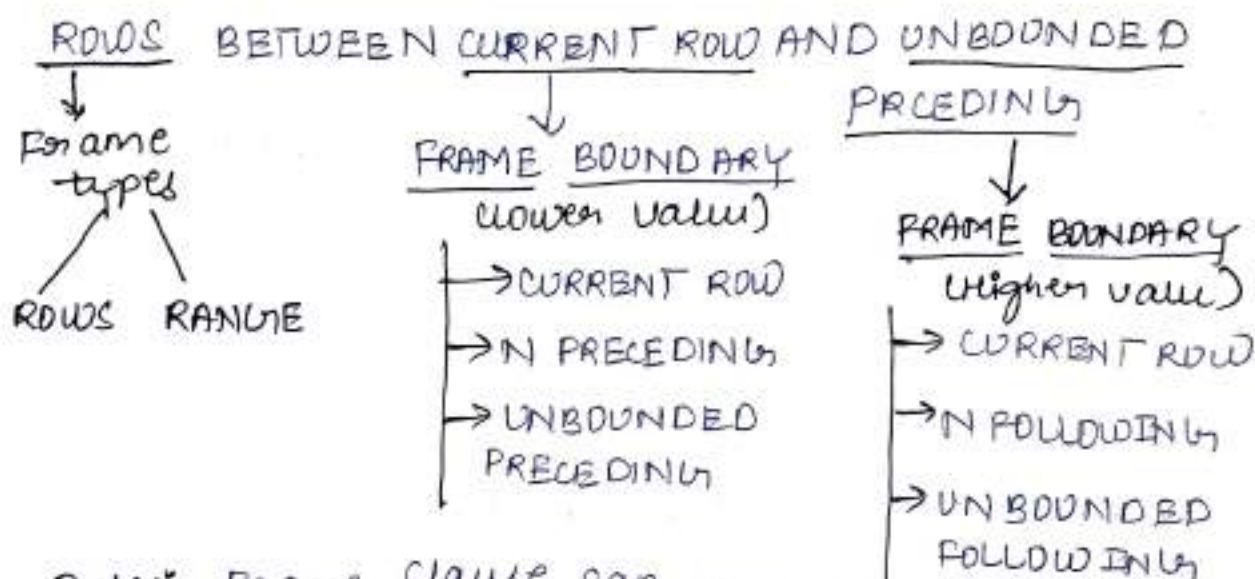
## #WINDOW FRAME

Define a subset of rows in a window. (each) that is relevant for calculation.

### ROW UNBOUNDED PRECEDING



AVG(Sales) OVER (PARTITION BY category ORDER BY orderdate)



Rules:- Frame clause can only be used together with order by clause.

- Lower value must be before the higher value

N-following : The N-th row before the current row.

unbounded following : last possible row within a window.

1-PRECEDING : n-th row before the current row.

unbounded preceding : first possible row within a window.

# compact frame.

For only preceding the current row can be skipped.

Short form ROWS 2 PRECEDING

ROWS UNBOUNDED PRECEDING.

Default frame : SQL uses default frame, if ORDER BY is used without FRAME.

↓  
ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW.

# RULE

- 1) window functions can be used ONLY in SELECT and ORDER BY CLAUSE.
- 2) Nesting window function NOT allowed.
- 3) SQL execute WINDOW function after WHERE clause.
- 4) window function can be used together with group by - in the same query, only if same columns are used.



## Window Aggregate Function.

AVG(Sales) OVER (PARTITION BY ProductID ORDER BY Sales)  
↓  
Required only Numeric)  
Both is optional

# only count is accept all data types. And all function only accept Numeric.

# COUNT() : Return the no. of rows within a window.

COUNT(Sales) : count numbers of Non-Null values in the column.

COUNT() = COUNT(\*) Both are same.  
[Count total numbers of rows including duplicates, not unique values.]

#1 use case :- (Overall analysis) Quick summary or Snapshot of entire dataset.

#2 use case :- (Total per groups) Group wise analysis, to understand within different categories.

↓  
SELECT  
orderID,  
orderDate,  
customerID,  
COUNT(\*) OVER ( ) Totalorder  
COUNT(\*) OVER (PARTITION BY customerID) orderbycust.  
FROM Sales.orders  
(1 use case)  
(2 use case)

#3 use case : (Data Quality check) Detecting numbers of nulls by comparing to total numbers of rows.



SELECT  
A,

COUNT(\*) OVER() Total cust,  
COUNT(subno) OVER() Total subno,  
COUNT(country) OVER() Total countries,  
FROM Sales.customers.

## # Data Quality Issue

- Duplicate leads to inaccurate in analysis.  
COUNT() can be used to identify duplicates.

```
SELECT  
    orderID,  
    COUNT(*) OVER(PARTITION BY orderID) checkPK  
FROM Sales.orders.
```

↓  
SELECT Duplicate id (bigger than 1)

```
SELECT  
* FROM (  
    SELECT  
        orderID,  
        COUNT(*) OVER(PARTITION BY orderID)  
        FROM Sales.orders AS ducive checkPK  
    )  
WHERE checkPK > 1
```

## # use case 4: Identify Duplicates :

Identify duplicate rows to improve data quality.

#SUM : Returns the sum of values within a window.

#use case 1 : Overall analysis :- Quick summary or snapshot of entire dataset.

#use case 2 : Total per groups :- group-wise analysis, to understand patterns within different categories.

↓ SELECT  
OrderID,  
OrderDate,  
ProductID,  
Sales,

SUM(Sales) OVER() TotalSales,

SUM(Sales) OVER(PARTITION BY ProductID) Total per Product  
By Sales.

FROM Sales.Orders

use case 1

use case 2

- part to whole analysis  
→ compare current sale to Total sales
- compare to Extreme analysis  
→ compare current sales to Highest or lowest
- compare to average analysis  
→ Help to evaluate whether a value is above or below the average.

Ex:- percentage contribution (part to whole analysis)

SELECT

OrderID,  
~~Order~~ ProductID,  
Sales,

SUM(Sales) OVER() TotalSales,

ROUND(CAST(Sales AS FLOAT)/SUM(Sales) OVER()  
\* 100, 2) percentage of Total.

FROM Sales.Orders



# AVG(): Return average values within a window.

\* If null in rows first remove NULL/  
Handle NULL and then do avg()

Ex: SELECT orderid, orderdate, sales, #1 use case: overall analysis: quick summary of entire data.

↓  
AVG(sales) OVER () AvgSales

AVG(sales) OVER (PARTITION BY productid) avgSales by products

↓  
#2 use case: Total per Groups: Group-wise analysis, to understand patterns within different categories

Ex: SELECT customerid, lastname, score, \* TO HANDLE NULLS  
COALESCE(score, 0),  
AVG(COALESCE(score, 0) OVER ()) AvgScoreWithNull  
FROM sales.customers

Ex: compare to average

SELECT -

\* FROM (

SELECT

#3 use case:

←

orderid,  
productid,  
sales,

AVG(sales) OVER () AvgSales  
FROM sales.orders

)

WHERE sales > AvgSales

compare to average:  
Help to evaluate whether  
a value is above below  
the average.



# MIN & MAX: Returns highest value in window.

↓  
Returns lowest value in the window.

Ex: SELECT  
orderID,  
productID,  
sales,

#1 use case: overall analysis:  
Quick summary of entire data.

MAX(sales) OVER() highest sales

MIN(sales) OVER() lowest sales

MAX(sales) OVER(PARTITION BY productID) highest sales by product  
MIN(sales) OVER(PARTITION BY productID) lowest sales by product  
FROM sales.orders

#2 use case: Total per  
groups: group-wise analysis, to  
understand patterns within different  
categories

Ex: Employee with highest salaries.

SELECT

\*  
FROM C

SELECT

\*,

MAX(salary) OVER() highest salary  
FROM sales.employees

) +

WHERE salary = highest salary

Ex: SELECT  
orderID,  
productID,  
sales,

MAX(sales) OVER() highest sales,

MIN(sales) OVER() lowest sales,

sales - MIN(sales) OVER() deviation from min

sales - MAX(sales) OVER() deviation from max

MAX(sales) - sales OVER() deviation from max

#3 use case: compare to extreme  
Help to evaluate how well  
a value is performing  
relative to the extremes

## # Distance from Extreme

- Lower the deviation, the closer the data point is to the extreme.

## # Analytical MSE case

- Running & Rolling Total.

• Tracking: Tracking current sales with Target sales.

- Trend analysis: providing insights into historical patterns.

\* They aggregate sequence of members, and the aggregation is updated each time a new member is added.

↳ Analysis over time

# Running Total: Aggregate all values from beginning up to the current point without dropping off older data.

# Rolling Total: aggregate all values within a fixed time window (eg 30 days) As new data is added, oldest data point will be dropped.

Rolling/shifting window.

# MOVING AVERAGE :

RUNNING AVERAGE

SELECT  
orderID,  
productID,  
orderDate,  
sales,

Avg(salu) OVER (PARTITION BY productID) avg, pro.

Avg(salu) OVER (PARTITION BY productID ORDER BY order  
Date) moving avg.

FROM sales.orders.

S:- calculate moving avg of sales for each product  
over time

Q. Including only next order.



customize frame : ROWS BETWEEN CURRENT ROW AND  
1 FOLLOWING.



# RANKING WINDOW FUNCTION

Task: Rank product based on their sales.

product	sales
E	70
B	30
A	20
C	10
D	5

percentage-based ranking

0	continuous value.
0.25	
0.5	
0.75	
1	

product	sales	Integ.-based Ranking	
E	70	1	discrete value
B	30	2	
A	20	3	
C	10	4	
D	5	5	

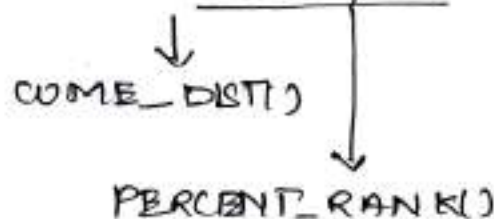
Find top 5

Find 20%

contribution of data

Top / Bottom N analysis

Distribution analysis



Syntax

RANK() OVER (PARTITION BY productID ORDER BY sales)

optional

Required

exp. must be ency

# ROW\_NUMBER. Assign a unique number to each row.

• It doesn't handle ties.

Two row sharing same value, they will not share same rank.

Sales	Rank	
100	1	→ unique Ranking without gap / skipping.
80	2	
80	3	ROW-NUMBER() OVER(ORDER BY sales DESC)
50	4	
20	5	→ doesn't handle ties.

Ba:

```

SELECT
    orderid,
    productid,
    sales,
    ROW-NUMBER() OVER(ORDER BY sales DESC) +
FROM sales.orderid.

```

### # RANK.

- Assign a Rank to each row.
- It handle ties
- It leaves gap in ranking.

Sales	Rank	
100	1	→ shared Ranking.
80	2	
80	2	→ leaving gaps (Ranking)
50	4	→ Handle ties
20	5	• ROW-NUMBER() OVER(ORDER BY sales DESC)

B<sub>2</sub>:

```

SELECT
    orderid,
    sales,
    RANK() OVER(ORDER BY sales DESC)
FROM sales.orderid.

```

### # DENSE-RANK()

- Assign a ~~new~~ rank to each row.
- Handle ties.
- doesn't leave gaps in ranking.

Sales	Rank	
100	1	• shared Ranking.
80	2	• Handle ties,
80	2	• Don't leave gap in Ranking.
50	3	
20	4	DENSE-RANK() OVER(ORDER BY sales DESC) +

Ex SELECT  
\*,  
DENSE\_RANK() OVER(ORDER BY sales DESC) +  
FROM ~~order~~ sales.orders.

## # TOP-IN-ANALYSIS

1) Top highest sales for each product.

```
SELECT
  order ID,
  product ID,
  sales,
  ROW_NUMBER() OVER(PARTITION BY product ID +
    ORDER BY sales DESC) Rankby product
FROM sales.orders
) +  $\xrightarrow{\text{Subquery}}$ 
SELECT
  *
FROM (
  +  $\xrightarrow{\text{Subquery}}$ 
```

WHERE Rankby product = 1

# use case: analyze Top performers to do Targeted Marketing.

2) BOTTOM-IN-ANALYSIS.

lowest 2 customers based on their sales.

# use case: help analyze the underperformance to manage risks and to do optimization.

Rule: columns used in group by and window function must be the same.

```
SELECT
  *
FROM (
  SELECT
    customer ID,
    SUM(sales) Totalsales,
    ROW_NUMBER() OVER(ORDER BY SUM(sales)) Rankcust.
  FROM sales.orders
) +
WHERE Rankcust <= 2
```



## # Generate unique IDs

• assign unique ID to row of 'Orders Archive' Table.

#use case: Helps to assign a unique identifier for each row to help paginating.

Ex:- SELECT  
ROW NUMBER() OVER (ORDER BY OrderID) uniqueID,  
\*  
FROM sales.Orders Archive.

Paginating: process of breaking down a large data into smaller, more manageable chunks.

## # Identify Duplicates.

#use case: - Identify and Remove Duplicates to improve data quality.

Ex: Identify duplicate row. return a clean

Result. SELECT  
\*  
FROM (

SELECT  
ROW NUMBER() OVER (PARTITION BY OrderID ORDER  
BY creationTime DESC) rn,

\*  
FROM sales.Orders Archive  
) t

WHERE rn = 1 / TO collect only bad data.

WHERE rn > 1

# NTILE Divides the rows into a specified number of approximately equal groups (buckets)

NTILE(2) OVER (ORDER BY sales DESC)

sales	NTILE
100	1
80	1
80	1
50	2
30	2

BUCKET SIZE =  $\frac{\text{NO. of ROW}}{\text{NO. of BUCKET}}$

$$2 = 5/2$$

SQL rule  
Larger group comes first.

Ex:

SELECT

orderid,  
Sales,

NTILE(3) OVER(ORDER BY Sales DESC) 3 bucket  
FROM Sales.Orders.

# \* NTILE USE CASE

# Data segmentation : Divides the datasets into distinct subsets based on certain criteria.

Ex: - segment all order into 'High, Medium, Low'

SELECT  
\*,

CASE WHEN buckets = 1 THEN 'High'  
WHEN buckets = 2 THEN 'Medium'  
WHEN buckets = 3 THEN 'Low'

END series segmentation

FROM ( SELECT

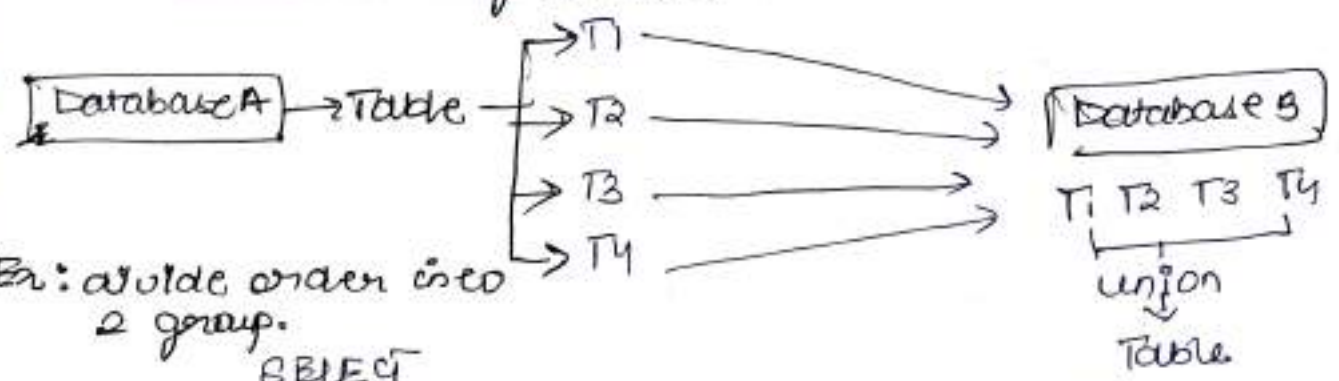
orderid,  
Sales,

NTILE(3) OVER(ORDER BY Sales DESC)  
buckets

FROM Sales.Orders

) +

## # Load Balancing NTILE(4)



Ex: divide order into 2 group.

SELECT

NTILE(2) OVER(ORDER BY orderid) buckets,  
\*

FROM Sales.Orders



## PERCENT AGE BASED RANKING

# CUMDIST: Cumulative distribution calculates the distribution of data points within a window.

Sales	DIST
100	0.2
80	0.6
80	0.6
50	0.8
30	1

$$\text{CUMDIST} = \frac{\text{position NA}}{\text{No. of Rows}}$$

$$= \frac{1}{5}$$

# Rule  
position of last occurrence of the same value

CUMDIST() OVER(ORDER BY sales DESC)

# PERCENT-RANK: calculate the relative position of each row.

Sales	Pos
100	0
80	0.25
80	0.25
50	0.75
30	1

$$\text{percent\_rank} = \frac{\text{position NR} - 1}{\text{NO. of ROW} - 1}$$

$$= \frac{1-1}{5-1} = \frac{0}{4} = \frac{3}{4}$$

$$= \frac{4-1}{5-1} = \frac{1}{4}$$

# Rule: position of first occurrence of the same value.

### Inclusive

- current row is included.

### Exclusive

- current row is excluded.

Ex- Find product that fall within the highest 10% of price.



SELECT  
\*,

CONCAT(DISTRANK \* 100, '%') DISTRANKPerc  
FROM

SELECT  
product,  
price,

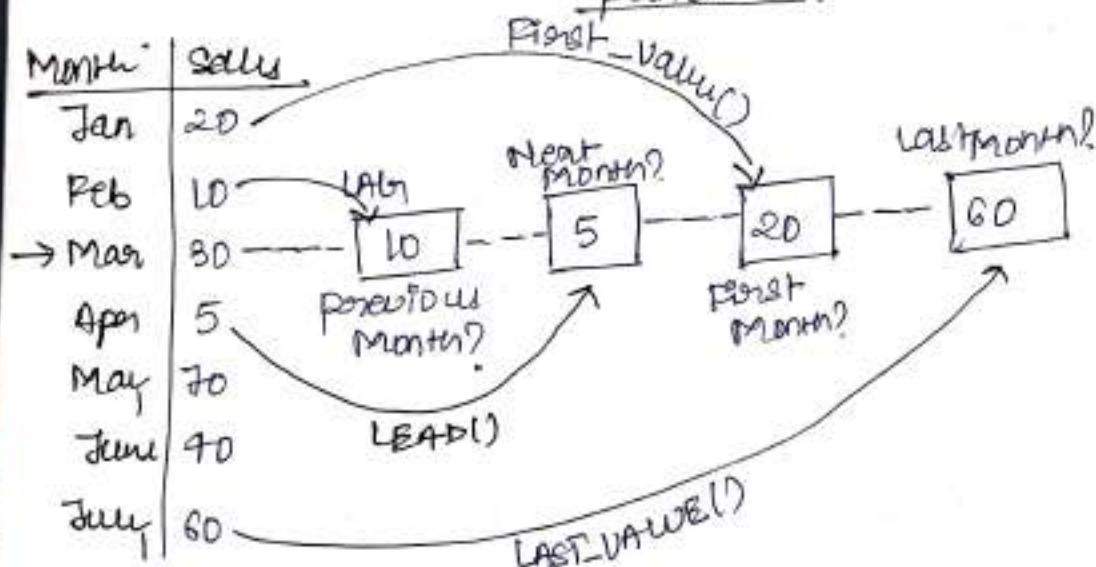
CUME-DIST() OVER(ORDER BY PRICE DESC)  
DISTRANK

FROM sales.products

+

WHERE DISTRANK < 0.4

Value window analytical  
function.



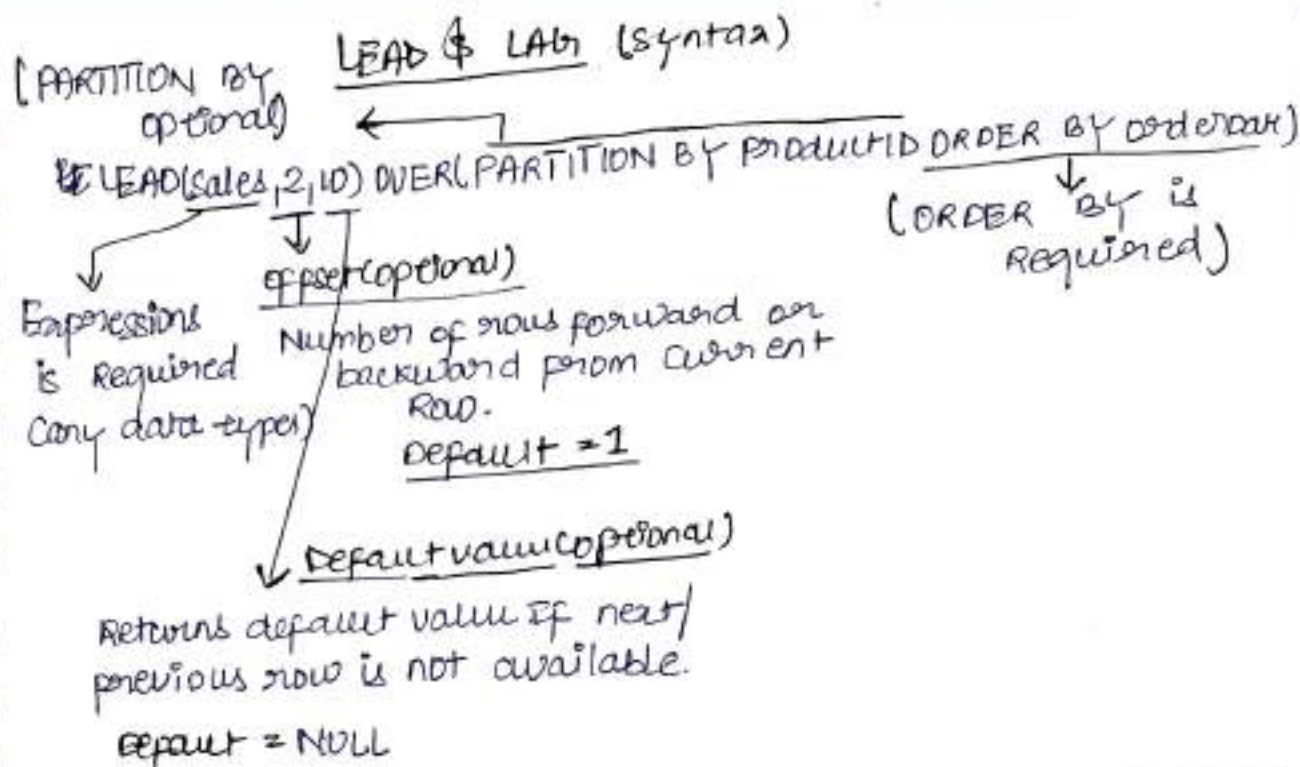
# Access value from other row.

# MIN/MAX

LEAD(): Access value from NEXT ROW.

LAG(): Access value from previous ROW.

Syntax.



LEAD(sales, 2, 0) OVER(ORDER BY MONTH) LAG(sales, 2, 0) OVER(ORDER BY MONTH)

Month	sales	LEAD
Jan	20	30
Feb	10	5
Mar	30	0
Apr	5	0

Month	sales	LAG
JAN	20	0
Feb	10	0
Mar	30	20
Apr	5	10

### #USE CASE

(MOM) MONTH-OVER-MONTH analysis.

Ens. MOM performance by finding % change in sales b/w current and previous month.

### #TIME SERIES ANALYSIS.

process to analyzing the data to understand pattern, trends and behaviours over time.

### \*YEAR-OVER-YEAR (YOY)

Analyse the overall growth or decline the business performance over years/time.