

SQL window functions

Window functions in SQL

A

Types of window functions

Let's have a closer look at some specific examples of window functions, and how to use them.

Aggregate

- AVG()
- MAX()
- MIN()
- SUM()
- COUNT()

Ranking

- ROW_NUMBER
- RANK()
- DENSE_RANK()
- NTILE()

Analytical

- LAG()
- LEAD()
- FIRST_VALUE()
- LAST_VALUE()

2

Data overview

We will use the following table called Employee that contains information about employees in a company located in South Africa. We assume the database is selected, so we don't specify it.

Start_date	Department	Province	First_name	Gender	Salary
2015-06-01	Finance	Gauteng	Lily	Female	35760
2020-08-01	Marketing	Western_Cape	Gabriel	Male	30500
2022-03-01	Data_analytics	Free_State	Maryam	Female	46200
2022-07-15	Marketing	Gauteng	Sophia	Female	36900
2019-05-01	Data_analytics	Western_Cape	Alex	Male	36200
2012-01-01	Finance	Free_State	Martha	Female	48500
2014-05-01	Finance	Western_Cape	Joshua	Male	35760
2017-06-15	Data_analytics	Gauteng	Emily	Female	37800
2016-01-01	Marketing	Western_Cape	David	Male	31000

A

Aggregate window functions

Suppose we want to see the **minimum salary** in each **department**. We can use MIN() as a window function and return the minimum value to each row.

```
SELECT
                   Department,
                   First_name,
                                                                             Using ORDER BY here
                   Salary,
                                                                             sorts the result set by
                   MIN(Salary) OVER (
Query
                                                                             department. This has no
                      PARTITION BY Department) AS Min_salary
                                                                             effect on how the
                FROM
                                                                             window function works.
                   Employee
                ORDER BY
                    Department;
```

Aggregate window functions

Department	First_name	Salary	Min_salary
Data_analytics	Maryam	46200	36200
Data_analytics	Alex	36200	36200
Data_analytics	Emily	37800	36200
Finance	Joshua	35760	35760
Finance	Lily	35760	35760
Finance	Martha	48500	35760
Marketing	Sophia	36900	30500
Marketing	David	31000	30500
Marketing	Gabriel	30500	30500

The lowest salary in each department is returned to the correct rows.

The minimum value resets for Finance and Marketing since each department is a different partition/window.

MIN(), MAX(), AVG(), SUM(), and COUNT() behave in a similar way when used as window functions.

5

/\

Aggregate window functions

Suppose we want to see the **total salary budget per department** as a **running total**. Using SUM() as a window function we can add up Salary in each partition separately. Adding ORDER BY will give a **running total** in each row.

```
SELECT
                 Department,
                 First_name.
                                                                            Adding ORDER BY here
                 Salary,
                                                                            enables the running sum
                 SUM(Salary) OVER (
Query
                                                                           calculation by date. All
                   PARTITION BY Department
                                                                            aggregate functions can
                   ORDER BY Date_Started) AS Dept_salary_budget
                                                                            be modified like this.
             FROM
                 Employee
             ORDER BY
                 Department:
```

Aggregate window functions

Date _started	Department	First _name	Salary	Dept_salary _budget
2012-01-01	Finance	Martha	48500	48500 01.
2014-05-01	Finance	Joshua	35760	84260 02.
2015-06-01	Finance	Lily	35760	120020 03.
2017-06-15	Data_analytics	Emily	37800	37800 04.
2019-05-01	Data_analytics	Alex	36200	74000
2022-03-01	Data_analytics	Maryam	46200	120200
2016-01-01	Marketing	David	31000	31000
2020-08-01	Marketing	Gabriel	30500	61500
2022-07-15	Marketing	Sophia	36900	98400

Employees' salaries are added together in each department in a running total.

The total salary budget in 2012 was R48500, which was only for Martha.

In 2014 we hired Joshua and the total salaries increased to R84260.

Finally, in 2015 we hired Lily so now we need to spend R120020 on salaries per month in Finance.

Note that the sum resets for each department and calculates a new running total.

7

Ranking window functions

Ranking window functions assign a **rank** or **row number** to each row within a specified window or subset of rows. Ranking window functions typically need an ORDER BY clause in order to work as intended.

RANK()

Ranks rows in each specified partition.

Duplicate rows are assigned the same rank and the next **rank is skipped**, taking duplicates into account.

e.g. 1, 2, 2, 4.

DENSE_RANK()

Ranks rows in each specified partition.

Duplicate rows are assigned the same rank but the ranks are sequential regardless of duplicates.

e.g. 1, 2, 2, 3.

ROW_NUMBER()

Assigns a unique number to each row within each partition, even if values in the partitioned column(s) are duplicated.

e.g. The sequence is 1, 2, 3, ..., regardless of duplicates.

NTILE()

Divides sorted partitions into n-number of equal **groups**. Each row in a partition is assigned a group number 1, 2, 3...

e.g. if you have 100 rows of data, NTILE(4) will divide the data into 4 x 25 row groups.

8

A

The RANK() function

The RANK() function assigns a rank to each row based on the order specified within the window. Rows with the same values receive the same rank, and the next rank is skipped.

```
SELECT

First_name,
Province,
RANK() OVER (
ORDER BY Province) AS Rank_assign

FROM
Employee
ORDER BY
Province;
```

Suppose we want to rank each row by Province.

Note that we didn't use PARTITION BY here, so we will rank **all** rows together.

The RANK() function

First_name	Province	Rank_assign		
Martha	Free_State	1	01.	
Maryam	Free_State	1		
Emily	Gauteng	3	02.	5 ro
Lily	Gauteng	3		
Sophia	Gauteng	3		
Alex	Western_Cape	6	03.	
David	Western_Cape	6		
Gabriel	Western_Cape	6		
Joshua	Western_Cape	6		

Rows with the same values
(Free_State) receive the same rank
(rank 1).

Rank 2 is skipped because Maryam fell under rank 1.

Rank 4 and 5 are skipped because there are **5 rows** above.



The DENSE_RANK() function

The DENSE_RANK() function operates similarly to the RANK() function except it **does not skip** any ranks even if rows have the same values.

```
SELECT

First_name,
Province,
DENSE_RANK() OVER (

ORDER BY Province) AS Rank_assign

FROM
Employee
ORDER BY
Province;
```

Suppose we want to rank each row by Province, but keep a sequential list, not skipping ranks.

Note that we didn't use PARTITION BY here, so we will rank **all** rows together.

The DENSE_RANK() function

First_name	Province	Rank_assign	
Martha	Free_State	1	01.
Maryam	Free_State	1	
Emily	Gauteng	2	02.
Lily	Gauteng	2	
Sophia	Gauteng	2	
Alex	Western_Cape	3	03.
David	Western_Cape	3	
Gabriel	Western_Cape	3	
Joshua	Western_Cape	3	

Rows with the same values (Free_State), receive the same rank (rank 1).

Ranks are sequential here, even though there are duplicate values.

Since there are three unique values for Province, our rank goes up to 3.



Use DENSE_RANK() to avoid rank gaps and potential confusion. Use RANK() when maintaining relative differences between ranks is essential.



The ROW_NUMBER() function

The ROW_NUMBER() function assigns a **unique sequential number** to each row within a partition, regardless of the column values. It makes sure that **no two rows can have the same row number** within a division.

```
SELECT

First_name,
Province,
ROW_NUMBER() OVER (

PARTITION BY Province
ORDER BY First_name) AS Row_assign

FROM
Employee
ORDER BY
Province;
```

Suppose we want to assign a unique number to each employee in a department based on their First_name.

Note that we add PARTITION BY Department here to reset the function for each partition.

The ROW_NUMBER() function

First_name	Province	Row_assign	
Martha	Free_State	1	01.
Maryam	Free_State	2	
Emily	Gauteng	1	02.
Lily	Gauteng	2	7
Sophia	Gauteng	3	
Alex	Western_Cape	1	03.
David	Western_Cape	2	7
Gabriel	Western_Cape	3	
Joshua	Western_Cape	4	

Sequential row assignment in the Free_State partition.

The row sequence resets in the Gauteng partition.

Sequential row assignment in the Western_Cape partition.

Each employee has a **unique** row **number** in their respective departments.

Query



The NTILE() function

NTILE() divides sorted partitions into n-number of equal groups. Each row in a partition is assigned a group number.

```
SELECT
   First_name,
   Province,
   NTILE(2) OVER (
        PARTITION BY Province
        ORDER BY First_name) AS Group_number
FROM
   Employee
ORDER BY
   Province;
```

Suppose we want to divide employees from each province into two groups.

Since we are using PARTITION BY here, the data are split into partitions, and then NTILE(2) divides each Province into 2 groups, and assigns a group number to each employee.

The NTILE() function

First_name	Province	Group_number
Martha	Free_State	1 01.
Maryam	Free_State	2
Emily	Gauteng	1 02
Lily	Gauteng	1
Sophia	Gauteng	2
Alex	Western_Cape	1 03
David	Western_Cape	1
Gabriel	Western_Cape	2
Joshua	Western_Cape	2

O1. Employees are assigned a group number per department.

If rows cannot split equally, they are always assigned to the first group.

Here the partition can be equally sub-divided.

The LAG() function

The LAG(column, n) function allows access of a value within a column from the **previous** n^{th} -row **relative** to the **current row**.

```
SELECT

Department,
First_name,
Salary,

LAG(Salary,1) OVER (
ORDER BY Date_started) AS Previous_salary

FROM
Employee
ORDER BY
Department;
```

Suppose we want to retrieve the previous salaries according to the employee's date of hire.

The LAG(Salary, 1) function is applied to the Salary column and ordered by Date_started. It returns the salary from the previous row, since n = 1, as the column Previous_salary.

The LAG() function

Department	First_name	Salary	Previous_salary
Data_analytics	Maryam	46200	NULL
Data_analytics	Alex	36200	46200
Data_analytics	Emily	37800	36200
Finance	Joshua	35760	37800
Finance	Lily	35760	35760
Finance	Martha	48500	35760
Marketing	Sophia	36900	48500
Marketing	David	31000	36900
Marketing	Gabriel	30500	31000

The LEAD() function

The LEAD(column, n) function allows access of a value within a column from the **following** n^{th} -row **relative** to the **current row**. It is the counterpart of the LAG() function.

```
SELECT

Department,
First_name,
Salary,

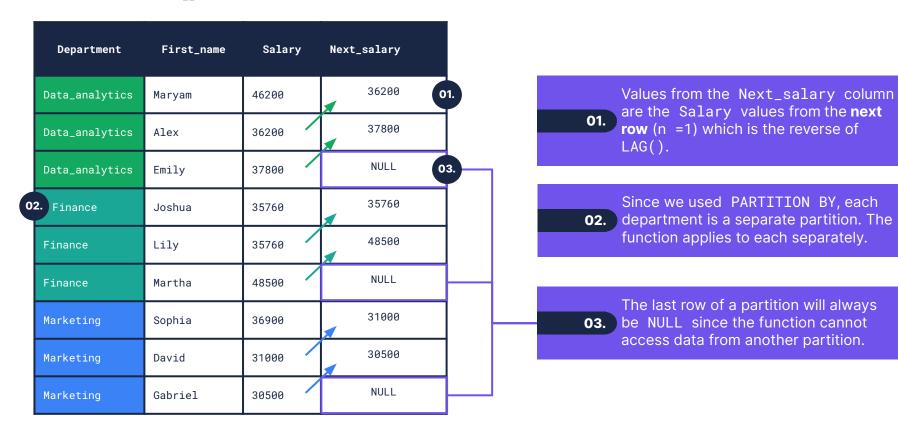
LEAD(Salary,1) OVER (
PARTITION BY Department
ORDER BY Date_started) AS Next_salary

FROM
Employee
ORDER BY
Department;
```

Suppose we want to retrieve the next employee's name according to the employee's date of hire.

The LEAD(Salary, 1) function is applied to the Salary column and ordered by Date_started. It returns the salary from the previous row, since n = 1, as the column Next_salary.

The LEAD() function



The FIRST_VALUE() function

The FIRST_VALUE() function allows the retrieval of the value of a column from the **first row** within a partition.

```
SELECT

Start_date,
Department,
First_name,
FIRST_VALUE(First_name) OVER (
ORDER BY Start_date
PARTITION BY Department) AS First_in_dept

FROM
Employee
ORDER BY
Department;
```

Suppose we want to retrieve the first employee the company hired.

The FIRST_VALUE() function is applied to the First_name column and ordered by Start_date. It returns the First_name from the first row as First_in_dept.

Window functions in SQL

The FIRST_VALUE() function

Start_date	Department	First _name	First _in_dept
2017-06-15	Data_analytics	Emily	Emily 01.
2019-05-01	Data_analytics	Alex	Emily
2022-03-01	Data_analytics	Maryam	Emily
2012-01-01	Finance	Martha	Martha 02.
2014-05-01	Finance	Joshua	Martha
2015-06-01	Finance	Lily	Martha
2016-01-01	Marketing	David	David
2020-08-01	Marketing	Gabriel	David
2022-07-15	Marketing	Sophia	David

Only the first value or the first hired employee Emily will be the output in the Data_analytics partition.

Since we used PARTITION BY

Department, the first employee from each department is returned.

The LAST_VALUE() function

The LAST_VALUE() function allows the retrieval of the value of a column from the **last row** within a window frame.

```
SELECT
Start_date,
Department,
First_name,
LAST_VALUE(First_name) OVER (
ORDER BY Start_date) AS Last_employee
FROM
Employee
ORDER BY
Department;
```

Suppose we want to retrieve the **last employee** the company hired.

The LAST_VALUE() function is applied to the First_name column and ordered by Start_date. It returns the First_name from the first row as Last_employee.

Window functions in SQL

The LAST_VALUE() function

Start_date	Department	First _name	Last _employee
2017-06-15	Data_analytics	Emily	Sophia
2019-05-01	Data_analytics	Alex	Sophia
2022-03-01	Data_analytics	Maryam	Sophia
2012-01-01	Finance	Martha	Sophia
2014-05-01	Finance	Joshua	Sophia
2015-06-01	Finance	Lily	Sophia
2016-01-01	Marketing	David	Sophia
2020-08-01	Marketing	Gabriel	Sophia
2022-07-15	Marketing	Sophia	Sophia

Only the last employee value, Sophia, will be the output. Since we didn't use PARTITION BY, the last employee among **all** rows is returned.