Adatacamp

SQL for Data Science

SQL Window Functions

Learn SQL online at www.DataCamp.com

Example dataset

We will use a dataset on the sales of bicycles as a sample. This dataset includes:

The [product] table

The product table contains the types of bicycles sold, their model year, and list price.

product_id	product_name	model_year	list_price
1	Treak 820 - 2016	2016	379.99
2	Ritchey Timberwolf Frameset - 2016	erwolf Frameset - 2016 2016 749.99	
3	Surly Wednesday Frameset - 2016	rly Wednesday Frameset - 2016 2016 999.9	
4	Trek Fuel EX 8 29 - 2016	el EX 8 29 - 2016 2016 2899.	
5	Heller Shagamaw Frame - 2016	2016 1320.99	

The [order] table

The order table contains the order_id and its date.

order_id	order_date	
1	2016-01-01T00:00:00.000Z	
2	2016-01-01T00:00:00.000Z	
3	2016-01-02T00:00:00.000Z	
4	2016-01-03T00:00:00.000Z	
5	2016-01-03T00:00:00.000Z	

The [order_items] table

The order_items table lists the orders of a bicycle store. For each order_id, there are several products sold (product_id). Each product_id has a discount value.

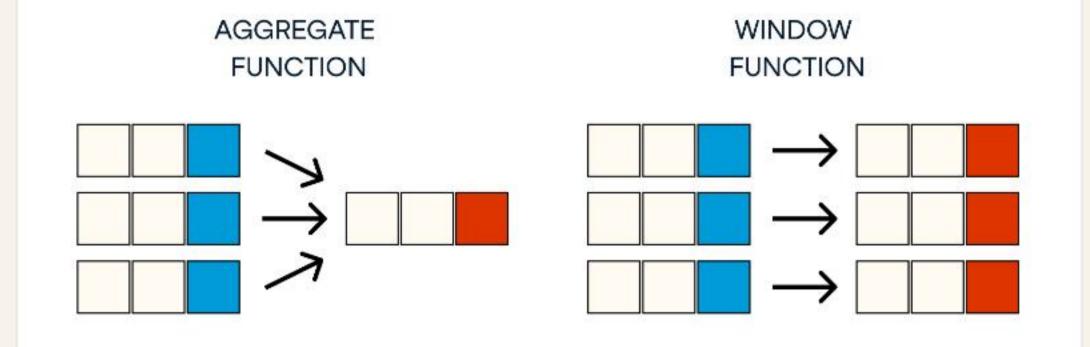
order_id	product_id	discount
1	20	0.2
1	8	0.07
1	10	0.05
1	16	0.05
1	4	0.2
2	20	0.07

What are Window Functions?

A window function makes a calculation across multiple rows that are related to the current row. For example, a window function allows you to calculate.

- Running totals (i.e. sum values from all the rows before the current row)
- 7-day moving averages (i.e. average values from 7 rows before the current row)
- Rankings

Similar to an aggregate function (GROUP BY), a window function performs the operation across multiple rows. Unlike an aggregate function, a window function does not group rows into one single row.



Syntax

Windows can be defined in the SELECT section of the query.

```
window_function() OVER(
     PARTITION BY partition_expression
     ORDER BY order_expression
    window_frame_extent
) AS window_column_alias
FROM table name
```

To reuse the same window with several window functions, define a named window using the WINDOW keyword. This appears in the query after the HAVING section and before the ORDER BY section.

```
window_function() OVER(window_name)
FROM table_name
[HAVING ...]
WINDOW window_name AS (
    PARTITION BY partition_expression
     ORDER BY order_expression
     window_frame_extent
[ORDER BY ...]
```

Order by

/* Rank price from HIGH->LOW */

ORDER BY is a subclause within the OVER clause. ORDER BY changes the basis on which the function assigns numbers to rows.

It is a must-have for window functions that assign sequences to rows, including RANK and ROW_NUMBER. For example, if we ORDER BY the expression 'price' on an ascending order, then the lowest-priced item will have the lowest rank.

Let's compare the following two queries which differ only in the ORDER BY clause.

product_name,			product_name,		
list_price,			list_price,		
RANK() OVER			RANK() OVER		
(ORDER BY list_pri	ce DESC) ran	(ORDER BY list_price ASC) ran			
FROM products			FROM products		
product_name ~	list_price ~	rank ~	product_name ~	list_price ∨	rank ~
Trek Domane SLR 9 Disc - 2018	11999.99	1	Strider Classic 12 Balance Bike - 2018	89.99	1
Trek Domane SLR 8 Disc - 2018	7499.99	2	Sun Bicycles Lil Kitt'n - 2017	109.99	2
Trek Domane SL Frameset - 2018	6499.99	3	Trek Boy's Kickster - 2815/2817	149.99	3

SELECT

/* Rank price from LOW->HIGH */

Partition by

We can use PARTITION BY together with OVER to specify the column over which the aggregation is performed.

Comparing PARTITION BY with GROUP BY, we find the following similarity and difference:

- Just like GROUP BY, the OVER subclause splits the rows into as many partitions as there are unique values in a column.
- · However, while the result of a GROUP BY aggregates all rows, the result of a window function using PARTITION BY aggregates each partition independently. Without the PARTITION BY clause, the result set is one single partition.

For example, using GROUP BY, we can calculate the average price of bicycles per model year using the following query.

SELECT	model_year ~	avg_price v
model_year,	20	6 980.29923
AVG(list_price) avg_price	26	7 1279.931176
1 products	20	1658.478441
BY model_year	20	9 2583.323333

What if we want to compare each product's price with the average price of that year? To do that, we use the AVG() window function and PARTITION BY the model year, as such.

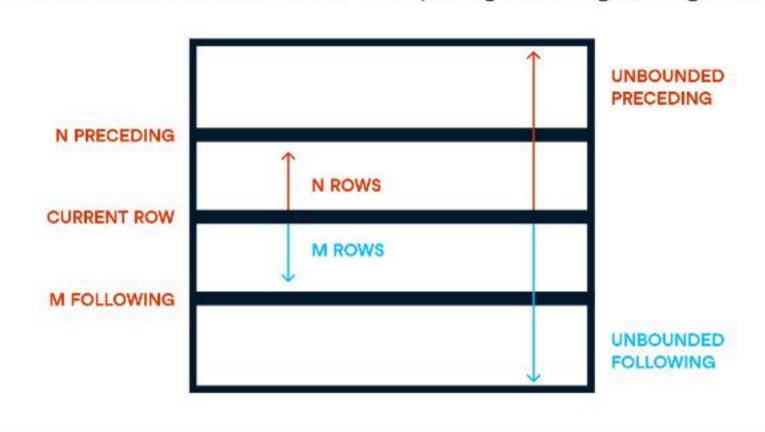
```
SELECT
model_year,
product_name,
list_price,
                                                          2017 Electra Amsterdam Fashion 7i Ladies' -
AVG(list_price) OVER
                                                         2017 Electra Ansterdan Original 3i -
                                                                                              659.99 1279.931176
  (PARTITION BY model_year)
avg_price
FROM products
```

Notice how the avg_price of 2018 is exactly the same whether we use the PARTITION BY clause or the GROUP BY clause.

Window frame extent

A window frame is the selected set of rows in the partition over which aggregation will occur. Put simply, they are a set of rows that are somehow related to the current row.

A window frame is defined by a lower bound and an upper bound relative to the current row. The lowest possible bound is the first row, which is known as UNBOUNDED PRECEDING. The highest possible bound is the last row, which is known as UNBOUNDED FOLLOWING. For example, if we only want to get 5 rows before the current row, then we will specify the range using 5 PRECEDING.



Accompanying Material

You can use this https://bit.ly/3scZtOK to run any of the queries explained in this cheat sheet.

U datacamp **SQL for Data Science** SQL Window Functions

Learn SQL online at www.DataCamp.com

Ranking window functions

There are several window functions for assigning rankings to rows. Each of these functions requires an ORDER BY sub-clause within the OVER clause.

The following are the ranking window functions and their description:

Function Syntax	Function Description	Additional notes
ROW_NUMBER()	Assigns a sequential integer to each row within the partition of a result set.	Row numbers are not repeated within each partition.
RANK()	Assigns a rank number to each row in a partition.	 Tied values are given the same rank. The next rankings are skipped.
PERCENT_RANK()	Assigns the rank number of each row in a partition as a percentage.	 Tied values are given the same rank. Computed as the fraction of rows less than the current row, i.e., the rank of row divided by the largest rank in the partition.
NTILE(n_buckets)	Distributes the rows of a partition into a specified number of buckets.	 For example, if we perform the window function NTILE(5) on a table with 100 rows, they will be in bucket 1, rows 21 to 40 in bucket 2, rows 41 to 60 in bucket 3, et cetera.
CUME_DIST()	The cumulative distribution: the percentage of rows less than or equal to the current row.	 It returns a value larger than 0 and at most 1. Tied values are given the same cumulative distribution value.

We can use these functions to rank the product according to their prices.

```
/* Rank all products by price */
   product_name,
   list_price,
   ROW_NUMBER() OVER (ORDER BY list_price) AS row_num,
   DENSE_RANK() OVER (ORDER BY list_price) AS dense_rank,
   RANK() OVER (ORDER BY list_price) AS rank,
   PERCENT_RANK() OVER (ORDER BY list_price) AS pct_rank,
  NTILE(75) OVER (ORDER BY list_price) AS ntile,
   CUME_DIST() OVER (ORDER BY list_price) AS cume_dist
FROM products
product_name ~
Strider Classic 12 Balance Bike
                                                                                       1 0.0031152648
- 2018
Sun Bicycles Lil Kitt'n - 2017
                                                                                       1 0.0062305296
                                                            2 2 0.003125
Trek Boy's Kickster - 2015/2017
                                                                                       1 0.0124610592
                                 149.99
                                                                         0.00625
Trek Girl's Kickster - 2017
                                 149.99
                                                                         0.00625
                                                                                       1 0.0124610592
Trek Kickster - 2018
                                                                           0.0125
                                                                                       1 0.015576324
                                 159.99
Trek Precaliber 12 Boys - 2017
                                                                         0.015625
                                                                                       2 0.0218068536
                                 189.99
```

0.015625

2 0.0218068536

189.99

Trek Precaliber 12 Girls - 2017

Value window functions

/* Find the difference in price from

the cheapest alternative */

FIRST_VALUE() and LAST_VALUE() retrieve the first and last value respectively from an ordered list of rows, where the order is defined by ORDER BY.

Value window function	Function
IRST_VALUE(value_to_return) OVER (ORDER BY value_to_order_by)	Returns the first value in an ordered set of values
LAST_VALUE(value_to_return) OVER (ORDER BY value_to_order_by)	Returns the last value in an ordered set of values
NTH_VALUE(value_to_return, n) OVER (ORDER BY value_to_order_by)	Returns the nth value in an ordered set of values.

To compare the price of a particular bicycle model with the cheapest (or most expensive) alternative, we can use the FIRST_VALUE (or LAST_VALUE).

/* Find the difference in price from

the priciest alternative */

```
product_name,
                                                          product_name,
                                                          list_price,
    list_price,
   FIRST_VALUE(list_price) OVER (
                                                         LAST_VALUE(list_price) OVER (
                                                           ORDER BY list_price
       ORDER BY list_price
       ROWS BETWEEN
                                                            ROWS BETWEEN
           UNBOUNDED PRECEDING
                                                               UNBOUNDED PRECEDING
           UNBOUNDED FOLLOWING
                                                               UNBOUNDED FOLLOWING
       ) AS cheapest_price,
                                                           ) AS highest_price
 FROM products
                                                      FROM products
                      HSt_price v cneapest_price v aim v
                                     89.99 8 Strider Classic 12 Balance Bike
Sun Bicycles Lil Kitt'n - 2017
                                                 Sun Bicycles Lil Kitt'n - 2017
```

Aggregate window functions

Aggregate functions available for GROUP BY, such as COUNT(), MIN(), MAX(), SUM(), and AVG() are also available as window functions.

Function Syntax	Function Description	
COUNT(expression) OVER (PARTITION BY partition_column)	Count the number of rows that have a non- null expression in the partition.	
MIN(expression) OVER (PARTITION BY partition_column)	Find the minimum of the expression in the partition.	
MAX(expression) OVER (PARTITION BY partition_column)	ARTITION BY Find the maximum of the expression in the partition.	
AVG(expression) OVER (PARTITION BY partition_column)	Find the mean (average) of the expression in the partition.	

Suppose we want to find the average, maximum and minimum discount for each product, we can achieve it as such.

SELECT						
order_id,						
product_id,						
discount,						
AVG(discount)	OVER	(PARTITION	BY	product_id)	AS	avg_discount,
MIN(discount)	OVER	(PARTITION	BY	product_id)	AS	min_discount,
MAX(discount)	OVER	(PARTITION	BY	product_id)	AS	max_discount
FROM order_items						

FROM order_it	cems				
order_id ∨	product_id ~	discount ~	avg_discount ~	min_discount ~	max_discount ~
2	16	0.05	0.113191	0.05	0.2
2	20	0.07	0.11253	0.05	0.2
3	20	0.05	0.11253	0.05	0.2
3	3	0.05	0.105581	0.05	0.2

LEAD, LAG

The LEAD and LAG locate a row relative to the current row.

Function Syntax	Function Description
LEAD(expression [,offset[,default_value]]) OVER(ORDER BY columns)	Accesses the value stored in a row after the current row.
LAG(expression [,offset[,default_value]]) OVER(ORDER BY columns)	Accesses the value stored in a row before the current row.

Both LEAD and LAG take three arguments:

- Expression: the name of the column from which the value is retrieved
- Offset: the number of rows to skip. Defaults to 1. Default_value: the value to be returned if the value retrieved is null.
- Defaults to NULL.

With LAG and LEAD, you must specify ORDER BY in the OVER clause.

LEAD and LAG are most commonly used to find the value of a previous row or the next row. For example, they are useful for calculating the year-on-year increase of business metrics like revenue.

Here is an example of using lag to compare this year's sales to last year's.

```
/* Find the number of orders in a year */
WITH yearly_orders AS (
 year(order_date) AS year,
 COUNT(DISTINCT order_id) AS num_orders
 FROM sales.orders
 GROUP BY year(order_date)
/* Compare this year's sales to last year's */
SELECT
```

LAG(num_orders) OVER (ORDER BY year) last_year_order, LAG(num_orders) OVER (ORDER BY year) - num_orders diff_from_last_year FROM yearly_orders

diff_from_last_year ~	last_year_order ~	num_orders ~	year ~
null	null	635	2016
-53	635	688	2017
396	688	292	2018

Similarly, we can make a comparison of each year's order with the next year's.

```
/* Find the number of orders in a year */
WITH yearly_orders AS (
 year(order_date) AS year,
 COUNT(DISTINCT order_id) AS num_orders
 FROM sales.orders
 GROUP BY year(order_date)
/* Compare the number of years compared to next year */
```

SELECT *,							
LEAD(num_orders)	OVER	(ORDER	BY	year)	next_year_ord	der,	
LEAD(num_orders) FROM yearly_orders	OVER	(ORDER	ВУ	year)	- num_orders	diff_from_next_year	
rkon year cy_orders							

diff_from_next_year ~	next_year_order ~	num_orders ~	year ~
53	688	635	2016
-396	292	688	2017
null	null	292	2018