

# SQL Window functions basics

With PostgreSQL



Clément Prévost [Follow](#)

May 26, 2017 · 6 min read



*For many developers, window functions are this mystical feature they only use once a year to build the famous Top-k query. In fact, they are part of the sql:2003 standard (yes, that's 10 years from now) and your DBMS should already have a great and optimised support for windowing functions.*

If you should keep only one relevant information from this article this must be that windowing functions are a select post-processing toolset. It allow us to iterate over the result of a select to compute values based on a wider view than just one row. You can use them to create groups of rows on which you may filter or create various aggregations. With this tool, you enable yourself to reason about the “next row” or “this group of rows” for each row in a select result, which is impossible in SQL 92.

## Window functions concepts

Window functions are calculated between the having and order by clauses. For each row of the partition, one can calculate a result based on the partition data. The processing order of the various SELECT parts is quite important to fully understand the intermediate data schema we are working with.

```
1  7 - SELECT
2  1 - [ FROM ... ]
3  2 - [ WHERE ... ]
4  3 - [ GROUP BY ... ]
5  4 - [ HAVING ... ]
6  5 - WINDOW window_name AS (
7      [ PARTITION BY ... ]
8      [ ORDER BY ... ]
9      [ RANGE ... ]
10 )
11 6 - [ ORDER BY ... ]
```

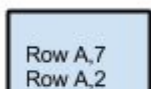
pg\_partition\_6.sql hosted with ❤ by GitHub

[view raw](#)

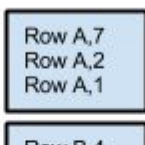
To be simplistic, window functions allow you to iterate over the result of your select before ordering the final result.

Window function works on row groups called partitions or windows, you can control the ordering and the window frame or range of these windows. A schema will explain it better than a thousand word:

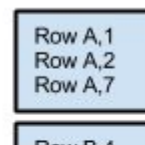
Select intermediate result  
(after having, before order by)

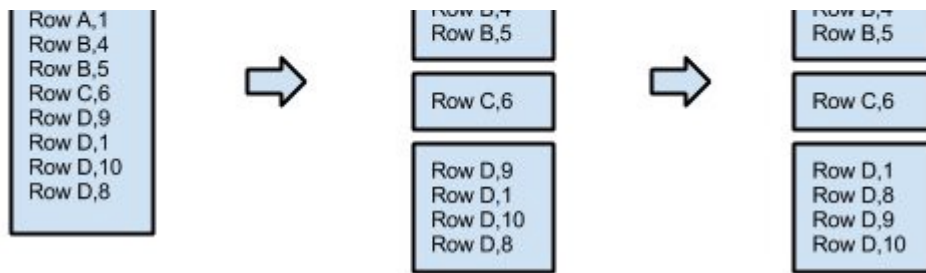


Partitioning

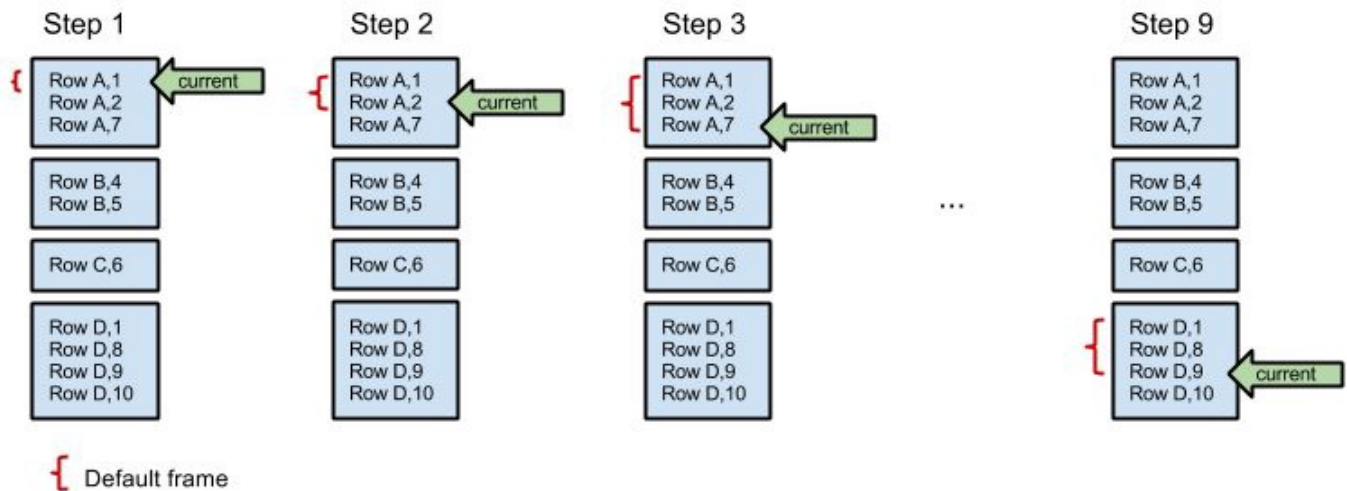


Ordering





On this first schema, the window function engine post-process the select result to break it into many parts called partitions. These partitions may have an inner ordering, different from the final result ordering.



This schema describe the behavior of the default window frame. The engine iterate through the result rows and store on each row the result of the aggregate function within the window frame.

By default, the frame goes from the first partition row to the current row, declaring an ORDER BY clause changes the default frame to the full partition. (cf. examples below) Due to this odd behaviour, you should always be aware of every aggregate function behaviour or set explicitly your window frame using the RANGE clause.

## What? Show me an example!

### The basics

Given a sample database containing 5000 employes randomly assigned to 5 categories from A to E:

```
SELECT * FROM employe LIMIT 10;
```

id	name	salary	category
1	Employee name 1	1919	D
2	Employee name 2	2188	A
3	Employee name 3	5907	A
4	Employee name 4	7554	D
5	Employee name 5	8136	B
6	Employee name 6	5933	A
7	Employee name 7	2550	D
8	Employee name 8	9539	E
9	Employee name 9	3217	E
10	Employee name 10	4925	B

Now, imagine that you have to display a table with the employee information, the min, max and average salary of the employee category along with every employee row. Without window functions, the dummy way to do this kind of group-by-but-keep-every-row is the self-join, or more plausible, using php/java/ruby/whatever post-processing script.

However, grouping data and aggregating it is typically a database job, so let's tell PostgreSQL to do the heavy lifting for us:

```
1  SELECT
2    *,
3    min(salary)
4    OVER report_by_category,
5    avg(salary)
6    OVER report_by_category,
7    max(salary)
8    OVER report_by_category
9  FROM employe
10 WINDOW report_by_category AS (
11   PARTITION BY category )
12 LIMIT 5;
```

window-function-1.sql hosted with ❤ by GitHub

[view raw](#)

id	name	salary	category	min	avg	max
1003	Employee name 1003	6515	C	1009	5550.74427480916	9996
1005	Employee name 1005	8187	C	1009	5550.74427480916	9996
1001	Employee name 1001	6106	D	1001	5507.41912512716	9976



5	1001	Employee name 1001	5100	E	1001	5551.63221884498	9997
6	1002	Employee name 1002	2491	E	1001	5551.63221884498	9997
7	1004	Employee name 1004	5130	E	1001	5551.63221884498	9997

pg\_export\_1.txt hosted with ❤ by GitHub [view raw](#)

Now, each row contains informations computed with the whole partition. This is powerful as we might use any window or aggregate functions. But this comes with great responsibility. As always in SQL, it is surprisingly easy to write unmaintainable code. You should always document your window definitions and give them a meaningful name.

## Playing with the window frame

While playing with window functions and using more advanced features, you may come to a point where everything goes crazy and your results are false. If your data is clean, then it's likely to be a window frame issue:

```

1  SELECT
2    *,
3    count(*)
4    OVER report_by_category AS count,
5    count(*)
6    OVER (report_by_category
7          ORDER BY salary) AS ordered_count
8  FROM
9    employe
10 WINDOW
11    report_by_category AS (
12    PARTITION BY
13    category
14    )
15 LIMIT 5;

```

window-function-2.sql hosted with ❤ by GitHub [view raw](#)

1	id	name	salary	category	count	ordered_count
2	-----+-----+-----+-----+-----+-----					
3	2456	Employee name 2456	1008	A	1035	1
4	2209	Employee name 2209	1015	A	1035	2
5	3692	Employee name 3692	1022	A	1035	3
6	3885	Employee name 3885	1031	A	1035	4
7	436	Employee name 436	1035	A	1035	5

Here, we ordered one of the two count function partitions by a salary column which should not impact our results. However, the two functions did not behaved the same way. This is because the unordered window frame range defaults to the whole partition, but the ordered window frame range is different. It goes from the first partition row (based on the order clause) to the current partition row.

Solving this issue is as simple as explicitly setting the default window frame to the whole partition using the RANGE clause.

```
1  SELECT
2      *,
3      count(*)
4      OVER report_by_category      AS count,
5      count(*)
6      OVER ordered_report_by_category AS ordered_count
7  FROM
8      employe
9  WINDOW
10     report_by_category AS (
11         PARTITION BY
12             category
13     ),
14     ordered_report_by_category AS (
15         PARTITION BY
16             category
17         ORDER BY
18             salary
19         RANGE
20             BETWEEN UNBOUNDED PRECEDING
21             AND UNBOUNDED FOLLOWING
22     )
23  LIMIT 5;
```

window-function-3.sql hosted with ❤ by GitHub

[view raw](#)

1	id		name		salary		category		count		ordered_count
2	-----+-----+-----+-----+-----+-----										
3	2456		Employee name 2456		1008		A		1035		1035
4	2209		Employee name 2209		1015		A		1035		1035

5	3692		Employee name 3692		1022		A		1035		1035
6	3885		Employee name 3885		1031		A		1035		1035
7	436		Employee name 436		1035		A		1035		1035

window-function-4.txt hosted with ❤ by GitHub

[view raw](#)

Here we explicitly set the ordered window frame to the whole partition and finally, our simple count function behave as expected.

Use this example as a headache avoidance guideline: ordering a window without defining the frame is equivalent to developing vicious and undetectable bugs.

## Know your enemy

This example is here to show you the limitations of windowing functions. As we explained in the introduction, you are playing with a post-processing tool, not a magical silver bullet. You have to understand the data flow between the various select clauses. As a result the windowing functions performance relies directly on the select result size after the HAVING clause.

To prove this point, we rely on the Top-K query. The goal is to retrieve the first K rows of each categories of products, of each employee department, the most productive shops of each country, ...

We Here is how you retrieve the 2 employes with the highest salaries of each category:

```

1  SELECT *
2  FROM (
3      SELECT
4          *,
5          row_number()
6          OVER ordered_report_by_category
7      FROM
8          employe
9      WINDOW ordered_report_by_category AS (
10         PARTITION BY
11             category
12         ORDER BY
13             salary DESC
14         RANGE
15         BETWEEN UNBOUNDED PRECEDING

```

```

16         AND UNBOUNDED FOLLOWING
17     )
18 ) tmp
19 WHERE
20     row_number <= 2
21 ;

```

window-function-4.sql hosted with ❤ by GitHub

[view raw](#)

1	id	name	salary	category	row_number
2	-----+-----+-----+-----				
3	4023	Employee name 4023	9996	A	1
4	3795	Employee name 3795	9993	A	2
5	1050	Employee name 1050	9990	B	1
6	2664	Employee name 2664	9988	B	2
7	956	Employee name 956	9996	C	1
8	2843	Employee name 2843	9986	C	2
9	607	Employee name 607	9976	D	1
10	4983	Employee name 4983	9970	D	2
11	1455	Employee name 1455	9997	E	1
12	1528	Employee name 1528	9990	E	2

window-function-5.txt hosted with ❤ by GitHub

[view raw](#)

The setup is a little bit different here. We are working with the same table but populated with 10M rows. I also added a btree index on the category and salary columns but the query still run in about 80 sec on my laptop.

What you are saying with this query is:

- Retrieve ALL my employee data
- Compute a row number based on the salary column for each category
- Pass this result to another query
- Only keep rows with row\_number >= 2

Any kind of index won't help you here, the query semantic told the DBMS to retrieve ALL your employee data without complains and he is just doing his job.



An efficient way of computing this result is by using the UNION technique:

```
1  SELECT *
2  FROM
3  (
4      (SELECT *
5          FROM employe
6          WHERE category = 'A'
7          ORDER BY salary DESC
8          LIMIT 2)
9      UNION (SELECT *
10          FROM employe
11          WHERE category = 'B'
12          ORDER BY salary DESC
13          LIMIT 2)
14      UNION (SELECT *
15          FROM employe
16          WHERE category = 'C'
17          ORDER BY salary DESC
18          LIMIT 2)
19      UNION (SELECT *
20          FROM employe
21          WHERE category = 'D'
22          ORDER BY salary DESC
23          LIMIT 2)
24      UNION (SELECT *
25          FROM employe
26          WHERE category = 'E'
27          ORDER BY salary DESC
28          LIMIT 2)
29  ) AS tmp
30  ORDER BY category, salary DESC;
```

window-function-5.sql hosted with ❤ by GitHub

[view raw](#)

1	id	name	salary	category
2	-----+-----+-----+-----			
3	4023	Employee name 4023	9996	A
4	3795	Employee name 3795	9993	A
5	1050	Employee name 1050	9990	B
6	2664	Employee name 2664	9988	B
7	956	Employee name 956	9996	C
8	2843	Employee name 2843	9986	C

9	607		Employee name 607		9976		D
10	4983		Employee name 4983		9970		D
11	1455		Employee name 1455		9997		E
12	1528		Employee name 1528		9990		E

window-function-5.txt hosted with  by GitHub

[view raw](#)

(Look at the raw version) Here, the semantic drastically differ from the window functions query even if the final result is the same. The various select on each category are able to use the salary index to reduce ASAP the number of working rows. On the same laptop this query run in 10 ms, about 8000 times faster than the window function version.

Keeping in mind that window function are just a post-processing toolset along with the select clauses computation order helps you avoid these kind of mistakes.

## What can I use this for?

If you are more productive writing a custom aggregation using your favorite scripting language, the only use case for you is performance improvement:

- Your DBMS is typically written in a low level language and have been optimized for years.
- Your queries are written in a declarative fashion that makes plenty of room for automatic optimizations.
- If you filter on window function results, you may significantly cut the data flow size between the database and you application (TopK query), resulting in less object creation and OOP overhead.

The other use cases being:

- Building complex reports directly in SQL
- Create complex migration scripts directly in SQL, avoiding comings and goings between application code and database.
- Showing off your incredible skills and mastery of the SQL language

# Great article, now give me some data to play with!

Here is a query that creates the table from the example section:

```
1 CREATE TABLE employe AS (  
2     SELECT  
3         n                                     AS id,  
4         'Employe name' || n                 AS name,  
5         trunc(random() * 9000 + 1000)       AS salary,  
6         chr((65 + trunc(random() * 5)) :: INTEGER) AS category  
7     FROM  
8         generate_series(1, 5000, 1) AS n  
9 );
```

window-function-6.sql hosted with ❤ by GitHub

[view raw](#)

Please have a quick look of the more in-depth tutorial from [depesz.com](http://www.depesz.com) and read the PostgreSQL documentation about the detailed list of window functions and don't forget that every aggregation function is available for windowing.

## Further reading

- <http://www.depesz.com/2012/11/20/window-window-on-the-wall/>
- <http://www.postgresql.org/docs/9.1/static/queries-table-expressions.html#QUERIES-WINDOW>
- <http://www.postgresql.org/docs/9.1/static/sql-select.html#SQL-WINDOW>
- <http://www.postgresql.org/docs/9.1/static/sql-expressions.html#SYNTAX-WINDOW-FUNCTIONS>
- <http://www.postgresql.org/docs/9.1/static/functions-window.html>
- <http://www.postgresql.org/docs/9.1/static/tutorial-window.html>

. . .

*Originally published at [inovia.fr](http://inovia.fr) on June 7, 2013.*