# SQL for Data Analysis – Tutorial – ep6 (advanced stuff)

*Written by Tomi Mester on August 7, 2017*

*Last updated on January 11, 2021*

Today I'll show you some more advanced SQL moves! Okay, I'd rather put them into the *intermediate* category. They are not beginner things — but it's also for sure that junior data scientists can't live without them.

In this article, I'll show you:

- subqueries,
- SQL `CASE` and
- SQL `HAVING`

These will definitely bring your SQL game to the next level and make you more efficient at your daily job. *(Or if you don't have a job yet, they will help impress the interviewer at your* SQL screening interview*.)*

**Note: to get the most out of this article, I recommend not just reading it, but actually doing the coding part with me!**

## Before we start…

… I suggest going through these articles first – if you haven't done so yet:

1. Set up your own data server to practice: How to set up Python, SQL, R and Bash (for non-devs)

2. Install SQL Workbench to manage your SQL stuff better: How to install SQL Workbench for postgreSQL

3. SQL for Data Analysis ep1 (SQL basics)

4. SQL for Data Analysis ep2 (SQL WHERE clause)

5. SQL for Data Analysis ep3 (SQL functions and GROUP BY)

6. SQL for Data Analysis ep4 (SQL best practices)

7. SQL for Data Analysis ep5 (SQL JOIN)

## SQL Cheat Sheet

Do you want to learn faster? Join the Data36 Inner Circle and download the SQL Cheat Sheet. Just enter your email address:

## SQL subqueries (query within a query)

There are some cases in which a simple SQL query isn't enough to answer your question. Open your SQL Workbench and let's start with an easy task (we will use our `flight_delays` data set again):

**Select the average departure delay by tail numbers ( `tailnum` column) from the table – and return the minimum and maximum values of these calculated averages.**
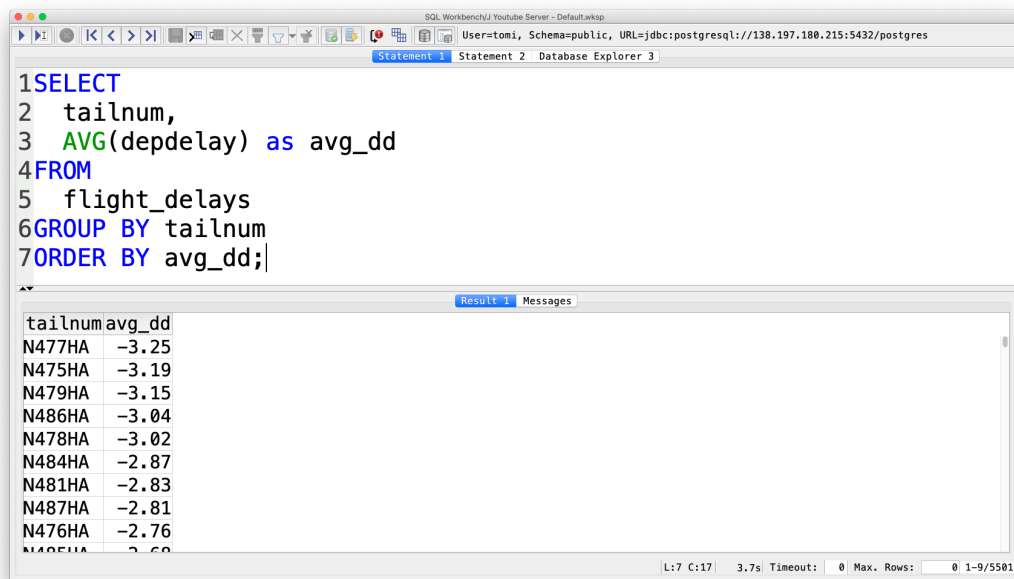
*Note: Let's try to solve it with the tools we have learned so far!*

This is the solution:

```
SELECT
  tailnum,
  AVG(depdelay) AS avg_dd
FROM
  flight_delays
GROUP BY tailnum
ORDER BY avg_dd;
```
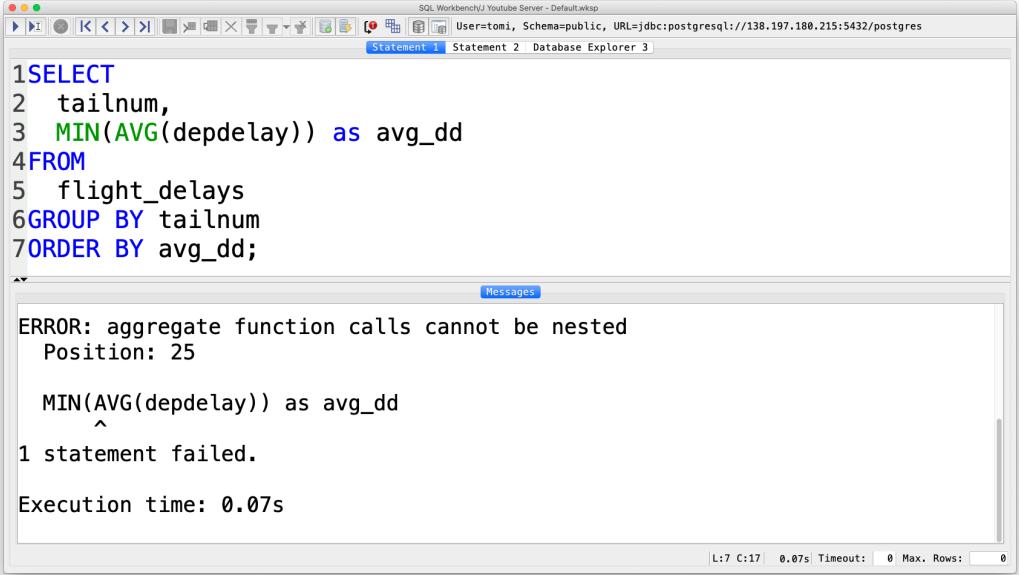
Well, at least… as far as we can get with the tools we've learned so far.

The result is:



Nice. Except that this is not the 100% correct answer for this question. Sure, you can scroll down to see the highest value in the `avg_dd` column and you can scroll back up for the lowest value. But I don't like that. It's not just inelegant, but in real life data science projects quite often it's not enough. If your automated script needs two specific values, you can't scroll through manually. You'll need to provide those two specific values. Period.

SQL-wise, the problem is that you have already run a function on your original table ( `AVG(depdelay)` ). And in SQL, you can't use *nested* functions, like `MIN(AVG(depdelay))` . You can even give it a try:

It'll just throw an error.

At this point you can't solve the problem with your original query, so you have two choices:
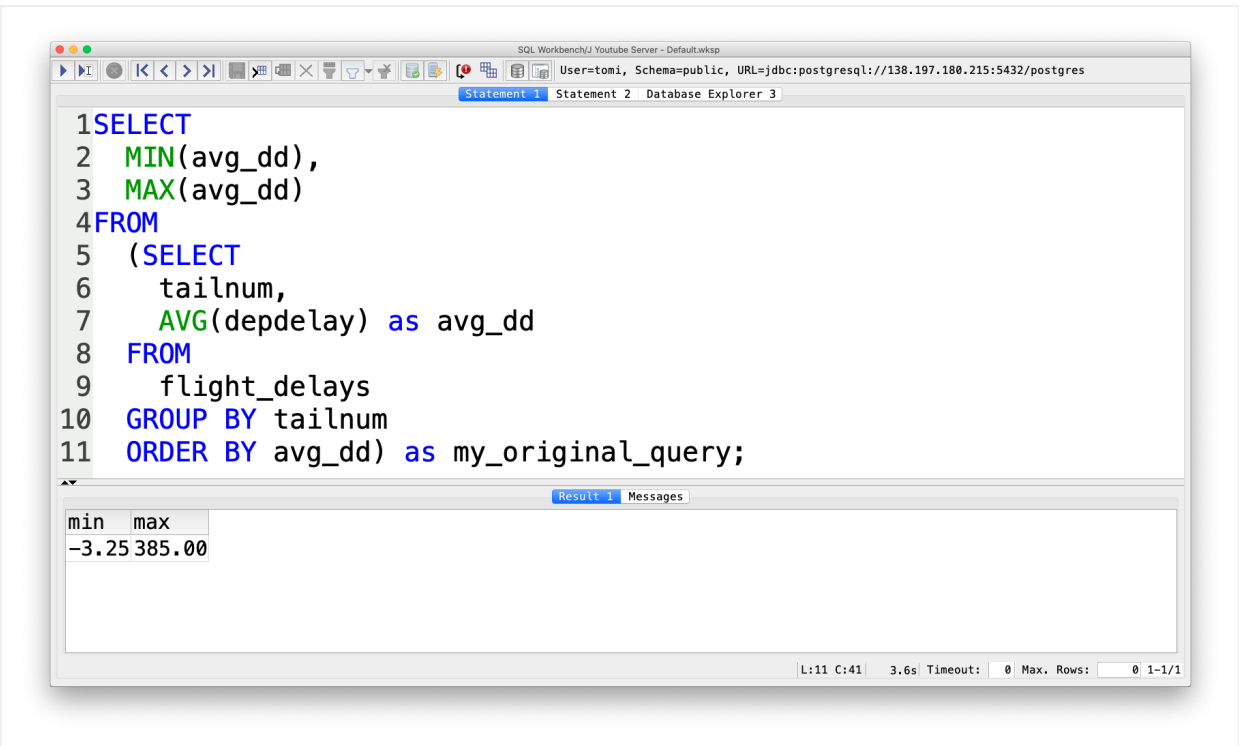
1. You can create a new table and put this result there. Then you can run your `MIN()` and `MAX()` functions on that freshly created table… This is not the preferred way, because it's a bit complicated – and it's not too flexible either.

2. A better solution is to *handle* the result of this query as if it *was* a new table – but in reality, put it as a subquery into another query.

Like this:

```
SELECT
  MIN(avg_dd),
  MAX(avg_dd)
FROM
  (SELECT
     tailnum,
     AVG(depdelay) AS avg_dd
   FROM
     flight_delays
   GROUP BY tailnum
   ORDER BY avg_dd) AS my_original_query;
```



If it looks complicated let me simplify it:

```
SELECT
  MIN(avg_dd),
  MAX(avg_dd)
FROM
  ([your_first_query]) AS [some_name];
```

Pretty cool. We have just gotten the results and you have just learned that after the `FROM` keyword you can always use subqueries instead of tables. And just to
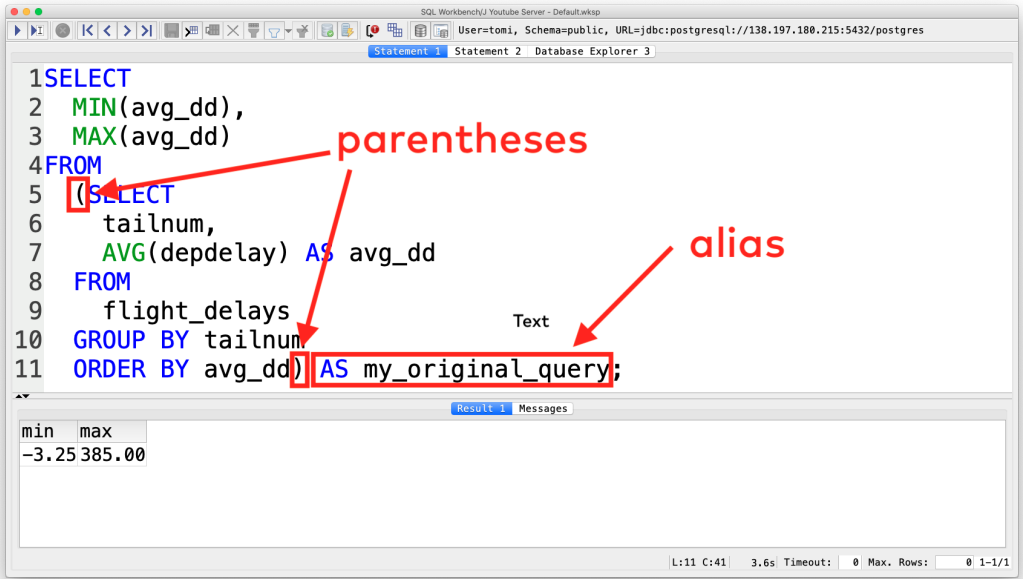
let you know, it doesn't apply only to the `FROM` keyword (e.g. you can do it with `WHERE`, too.)

*Note 1: While this solution is very handy, I have to admit that sometimes creating new tables is more efficient in terms of CPU-time. That would be a bit too advanced of a topic, so I won't go deeper into it in this article.*

## SQL subquery syntax requirements

When you use subqueries in SQL, you have to be very careful with the syntax. To be more specific, watch out for these two things:

- your inner query must go between parentheses (`(` and `)`)
- your inner query must get a new "name" with aliasing — so after the parentheses add this: `AS new_name` ... except that of course you should replace `new_name` with something more meaningful. 🙂



*SQL subquery syntax requirements*

# Test yourself #1 (with subqueries)

If you get the concept, here's another SQL exercise to test yourself with!

- **Calculate the total of the distances flown by each plane... *(Note: do the grouping based on the `tailnum` column — and `SUM` the values in the `distance` column)***
- **...then take only those planes for which this total distance is greater than 1,000,000...**
- **then calculate the average of the total distances flown by these planes!**

*Hint 1: this might be a somewhat complex query. I recommend first trying to sketch and design on paper how it will work!*

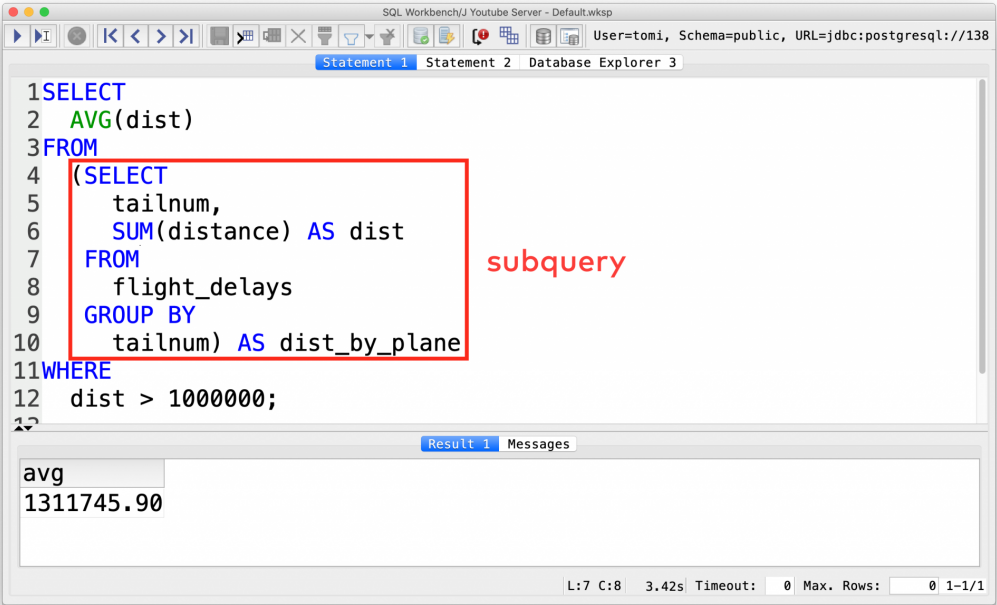*Hint 2: The result will be one number:* `1311745.90`

3... 2... 1... Go!

.

.

.

Here's my solution...

*(… but I'm pretty sure that there are other SQL queries that would deliver the same result. And you are more than welcome to post them in the comment section below the article!)*

```
SELECT
  AVG(dist)
FROM
  (SELECT
    tailnum,
    SUM(distance) AS dist
  FROM
    flight_delays
  GROUP BY
    tailnum) AS dist_by_plane
WHERE dist > 1000000;
```



*the solution (inner query highlighted)*

See? Nothing new has happened here – only that we put a query into another query. That's what SQL subqueries are about.
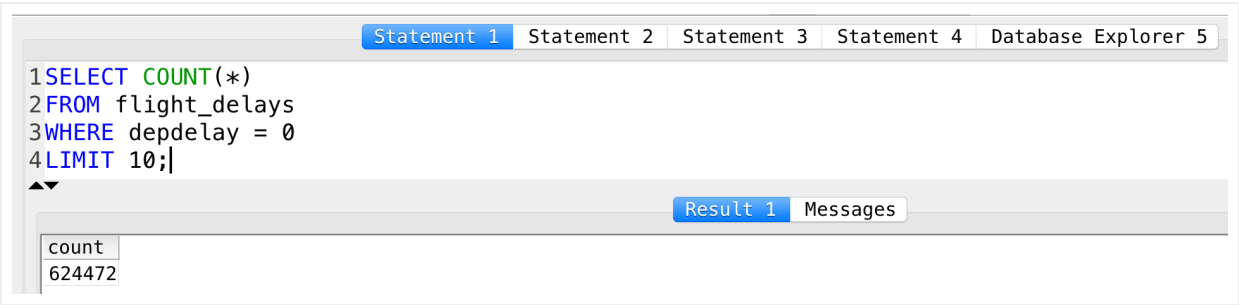
By the way, there is an alternative solution for this task using SQL `HAVING` … And I'll get back to that soon.

## SQL CASE (the "if statement" of SQL)

Let's say we want to see how many planes departed early, late or on time. With our current toolset we can do this by running 3 queries – one by one.

E.g. the number of flights that departed right on time would look like this:

```
SELECT COUNT(*)
FROM flight_delays
WHERE depdelay = 0
LIMIT 10;
```



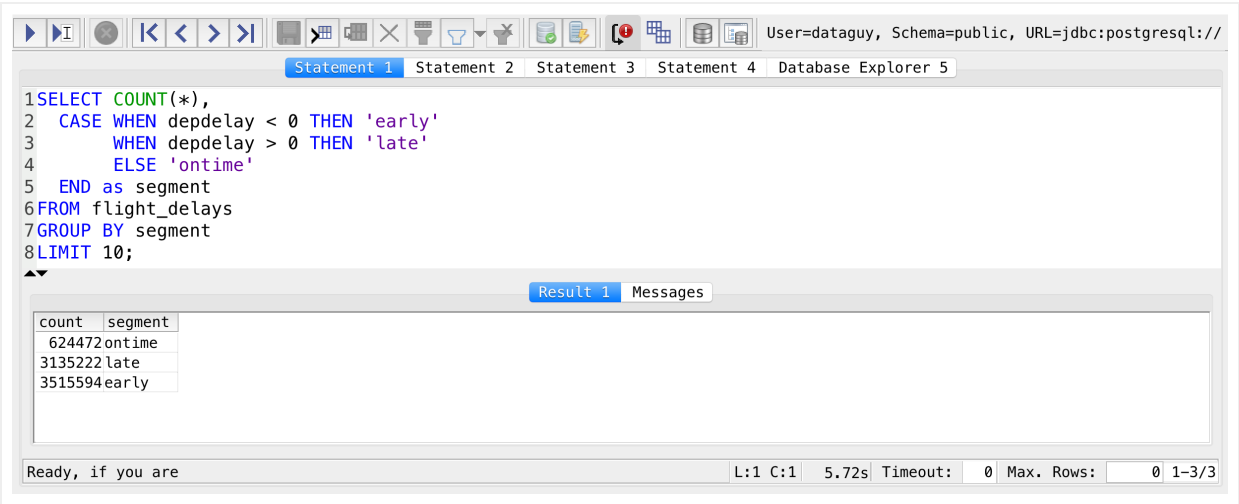To have all the variations you can change the `=` in the `WHERE` filter to `<` and to `>`.

The results will be:

- Departed early: 3515594

- **Departed on time:** `624472`

- **Departed late:** `3135222`

Good news! You can also do this using a single (slightly more advanced) SQL query if you use the `CASE` SQL keyword. Try this:

```
SELECT COUNT(*),
   CASE WHEN depdelay < 0 THEN 'early'
        WHEN depdelay > 0 THEN 'late'
        ELSE 'ontime'
   END as segment
FROM flight_delays
GROUP BY segment
LIMIT 10;
```
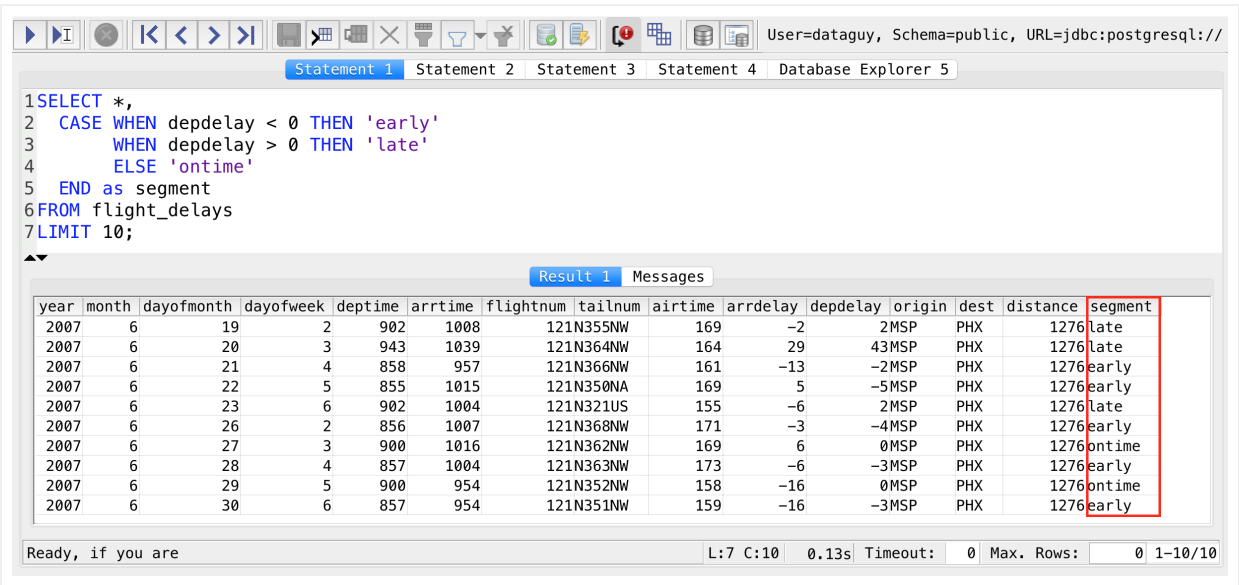


Nice, huh?

Let's take a look at the new part:

```
CASE WHEN depdelay < 0 THEN 'early'
     WHEN depdelay > 0 THEN 'late'
     ELSE 'ontime'
END as segment
```

If you are familiar with if statements in bash or in Python, it's not too hard to understand what's happening here: the `CASE WHEN` part of your query categorizes the flights, based on the `depdelay` column, into 3 categories.

*Note: if it's not clear to you at first sight, I recommend running the query without the `GROUP BY`, like this:*

```
SELECT *,
   CASE WHEN depdelay < 0 THEN 'early'
        WHEN depdelay > 0 THEN 'late'
        ELSE 'ontime'
   END as segment
FROM flight_delays
LIMIT 10;
```



*See? The trick is that the `CASE WHEN` statement is creating a new column at the end of the table.*

To be honest I don't use the SQL `CASE` statement too often during my daily job… but sometimes it's very handy when I have to do some quick and dirty data cleaning/transformation, turn a continuous value into a categorical value (like we did in the above example), and so on…

## SQL HAVING

You will need SQL HAVING for complex queries.

`HAVING` is for filtering on the results of SQL functions like `COUNT()`, `SUM()`, `AVG()`, etc.

But wait? Isn't the WHERE keyword for filtering in SQL?

Yes but here's the thing: in SQL, you can't use `WHERE` with the results of aggregate functions. Why is that? Here's a short explanation!

Let's say we want to query something simple from our table: how many times an airport shows up in the table. You should go with this query:

```
SELECT COUNT(*) as nmbr,
       dest
FROM flight_delays
GROUP BY dest;
```
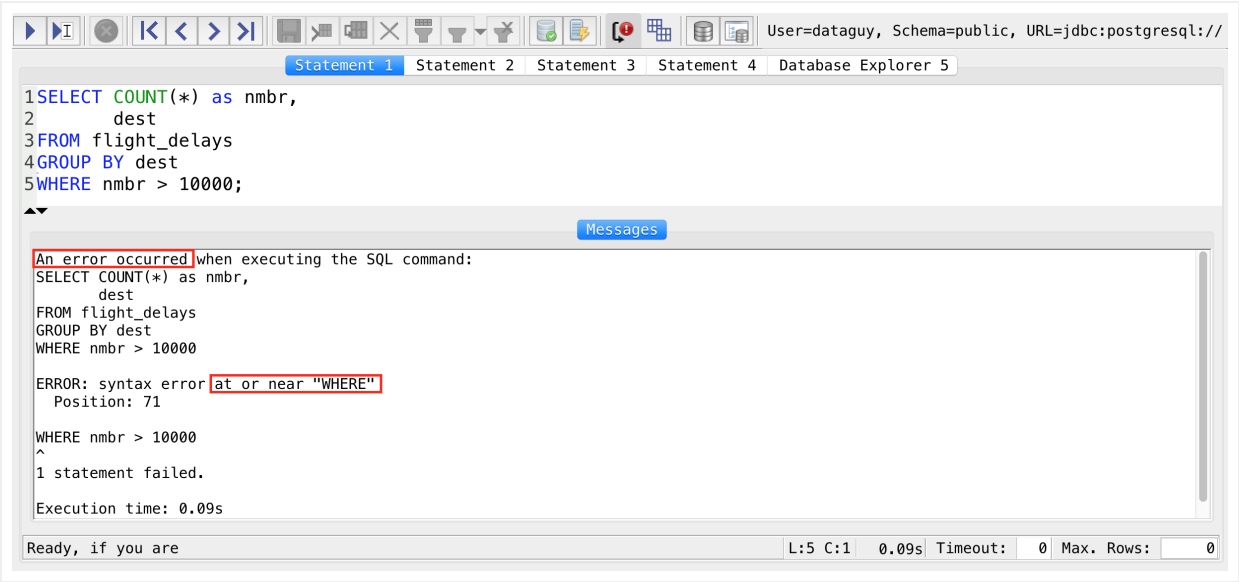
Good.

Now we want to see only those airports that show up more than 10,000 times (means that the result of `COUNT(*)` is greater than `10000`).

Someone new to SQL would try this *(and I have to admit it might seem to be highly logical at first…)*

***BUT THIS QUERY WON'T(!) WORK:***

```
SELECT COUNT(*) as nmbr,
       dest
FROM flight_delays
GROUP BY dest
WHERE nmbr > 10000;
```



Yes, it looks logical, but again: this query won't work. Why? The answer is in my previous article (SQL best practices)… But I am going to highlight it here as well. This is the order of the SQL keywords your computer sees when processing your query:
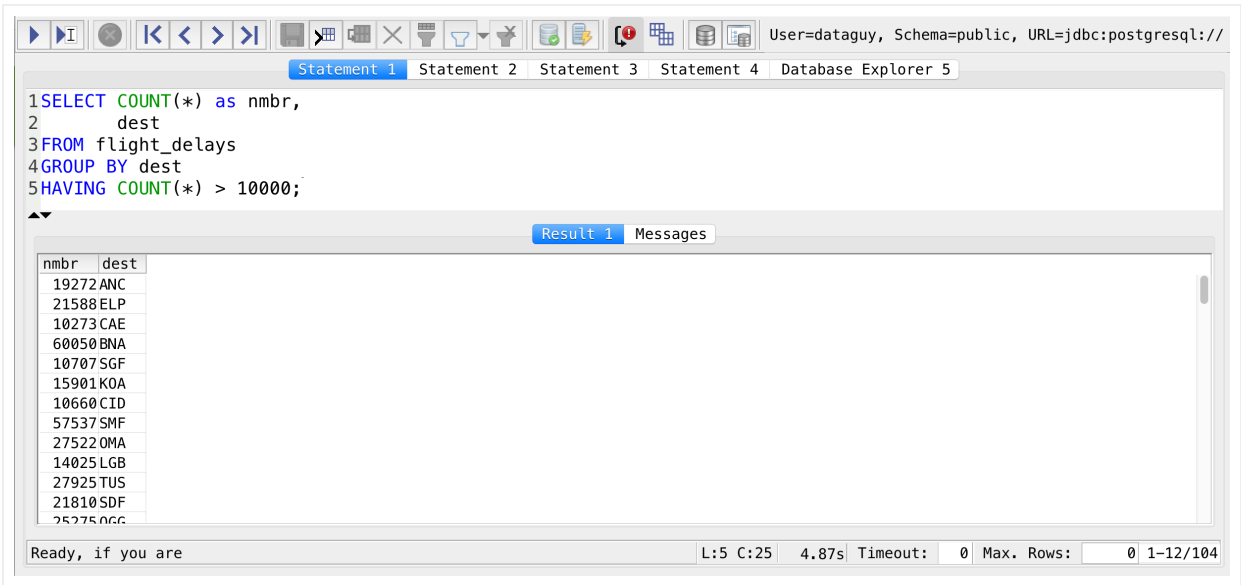
1. `FROM`
2. `ON`

3. `JOIN`

4. `WHERE`

5. `GROUP BY`

6. `SELECT`

7. `DISTINCT`

8. `ORDER BY`

9. `LIMIT`

As you can see, `WHERE` comes before `GROUP BY` and `SELECT` . This means that any aggregation column created by your `GROUP BY` clause can't be part of your `WHERE` filter. And that's where the SQL HAVING keyword comes into the picture. `HAVING` is the new `WHERE` – at least for the aggregated values in your SQL query. (Okay, this sounds very nerdy.)

You have to change your previous *wrong* query to this:

```
SELECT COUNT(*) as nmbr,
       dest
FROM flight_delays
GROUP BY dest
HAVING COUNT(*) > 10000;
```



And it'll just work!

The only unfortunate thing is that you can't use your aliases with your `HAVING` , so you need to type the whole `COUNT(*)` function again… But we can deal with that!

# Test yourself #2 (alternative solution with HAVING)

Let's try to solve the same task that we had in *TEST YOURSELF #1*… Only this time, try to take advantage of `HAVING` , as well! So:
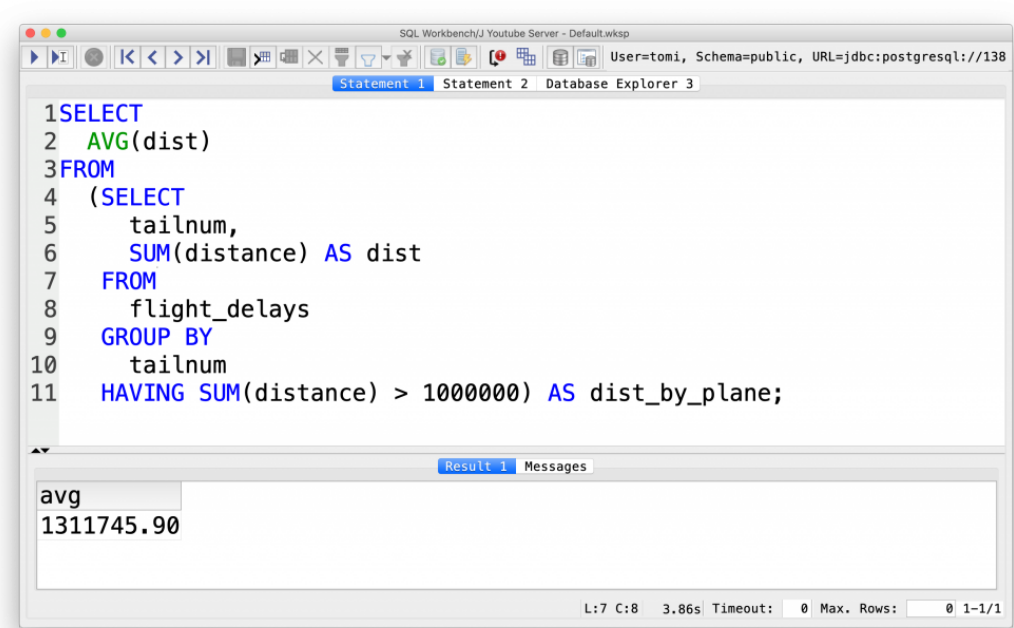
- **Calculate the total of the distances flown by each plane…** *(Note: do the grouping based on the* `tailnum` *column — and* `SUM` *the values in the* `distance` *column)*

- **…then take only those planes for which this total distance is greater than 1,000,000…**

- **then calculate the average of the total distances flown by these planes!**

*Hint 1: you need to tweak your previous query just a little bit!*

*Hint 2: The expected result is still:* `1311745.90` .

.

.

.

The solution is:

```
SELECT
  AVG(dist)
FROM
  (SELECT
     tailnum,
     SUM(distance) AS dist
   FROM
     flight_delays
   GROUP BY
     tailnum
   HAVING SUM(distance) > 1000000) AS dist_by_plane;
```



In this version, I've changed only one thing — instead of the `WHERE` in the outer query, I put a `HAVING` in the inner query. The result is the same. There might be slight differences in the computing time of the two solutions… But in this case, it won't really matter.

*Challenge: Is there a possible better solution by using* `CASE` *?* 😉 *Let me know in the comments!*

## Conclusion

Nice! These were some more advanced SQL concepts! And you have gone through them! Even better: you are done with the querying part of the SQL for data analysis tutorial series! Congrats!

If you have read and worked through all 6 episodes, you can proudly say that you know the basics of SQL. (You might also want to check out how to CREATE SQL tables.)

If you want to be better and more confident, let's do these 3 things now:

- practice,
- **practice** and
- <u>PRACTICE</u>!

I've created a 7-day SQL course for that, by the way. It doesn't just make you practice everything we learned so far but it also contains 20 job-interview-like SQL questions (and the solutions with explanations). Find more info here:

# SQL for Aspiring Data Scientists (7-day online course)

I've created an online course that will take you from zero to intermediate level with SQL in 7 days. Go ahead and check it out here:



More info...

- If you want to learn more about how to become a data scientist, take my 50-minute video course: How to Become a Data Scientist. (It's free!)

- Also check out my 6-week online course: The Junior Data Scientist's First Month video course.

*Cheers,*

**Tomi Mester**

August 7, 2017    In Coding In Data Science and Analytics

#analytics   #data coding   #learn data science   #learn to code   #small data   #sql   #sql for data analysis   #tomi mester

← PREVIOUS POST      NEXT POST →