

# SQL for Data Science

## Contents

<b>1</b>	<b>PostgreSQL For Data Science</b>	<b>1</b>
1.1	SQL Query Basics (Filtering) . . . . .	1
1.2	Using Functions . . . . .	2
1.2.1	String Functions . . . . .	3
1.2.2	Grouping Functions . . . . .	3
1.3	Grouping Data / Computing Aggregates . . . . .	4
1.4	Using Subqueries . . . . .	5
1.5	The CASE Clause . . . . .	6
1.6	Practice Problems with Solutions . . . . .	7

# 1 PostgreSQL For Data Science

## 1.1 SQL Query Basics (Filtering)

The SQL language is case insensitive, but the data within the cells is case sensitive. However, the industry standard is to capitalize the SQL keywords

```
1 select * from employees;  
2 SELECT * FROM employees; -- Industry Standard  
3 -- These two statements are the same and will produce same output
```

SQL Keywords:

- SELECT - this will specify which column/attribute you want to see in a query
- FROM - specifies which table you want to select the column/attribute from
- WHERE - a condition that lets you set parameters on a given column
- LIKE - used in addition with WHERE and '% %', lets you search for parts of a word

```
1 SELECT * FROM employees  
2 WHERE department LIKE 'F%nit%' -- such as Furniture (case sensitive)  
3 -- word must start with 'F' and contains 'nit' somewhere in middle
```

We can filter on multiple conditions by using the AND operator or the OR operator

We can also filter with these operators together by using ( )

```
1 SELECT * FROM employees  
2 WHERE department = 'Clothing' AND salary > 90000 AND region_id = 2  
3 -- all 3 conditions must be met in order to be added to query  
4  
5 SELECT * FROM employees  
6 WHERE department = 'Clothing' OR salary > 150000  
7 -- will add any row that meets either of these conditions to the query  
8  
9 SELECT * FROM employees  
10 WHERE salary < 40000 AND (department = 'Clothing' OR department = 'Pharmacy')  
11 -- will add any row where salary is less than 40k and in either department to query
```

We can use filtering operators to further process our data. Some key operators are:

- NOT - gives us value that are not in the given search ('NOT =' is same as '!=').
- IS NULL - will return all values that are null (use IS NOT NULL to get all non-null values).
- IN - will return all values that are found within the given parameters.
- BETWEEN - will search for values within a given range (inclusive).

```
1 SELECT * FROM employees  
2 WHERE NOT department = 'Clothing'; -- gives us all departments besides clothing  
3  
4 SELECT * FROM employees  
5 WHERE email IS NOT NULL; -- gives all emails that don't have null value  
6  
7 SELECT * FROM employees  
8 WHERE department IN ('Sports', 'Firs Aid', 'Garden');  
9 -- return values in any of the given departments  
10  
11 SELECT * FROM employees  
12 WHERE salary BETWEEN 80000 AND 100000; -- all salaries in this range
```

We can change the way the information is displayed in the output of our queries with certain keywords:

- ORDER BY - sorts the data by given column, either ascending (ASC) or descending (DESC)
- DISTINCT - will return all unique values for a given column
- LIMIT - will set how many records to show from our query (same as FETCH FIRST \_ ROWS ONLY)
- AS - allows us to rename a column to a given parameter name (good for exporting queries)

```
1 SELECT * FROM employees
2 ORDER BY employee_id DESC; -- orders from employee_id 1000 down to 1
3
4 SELECT DISTINCT department FROM employees -- selects unique departments in table
5 ORDER BY 1 -- orders by first parameter for select statement
6 LIMIT 10; -- only show first 10 records
7
8 SELECT first_name AS "First Name", salary AS yearly_salary -- use " " when spaces
9 FROM employees
10 FETCH FIRST 10 ROWS ONLY; -- only show first 10 records (same as LIMIT 10)
```

## 1.2 Using Functions

We can use functions with the SELECT statement to alter our data in output queries (not in the table).

- UPPER( ) - converts output to all uppercase
- LOWER( ) - converts output to all lowercase
- LENGTH( ) - gives you the length for a each string in all rows for a given column
- TRIM( ) - will take off any extra space in beginning or end (good for cleaning data)

```
1 -- all uppercase first name, all lowercase department
2 SELECT LENGTH(first_name), LOWER(department)
3 FROM employees;
4
5 -- we can test a function without selecting from a table
6 SELECT LENGTH(TRIM(' HELLO ')) -- trim extra space, return length = 5
```

We can combine multiple columns together by using concatenation || or a Boolean expression

```
1 -- combine first name and last name into new column called full_name
2 SELECT first_name || ' ' || last_name AS full_name
3 FROM employees;
4
5 -- use a boolean to create a new column of True or False
6 SELECT (salary > 140000) AS highly_paid
7 FROM employees
8 ORDER BY salary DESC; -- show all true values first
9
10 SELECT department, ('Clothing' IN department) -- creates a new column of T/F
11 FROM employees;
12
13 SELECT department, (department LIKE '%oth%') -- creates a new column of T/F
14 FROM employees; -- T if department has 'oth' in its name, otherwise F
```

### 1.2.1 String Functions

We can perform functions on strings in a given table and output them in our query (not in the table).

- SUBSTRING( ) - lets us extract part of a string (FROM start FOR length)
- REPLACE( ) - lets us change the name of strings to a new string
- POSITION( ) - lets us find a position of a given character
- COALESCE( ) - lets us add a value into any null cells for the output

```
1  -- test the substring function with a given string
2  SELECT SUBSTRING('This is test data' FROM 1 FOR 4) AS test_data; -- returns 'This'
3  -- excluding the FOR will go from position 1 to end of string
4
5  -- replace every occurrence of Clothing with Attire in a new column
6  SELECT department, REPLACE(department, 'Clothing', 'Attire') AS modified_depart
7  FROM employees;
8
9  -- return a column of integers that contain the position of @ in the email column
10 SELECT POSITION('@' IN email)
11 FROM employees;
12
13 -- use position and substring to extract email domain names into a new column
14 SELECT email, SUBSTRING(email FROM POSITION('@' IN email)+1) AS email_domain
15 FROM employees; -- for above statement, +1 ignores @ sign (increment start position)
16
17 -- create new column from email with null values filled in as NONE
18 SELECT COALESCE(email, 'NONE') AS email_filled
19 FROM employees;
```

### 1.2.2 Grouping Functions

We can perform calculations on data and get statistical insight from these queries (for numeric data).

Note that grouping functions take in multiple rows, but only output one row (the calculation).

- MAX( ) - returns the highest numeric value in a given column
- MIN( ) - returns the lowest numeric value in a given column
- AVG( ) - returns the average numeric value in a given column
- ROUND( , ) - rounds our data to a given decimal place value
- COUNT( ) - gives the total number of records in a column (excludes null values)
- SUM( ) - sums the numeric values in a given column

```
1  -- find the highest, lowest, and average paid salary in our table
2  SELECT MAX(salary) FROM employees;
3  SELECT MIN(salary) FROM employees;
4  SELECT ROUND(AVG(salary),3) FROM employees; -- average rounded to 3 decimal places
5
6  -- find total number of records in our table
7  SELECT COUNT(*) FROM employees;
8  -- note: we often use count(*) to make sure all records are counted incase of null's
9
10 -- find the yearly amount paid to employees
11 SELECT SUM(salary) FROM employees;
```

## 1.3 Grouping Data / Computing Aggregates

We can use the GROUP BY command to group records that are the same in a given column.

- Note that 'group by' comes after *where* clause but before *order by*.
- We can call the parameter for 'group/order by' with the column position in the select statement.
- Any non-aggregate columns in select list must also be mentioned in group by clause.

```
1  -- group salaries by department for region's 4, 5, 6, and 7
2  SELECT department, SUM(salary)
3  FROM employees
4  WHERE region_id IN (4,5,6,7)
5  GROUP BY department;
6
7  -- get number of employees, average/min/max salary for each department and
8  -- order from highest employee count to lowest
9  SELECT department, COUNT(*) AS num_of_employees,
10     ROUND(AVG(salary), 3) AS average_salary,
11     MIN(salary) AS lowest_salary,
12     MAX(salary) AS highest_salary
13  FROM employees
14  GROUP BY 1 -- department
15  ORDER BY 2 DESC; -- num_of_employees
16
17  -- get the number of employees by gender for each department
18  SELECT department, gender, COUNT(*) AS num_of_employees
19  FROM employees
20  GROUP BY 1, 2 -- gender is non-aggregate function (must be included)
21  ORDER BY 1; -- order by department (easier to read)
22
23  -- how many people have the same first name (name and count)
24  SELECT first_name, COUNT(*) AS occurrences
25  FROM employees
26  GROUP BY 1
27  HAVING COUNT(*) > 1;
28
29  -- get unique domain names for email and the number of employees having that domain
30  SELECT SUBSTRING(email FROM POSITION('@' IN email)+1) AS email_domain,
31     COUNT(*) AS num_of_employees
32  FROM employees
33  WHERE email IS NOT NULL
34  GROUP BY 1
35  ORDER BY 2 DESC;
36
37  -- get the min/max/avg salary for each gender in every region
38  SELECT gender, region_id, MIN(salary) AS min_salary,
39     MAX(salary) AS max_salary,
40     ROUND(AVG(salary)) AS avg_salary -- round nearest whole num
41  FROM employees
42  GROUP BY 1,2
43  ORDER BY 1,2 ASC;
```

We can use the HAVING command to filter data that has been grouped (similar to where clause).

- This command is used to filter *aggregated* data.
- Note that the having command comes after the *group by* and before the *order by* clauses

```
1  -- get all department names that have more than 35 employees
2  SELECT department, COUNT(*)
3  FROM employees
4  GROUP BY 1
5  HAVING COUNT(*) > 35
6  ORDER BY 1;
```

## 1.4 Using Subqueries

We can refer to columns of specific sources by specifying the table name before the column in the SELECT statement. We can make this even more simple by giving the sources aliases as their reference name.

- Note for future - we can have select statements nested in from statements (new source of data).

```
1 -- alias employees table as 'e' and departments as 'd'
2 SELECT d.department -- get department column from departments table
3 FROM employees e, departments d;
```

We can use subqueries in the WHERE/FROM clause (acts as a source from which data can be pulled from). In the from clause, we need to give the subquery an alias.

- The inner query (subquery) returns a list of data that we can pull from.
- Renaming a subqueries column names means we must reference these names in the outer query.
- We can have multiple sources of data (subqueries) in our FROM clause.

```
1 -- select all employees who work in a department not listed in departments table
2 SELECT * FROM employees
3 WHERE department NOT IN (SELECT department from departments);
4
5 -- select first name and salary from subquery 'a' of employees with salary > 150k
6 SELECT a.first_name, a.salary
7 FROM (SELECT * FROM employees WHERE salary > 150000) a
8
9 -- renaming inner query names (same as above problem)
10 SELECT a.employee_name, a.yearly_salary
11 FROM (SELECT first_name AS employee_name, salary AS yearly_salary
12       FROM employees WHERE salary > 150000) a
13
14 -- select all employees who work in the electronics division
15 SELECT * FROM employees
16 WHERE department IN (SELECT department
17 FROM departments WHERE division = 'Electronics');
18
19 -- select all employees that work in Asia or Canada and make over 130k
20 SELECT * FROM employees
21 WHERE region_id IN (SELECT region_id FROM regions WHERE country IN ('Asia','Canada'))
22 AND salary > 130000;
```

We can use subqueries in the SELECT statement, but we must make sure that the subquery only returns one record. This can be used to compare all records of outer query to one value from inner query (see example below).

```
1 -- INVALID SYNTAX EXAMPLE
2 SELECT first_name, salary, (SELECT first_name FROM employees)
3 FROM employees;
4 /* this will not run, the inner query will try to return all 1000 first_names
5    for each record in outer query. */
6
7 /* select first_name and department of employee and how much less they make than
8    the highest paid employee (in Asia and Canada) */
9 SELECT first_name, department,
10        ((SELECT MAX(salary) from employees) - salary) AS less_income
11 FROM employees
12 WHERE region_id IN (SELECT region_id FROM regions
13                    WHERE country IN ('Asia','Canada'));
```

We can use the ALL/ANY clause in the where/having clause, and they have the following function:

- ANY - will return true if any of the subquery values meet the condition.
- ALL - will return true if all subquery values meet the condition.

You can use all operator >, <, >=, <= when using these clauses.

Note that ANY will still return values inside the subquery (tricky to use).

```
1  -- select all employees who's region_id is greater than US region id's (1,2,3)
2  SELECT * FROM employees
3  WHERE region_id > ALL (SELECT region_id FROM regions WHERE country='United States');
4
5  /* select all employees that work in the kids division and the dates
6  at which they were hired is greater than all hire dates of employees
7  in the maintenance department */
8  SELECT * FROM employees
9  WHERE department IN (SELECT department FROM departments WHERE division = 'Kids')
10 AND hire_date > ALL (SELECT hire_date FROM employees
11                      WHERE department = 'Maintenance');
12
13 -- select the salary that occurs the most often (and is the highest value)
14 SELECT salary, COUNT(*) FROM employees
15 GROUP BY 1
16 ORDER BY 2 DESC, 1 DESC
17 LIMIT 1;
18 -- same as...
19 SELECT salary FROM employees
20 GROUP BY salary
21 HAVING COUNT(*) >= ALL(SELECT COUNT(*) FROM employees GROUP BY salary)
22 ORDER BY 1 DESC
23 LIMIT 1;
24
25 -- find the average salary excluding the lowest and highest paid employee
26 SELECT ROUND(AVG(salary)) FROM employees
27 WHERE salary < ALL (SELECT MAX(salary) FROM employees)
28 AND salary > ALL (SELECT MIN(salary) FROM employees);
```

## 1.5 The CASE Clause

## 1.6 Practice Problems with Solutions

### (1.1 SQL Query Basics)

```
1  /* First name and email of females that work in the tools department having a
2     salary greater than 110,000 */
3  SELECT first_name, email FROM employees
4  WHERE gender = 'F' AND department = 'Tools' AND salary > 110000;
5
6  /* First name and hire date of employees who earn more than 165,000 as well as
7     employees that work in the sports department and are men */
8  SELECT first_name, hire_date FROM employees
9  WHERE salary > 165000 OR (department = 'Sports' and gender = 'M');
10
11 /* First name and hire date of employees hired during Jan 1, 2002 and Jan 1, 2004 */
12 SELECT first_name, hire_date FROM employees
13 WHERE hire_date BETWEEN '2002-01-01' AND '2004-01-01';
14
15 /* All columns from male employees who work in the automotive department and earn
16     more than 40,000 and less than 100,000 as well as females that work in the
17     toy department */
18 SELECT * FROM employees
19 WHERE gender = 'M' AND department = 'Automotive'
20 AND (salary BETWEEN 40000 and 100000)
21 OR (gender = 'F' AND department = 'Toys');
```

### (1.2 Using Functions)

```
1  /* Write a query against the professors table that can output the following in the
2     result: "Chong works in the Science department" */
3  SELECT last_name || ' works in the ' || department || ' department'
4  FROM professors
5  WHERE last_name = 'Chong';
6
7  /* Write a query that says if a professor is highly paid (above 95000)
8     in the format "It is false that professor Chong is highly paid" */
9  SELECT 'It is ' || (salary > 95000) || ' that professor ' || last_name ||
10     ' is highly paid'
11 FROM professors;
12
13 /* Write a query that returns all of the records and columns from the professors
14     table but shortens the department names to only the first three characters in
15     upper case. */
16 SELECT last_name, UPPER(SUBSTRING(department FROM 1 FOR 3)) AS department_abrv,
17     salary, hire_date
18 FROM professors;
19
20 /* Write a query that returns the highest and lowest salary from the professors
21     table excluding the professor named 'Wilson'. */
22 SELECT MAX(salary) AS max_salary, MIN(salary) AS min_salary
23 FROM professors
24 WHERE last_name != 'Wilson';
25
26 /* Write a query that will display the hire date of the professor that has been
27     teaching the longest. */
28 SELECT MIN(hire_date) FROM professors;
```



### (1.3 Grouping Data / Computing Aggregates)

```
1  -- Write a query that displays only the state with the largest amount of fruit supply
2  SELECT state
3  FROM fruit_imports
4  GROUP BY 1
5  ORDER BY SUM(supply) DESC
6  LIMIT 1;
7
8  -- Write a query that returns the most expensive cost_per_unit of every season.
9  SELECT season, MAX(cost_per_unit)
10 FROM fruit_imports
11 GROUP BY 1;
12
13 -- Write a query that returns the state that has more than 1 import of the same fruit
14 SELECT state
15 FROM fruit_imports
16 GROUP BY 1, name
17 HAVING COUNT(name) > 1;
18
19 -- Write a query that returns the seasons that produce either 3 fruits or 4 fruits.
20 SELECT season
21 FROM fruit_imports
22 GROUP BY 1
23 HAVING ((COUNT(name) = 3) OR (COUNT(name) = 4));
24
25 /* Write a query that takes into consideration the supply and cost_per_unit columns
26    for determining the total cost and returns the most expensive state with the total
27    cost. */
28 SELECT state, SUM(supply * cost_per_unit) AS total_cost
29 FROM fruit_imports
30 GROUP BY 1
31 ORDER BY 2 DESC
32 LIMIT 1;
33
34 -- Execute the below SQL script and write a query that returns the count of 4.
35 CREATE table fruits (fruit_name varchar(10));
36 INSERT INTO fruits VALUES ('Orange');
37 INSERT INTO fruits VALUES ('Apple');
38 INSERT INTO fruits VALUES (NULL);
39 INSERT INTO fruits VALUES (NULL);
40
41 SELECT COUNT(COALESCE(fruit_name, 'SOMEVALUE'))
42 FROM fruits;
```

### (1.4 Using Subqueries)