

Data Science

Contents

1	PostgreSQL For Data Science	1
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	Correlated Subqueries	7
1.7	Working with Multiple Tables	8
1.7.1	Table Unions	8
1.7.2	Table Joins	8
1.7.3	Creating Views	9
1.8	Window Functions	9
1.9	Practice Problems with Solutions	11

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

We can use the CASE clause as a conditional expression to create a new column in the SELECT clause.

We can also use CASE statements to transpose tables (rows to columns, columns to rows).

```
1  -- select the total number for each category of the pay scale
2  SELECT a.pay_category, COUNT(*) FROM (
3      SELECT first_name, salary, -- subquery case statement
4          CASE
5              WHEN salary < 100000 THEN 'UNDER PAID'
6              WHEN salary > 100000 AND salary < 160000 THEN 'PAID WELL'
7              WHEN salary > 160000 THEN 'EXECUTIVE'
8              ELSE 'UNPAID'
9          END AS pay_category -- m name
10     FROM employees
11     ORDER BY salary DESC ) a -- access variable
12 GROUP BY 1;
13
14 -- transpose the above table using case statements (rows will have employee count)
15 SELECT SUM( CASE WHEN salary < 100000 THEN 1 ELSE 0 END) as under_paid,
16        SUM( CASE WHEN salary > 100000 AND salary < 160000 THEN 1 ELSE 0 END) as paid_well,
17        SUM( CASE WHEN salary > 160000 THEN 1 ELSE 0 END) as executive
18 FROM employees;
```

```

1  -- find total number of employees for the given departments
2  SELECT department, count(*)
3  FROM employees
4  WHERE department IN ('Sports', 'Tools', 'Clothing', 'Computers')
5  GROUP BY 1;
6
7  -- transpose the above table using case statements (rows will have employee count)
8  SELECT SUM( CASE WHEN department = 'Sports' THEN 1 ELSE 0 END) as sports_employees,
9         SUM( CASE WHEN department = 'Tools' THEN 1 ELSE 0 END) as tools_employees,
10        SUM( CASE WHEN department = 'Clothing' THEN 1 ELSE 0 END) as clothing_employees,
11        SUM( CASE WHEN department = 'Computers' THEN 1 ELSE 0 END) as computers_employees
12  FROM employees;
13
14  -- count how many employees are in each country
15  SELECT COUNT(a.region_1) + COUNT(a.region_2) + COUNT(a.region_3) AS united_states,
16         COUNT(a.region_4) + COUNT(a.region_5) AS asia,
17         COUNT(a.region_6) + COUNT(a.region_7) AS canada
18  FROM (SELECT first_name,
19         CASE WHEN region_id = 1 THEN (SELECT country FROM regions WHERE region_id = 1)
20         END AS region_1,
21         CASE WHEN region_id = 2 THEN (SELECT country FROM regions WHERE region_id = 2)
22         END AS region_2,
23         CASE WHEN region_id = 3 THEN (SELECT country FROM regions WHERE region_id = 3)
24         END AS region_3,
25         CASE WHEN region_id = 4 THEN (SELECT country FROM regions WHERE region_id = 4)
26         END AS region_4,
27         CASE WHEN region_id = 5 THEN (SELECT country FROM regions WHERE region_id = 5)
28         END AS region_5,
29         CASE WHEN region_id = 6 THEN (SELECT country FROM regions WHERE region_id = 6)
30         END AS region_6,
31         CASE WHEN region_id = 7 THEN (SELECT country FROM regions WHERE region_id = 7)
32         END AS region_7
33  FROM employees) a; /* subquery 'a' has a column for the name of each employee and a
34  column for each region, if an employee works in the region it has the country name,
35  otherwise it has a null value */

```

1.6 Correlated Subqueries

A correlated subquery means that the subquery portion is correlated to the outer query (in other words, the nested subquery uses values from the outer query). This also means that the nested query will run for every single record of the outer query since we are forming a link between the two. Note that as the number of records grow, this method becomes more and more time consuming.

```

1  -- get all employees who make over the average salary in their department
2  SELECT first_name, salary
3  FROM employees e1
4  WHERE salary > (SELECT ROUND(AVG(salary)) -- subquery for average of each department
5                  FROM employees e2 WHERE e1.department = e2.department);
6
7  -- get first name, department, salary, and average department salary
8  SELECT first_name, department, salary,
9         (SELECT ROUND(AVG(salary)) FROM employees e2 WHERE e1.department = e2.department)
10  FROM employees e1;
11
12  -- get all departments that have more than 38 employees and their max salary
13  SELECT department, (SELECT MAX(salary) FROM employees WHERE department= d.department)
14  FROM departments d
15  WHERE 38 < (SELECT COUNT(*) FROM employees e2 WHERE e2.department = d.department);

```


1.7 Working with Multiple Tables

1.7.1 Table Unions

We can use the UNION clause to stack data on top of each other (note that this will remove any duplicates and only show unique values). However, we can use the UNION ALL clause to stack data on top of each other without removing any duplicates (show all records from both data sources). Also note that the columns you are using UNION on must match for both tables.

```
1  -- select the number of employees for each department and a total count as the bottom
2  SELECT department, COUNT(*)
3  FROM employees
4  GROUP BY 1
5  UNION ALL
6  SELECT 'TOTAL', COUNT(*)
7  FROM employees;
```

We can use the EXCEPT clause to take the first result set and removes from it all the rows found in the second result set (records that exist in first data source but not second data source).

```
1  -- select all department from employees not in departments table
2  SELECT DISTINCT department
3  FROM employees
4  EXCEPT
5  SELECT department
6  FROM departments;
```

1.7.2 Table Joins

Joining is a way to link two or more tables together through a common column. We can do a basic join through the WHERE clause, and specifying which matching columns we want to join on. Note that calling a column that is ambiguous (same in multiple tables) must be specified which table you want to pull it from (by table name or alias). We can join subqueries since they are a source of data.

```
1  /* get first name, email (non-null values), department, division, and country for
2  all employees */
3  SELECT first_name, email, e.department division, country -- specify department from e
4  FROM employees e, departments d, regions r
5  WHERE (e.department = d.department) -- join e with d
6        AND (e.region_id = r.region_id) -- join (e/d) with r
7  AND email IS NOT NULL;
8
9  -- get the number of employees in each country (make regions a subquery)
10 SELECT country, COUNT(employee_id) AS num_of_employees
11 FROM (SELECT * from regions) r, employees e
12 WHERE e.region_id = r.region_id
13 GROUP BY 1;
```

We can use the INNER JOIN... ON clauses to combine two or more tables on a specific column with *matching* data. Note that for all data that does not match, it will be left out of the combined table.

```
1  SELECT first_name, email, division, country
2  FROM employees INNER JOIN departments
3        ON employees.department = departments.department
4  INNER JOIN regions ON employees.region_d = regions.region_id
5  WHERE email IS NOT NULL;
```

We can use the LEFT JOIN... ON clauses to combine two or more tables, and keep all data from the first table even if it does not match any data in the second table. Similarly, we can use RIGHT JOIN... ON clauses to keep all data from the second table even if it's not in the first table. Note that for both, any missing values will be filled in with *null*.

```

1  -- get all department records that are in employees but not in departments
2  SELECT DISTINCT employees.department
3  FROM employees LEFT JOIN departments
4  ON employees.department = departments.department
5  WHERE departments.department IS NULL;

```

The CROSS JOIN clause lets us find the Cartesian product of two data sources (multiply each data record in the first source with every data record in the second source, all possible combinations).

```

1  -- employees table (1000 records) x departments table (24 records) = 24000 records
2  SELECT * FROM employees CROSS JOIN departments;

```

Join, Union, and Subquery combined example:

```

1  -- select the first and last employee hired along with department and country
2  (SELECT first_name, department, hire_date, country
3  FROM employees INNER JOIN regions
4  ON employees.region_id = regions.region_id
5  WHERE hire_date = (SELECT MIN(hire_date) FROM employees e2)
6  LIMIT 1) -- closing in parentheses lets us use the LIMIT clause
7  UNION ALL
8  SELECT first_name, department, hire_date, country
9  FROM employees INNER JOIN regions
10 ON employees.region_id = regions.region_id
11 WHERE hire_date = (SELECT MAX(hire_date) FROM employees e2)

```

1.7.3 Creating Views

Views are created based on queries (can't insert or delete data from a view) and are saved to the database we are working within. We can access columns from the view just like other data sources. The standard to name these views are to put a v_ before the name so that anyone accessing knows it is a view. Also, every time you access a view, it is running the query that the view is made of.

```

1  -- save view to be accessed later
2  CREATE VIEW v_employee_information AS
3  SELECT first_name, email, e.department, salary, division, region, country
4  FROM employees e, departments d, regions r
5  WHERE e.department = d.department
6  AND e.region_id = r.region_id;
7
8  SELECT first_name, region FROM v_employee_information;

```

1.8 Window Functions

The OVER() keyword allows us to detach a given aggregate function from a subquery and group it on a particular window by passing PARTITION BY inside the parenthesis. Note that leaving parenthesis blank will cause the function to use the entire data source as the window.

```

1  -- select the first name, department, and how many employees in the given department
2  SELECT first_name, department, COUNT(*) OVER(PARTITION BY department)
3  FROM employees
4
5  /* select name, department, number of employees in the current department, and
6  the number of employees in the current region */
7  SELECT first_name, department, region_id,
8  COUNT(*) OVER(PARTITION BY department) AS dept_count,
9  COUNT(*) OVER(PARTITION BY region_id) AS region_count
10 FROM employees;

```

We can use the ORDER BY and RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW (*this is optional since it is default*) clause, meaning sum all values from first record to the current row, in the OVER() function to compute data in a given order.

```

1  -- select name, hire date, salary, and the running total of salaries
2  SELECT first_name, hire_date, salary,
3  SUM(salary) OVER(ORDER BY hire_date RANGE BETWEEN UNBOUNDED PRECEDING AND
4                  CURRENT ROW) AS running_total_of_salaries
5  FROM employees;
6
7  -- find the running total of salaries for each department
8  SELECT first_name, hire_date, department, salary,
9  SUM(salary) OVER(PARTITION BY department ORDER BY hire_date) AS salary_running_total
10 FROM employees;

```

We can pass ROWS BETWEEN _ PRECEDING AND CURRENT ROW to add the value from the _ previous row(s) and current row together (compute calculations on adjacent rows).

```

1  SELECT first_name, hire_date, department, salary,
2  SUM(salary) OVER(ORDER BY hire_date ROWS BETWEEN 1 PRECEDING AND
3                  CURRENT ROW) AS salary_running_total
4  FROM employees;

```

We can use the RANK() clause to rank a given partition by on the parameters within the OVER() clause.

```

1  -- get all employees who are in the 8th rank for their given department
2  SELECT * FROM (
3      SELECT first_name, email, department, salary,
4      RANK() OVER(PARTITION BY department ORDER BY salary DESC)
5      FROM employees ) a
6  WHERE rank = 8;

```

We can use the NTILE() function, where we pass a parameter of how many groups to split the data into, to rank groups of rows (think of as creating brackets).

```

1  -- split each department into 5 brackets based on their salary
2  SELECT first_name, email, department, salary,
3  NTILE(5) OVER(PARTITION BY department ORDER BY salary DESC) AS salary_bracket
4  FROM employees;

```

We can use the FIRST_VALUE() function to get the first value of a specified column, which is passed as a parameter, for each group that we are partitioning by. We can also use NTH_VALUE(,) by passing the column name and the value position we wish to get.

```

1  -- get the first salary for each department (highest paid)
2  SELECT first_name, email, department, salary,
3  FIRST_VALUE(salary) OVER(PARTITION BY department ORDER BY salary DESC)
4  FROM employees;

```

We can use the LAG() function to get the previous value of a column that is passed as a parameter. Similarly, we can use the LEAD() function to get the next value of a column that is passed as a parameter.

```

1  -- get the salary of the employee that is one below the current employee
2  SELECT first_name, department, salary,
3  LEAD(salary) OVER(PARTITION BY department ORDER BY salary DESC) AS lower_salary
4  FROM employees;

```

We can group by multiple values to see each result in one query by using the GROUPING SETS() clause and passing the columns we wish to group by (it will fill in null for all missing records in the column not being grouped by). Note that passing () to the grouping sets clause will execute the aggregate function on all of the data.

```

1  /* group by continent, country, and city to see total units sold for each and also
2     the total units sold in the entire table */
3  SELECT continent, country, city, sum(units_sold)
4  FROM SALES
5  GROUP BY GROUPING SETS(continent, country, city, ());

```

Similar to grouping sets, we can use the ROLLUP() clause. Lets assume we want to group by 3 columns, then rollup will show us the results if we grouped by nothing (only shows the aggregate function value), all columns, the first 2 columns, and the first column all in one query. We can also use the CUBE() clause to show all possible combinations of groupings for any columns that are passed as parameters.

```

1  SELECT continent, country, city, sum(units_sold)
2  FROM SALES
3  GROUP BY ROLLUP(continent, country, city);

```

1.9 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
14
15
16

```

```

17  /* Write a query that returns all of the records and columns from the professors
18     table but shortens the department names to only the first three characters in
19     upper case. */
20  SELECT last_name, UPPER(SUBSTRING(department FROM 1 FOR 3)) AS department_abrv,
21         salary, hire_date
22  FROM professors;
23
24  /* Write a query that returns the highest and lowest salary from the professors
25     table excluding the professor named 'Wilson'. */
26  SELECT MAX(salary) AS max_salary, MIN(salary) AS min_salary
27  FROM professors
28  WHERE last_name != 'Wilson';
29
30  /* Write a query that will display the hire date of the professor that has been
31     teaching the longest. */
32  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 state that has more than 1 import of the same fruit
9  SELECT state
10 FROM fruit_imports
11 GROUP BY 1, name
12 HAVING COUNT(name) > 1;
13
14 -- Write a query that returns the seasons that produce either 3 fruits or 4 fruits.
15 SELECT season
16 FROM fruit_imports
17 GROUP BY 1
18 HAVING ((COUNT(name) = 3) OR (COUNT(name) = 4));
19
20 /* Write a query that takes into consideration the supply and cost_per_unit columns
21    for determining the total cost and returns the most expensive state with the total
22    cost. */
23 SELECT state, SUM(supply * cost_per_unit) AS total_cost
24 FROM fruit_imports
25 GROUP BY 1
26 ORDER BY 2 DESC
27 LIMIT 1;
28
29 -- Execute the below SQL script and write a query that returns the count of 4.
30 CREATE table fruits (fruit_name varchar(10));
31 INSERT INTO fruits VALUES ('Orange');
32 INSERT INTO fruits VALUES ('Apple');
33 INSERT INTO fruits VALUES (NULL);
34 INSERT INTO fruits VALUES (NULL);
35
36 SELECT COUNT(COALESCE(fruit_name, 'SOMEVALUE'))
37 FROM fruits;

```

(1.4 Using Subqueries)

```
1  /* Write a query that returns the names of those students that are taking
2     Physics and US History. */
3  SELECT student_name FROM students
4  WHERE student_no IN (SELECT student_no FROM student_enrollment
5                       WHERE course_no IN (SELECT course_no FROM courses
6                                           WHERE course_title
7                                           IN ('Physics', 'US History')));
8
9  -- Write a query that returns the name of the student that is taking the most courses
10 SELECT student_name FROM students
11 WHERE student_no = (SELECT student_no FROM student_enrollment
12                    GROUP BY 1
13                    ORDER BY COUNT(student_no) DESC
14                    LIMIT 1);
15
16 -- Write a query to find the student that is the oldest (no LIMIT or ORDER BY)
17 SELECT student_name FROM students
18 WHERE age = (SELECT MAX(age) FROM students);
```

(1.5 The CASE Clause)

```
1  /* Write a query that displays 3 columns. The query should display the fruit and
2     it's total supply along with a category of either LOW, ENOUGH or FULL */
3  SELECT name, total_supply,
4  CASE
5  WHEN total_supply < 20000 THEN 'LOW'
6  WHEN total_supply > 20000 AND total_supply < 50000 THEN 'ENOUGH'
7  WHEN total_supply > 50000 THEN 'FULL'
8  END AS supply_category
9  FROM (SELECT name, sum(supply) AS total_supply
10        FROM fruit_imports
11        GROUP BY name) a;
12
13 /* Taking into consideration the supply column and the cost_per_unit column, you
14 should be able to tabulate the total cost to import fruits by each season. Write
15 a query that would transpose this data so that the seasons become columns and
16 the total cost for each season fills the first row */
17 SELECT SUM( CASE WHEN season = 'Winter' THEN total_cost ELSE 0 END) AS "Winter",
18        SUM( CASE WHEN season = 'Summer' THEN total_cost ELSE 0 END) AS "Summer",
19        SUM( CASE WHEN season = 'All Year' THEN total_cost ELSE 0 END) AS "All Year",
20        SUM( CASE WHEN season = 'Spring' THEN total_cost ELSE 0 END) AS "Spring",
21        SUM( CASE WHEN season = 'Fall' THEN total_cost ELSE 0 END) AS "Fall"
22 FROM (SELECT season, SUM(supply * cost_per_unit) AS total_cost
23        FROM fruit_imports
24        GROUP BY season) a;
```

(1.6 Correlated Subqueries)

```
1  /* Write a query that gets the min and max salary for each department and says which
2     is the HIGHEST and LOWEST salary (should have 2 records for each department in
3     final query */
4  SELECT department, first_name, salary,
5  CASE WHEN salary = d_max THEN 'HIGHEST SALARY'
6       WHEN salary = d_min THEN 'LOWEST SALARY'
7  END AS salary_in_department
8  FROM (SELECT department, first_name, salary,
9         (SELECT MAX(salary) FROM employees WHERE department = e1.department) AS d_max,
10        (SELECT MIN(salary) FROM employees WHERE department = e1.department) AS d_min
11        FROM employees e1
12        ORDER BY 1) a
13  WHERE salary = dep_max OR salary = dep_min
14  ORDER BY 1;
```

(1.7 Multiple Tables)

```
1  /* Write a query that shows the student's name, the courses the student is taking
2     and the professors that teach that course.*/
3  SELECT student_name, se.course_no, p.last_name
4  FROM students s
5  INNER JOIN student_enrollment se
6  ON s.student_no = se.student_no
7  INNER JOIN teach t
8  ON se.course_no = t.course_no
9  INNER JOIN professors p
10 ON t.last_name = p.last_name
11 ORDER BY student_name;
12
13 /* In the above problem, you discovered why there is repeating data. How can we
14    eliminate this redundancy? Make every record distinct. */
15 SELECT student_name, course_no, min(last_name)
16 FROM (
17     SELECT student_name, se.course_no, p.last_name
18     FROM students s
19     INNER JOIN student_enrollment se
20     ON s.student_no = se.student_no
21     INNER JOIN teach t
22     ON se.course_no = t.course_no
23     INNER JOIN professors p
24     ON t.last_name = p.last_name
25 ) a
26 GROUP BY student_name, course_no
27 ORDER BY student_name, course_no;
28
29 /* write a query that returns employees whose salary is above average for
30    their given department. */
31 SELECT first_name
32 FROM employees outer_emp
33 WHERE salary > (
34     SELECT AVG(salary)
35     FROM employees
36     WHERE department = outer_emp.department);
```