

Juan Ignacio Galvalisi

Exercise 3.9: Common Table Expressions (CTE)

Step 1: Answer the business questions from step 1 and 2 of task 3.8 using CTEs

1. Rewrite your queries from steps 1 and 2 of task 3.8 as CTEs.

3.8 - Step 1

Rockbuster/postgres@PostgreSQL 14 ▾


Query Editor Query History

```
1 SELECT
2     ROUND(AVG(total_amount_paid),2) AS average
3 FROM
4     (SELECT
5         A.customer_id,
6         A.first_name,
7         A.last_name,
8         D.city,
9         E.country,
10        SUM(B.amount) AS total_amount_paid
11     FROM customer A
12     INNER JOIN payment B ON A.customer_id = B.customer_id
13     INNER JOIN address C ON A.address_id = C.address_id
14     INNER JOIN city D ON C.city_id = D.city_id
15     INNER JOIN country E ON D.country_id = E.country_id
16     WHERE D.city IN ('Aurora', 'Atlixco', 'Xintai', 'Adoni', 'Dhule (Dhulia)',
17                     'Kurashiki', 'Pingxiang', 'Sivas', 'Celaya', 'So Leopoldo')
18     GROUP BY
19         A.customer_id,
20         A.first_name,
21         A.last_name,
22         D.city,
23         E.country
24     ORDER BY total_amount_paid DESC
25     LIMIT 5) AS total_amount_paid;
26
```

Data Output Explain Messages Notifications

	average numeric	
1	107.35	

3.8 - Step 2


Rockbuster/postgres@PostgreSQL 14

Query Editor
Query History

```

1  SELECT
2      DISTINCT (A.country),
3      COUNT(DISTINCT D.address_id) AS all_customer_count,
4      COUNT(DISTINCT top_5_customers) AS top_customer_count
5  FROM country A
6  INNER JOIN city B ON A.country_id = B.country_id
7  INNER JOIN address C ON B.city_id = C.city_id
8  INNER JOIN customer D ON C.address_id = D.address_id
9  LEFT JOIN
10 (SELECT
11     A.customer_id,
12     A.first_name,
13     A.last_name,
14     D.city,
15     E.country,
16     SUM(B.amount) AS total_amount_paid
17 FROM customer A
18 INNER JOIN payment B ON A.customer_id = B.customer_id
19 INNER JOIN address C ON A.address_id = C.address_id
20 INNER JOIN city D ON C.city_id = D.city_id
21 INNER JOIN country E ON D.country_id = E.country_id
22 WHERE
23 e.country IN ('India', 'China', 'United States', 'Japan', 'Mexico', 'Brazil',
24               'Russian Federation', 'Philippines', 'Turkey', 'Indonesia')
25 AND D.city IN ('Aurora', 'Atlixco', 'Xintai', 'Adoni', 'Dhule (Dhulia)',
26               'Kurashiki', 'Pingxiang', 'Sivas', 'Celaya', 'So Leopoldo')
27 GROUP BY
28     A.customer_id,
29     A.first_name,
30     A.last_name,
31     D.city,
32     E.country
33 ORDER BY total_amount_paid DESC
34 LIMIT 5) AS top_5_customers
35 ON A.country = top_5_customers.country
36 GROUP BY A.country
37 ORDER BY top_customer_count DESC
38 LIMIT 4;


```

Data Output
Explain
Messages
Notifications

	country character varying (50)	all_customer_count bigint	top_customer_count bigint
1	Mexico	30	2
2	India	60	1
3	Turkey	15	1
4	United States	36	1

2. Copy-paste your CTEs and their outputs into your answers document.

3.9 - Step 1


Rockbuster/postgres@PostgreSQL 14

Query Editor
Query History

```


1  WITH
2  final_average (customer, first_name, last_name, city, country, total_amount_paid) AS
3  (SELECT
4      A.customer_id,
5      A.first_name,
6      A.last_name,
7      D.city,
8      E.country,
9      SUM(B.amount) AS total_amount_paid
10 FROM customer A
11 INNER JOIN payment B ON A.customer_id = B.customer_id
12 INNER JOIN address C ON A.address_id = C.address_id
13 INNER JOIN city D ON C.city_id = D.city_id
14 INNER JOIN country E ON D.country_id = E.country_id
15 WHERE D.city IN ('Aurora', 'Atlixco', 'Xintai', 'Adoni', 'Dhule (Dhulia)',
16 'Kurashiki', 'Pingxiang', 'Sivas', 'Celaya', 'So Leopoldo')
17 GROUP BY
18     A.customer_id,
19     A.first_name,
20     A.last_name,
21     D.city,
22     E.country
23 ORDER BY total_amount_paid DESC
24 LIMIT 5)
25 SELECT
26     ROUND(AVG(total_amount_paid),2) AS average
27 FROM final_average

```

Data Output
Explain
Messages
Notifications

	average numeric
1	107.35

3.9 - Step 2


Rockbuster/postgres@PostgreSQL 14

Query Editor
Query History

```

1  WITH customers_top_5 (customer_id, first_name, last_name, city, country, total_amount_paid) AS
2      (SELECT
3          A.customer_id,
4          A.first_name,
5          A.last_name,
6          D.city,
7          E.country,
8          SUM(B.amount) AS total_amount_paid
9      FROM customer A
10     INNER JOIN payment B ON A.customer_id = B.customer_id
11     INNER JOIN address C ON A.address_id = C.address_id
12     INNER JOIN city D ON C.city_id = D.city_id
13     INNER JOIN country E ON D.country_id = E.country_id
14     WHERE
15         e.country IN ('India', 'China', 'United States', 'Japan', 'Mexico', 'Brazil',
16                     'Russian Federation', 'Philippines', 'Turkey', 'Indonesia')
17         AND D.city IN ('Aurora', 'Atlixco', 'Xintai', 'Adoni', 'Dhule (Dhulia)',
18                     'Kurashiki', 'Pingxiang', 'Sivas', 'Celaya', 'So Leopoldo')
19     GROUP BY
20         A.customer_id,
21         A.first_name,
22         A.last_name,
23         D.city,
24         E.country
25     ORDER BY total_amount_paid DESC
26     LIMIT 5)
27  SELECT
28      DISTINCT (A.country),
29      COUNT(DISTINCT D.address_id) AS all_customer_count,
30      COUNT(DISTINCT customers_top_5) AS top_customer_count
31  FROM country A
32  INNER JOIN city B ON A.country_id = B.country_id
33  INNER JOIN address C ON B.city_id = C.city_id
34  INNER JOIN customer D ON C.address_id = D.address_id
35  LEFT JOIN customers_top_5 ON A.country = customers_top_5.country
36  GROUP BY A.country
37  HAVING COUNT(DISTINCT customers_top_5) > 0
38  ORDER BY
39      top_customer_count DESC,
40      all_customer_count DESC;

```

Data Output
Explain
Messages
Notifications

	country character varying (50)	all_customer_count bigint	top_customer_count bigint
1	Mexico	30	2
2	India	60	1
3	United States	36	1
4	Turkey	15	1

3. Write 2 to 3 sentences explaining how you approached this step, for example, what you did first, second, and so on.

In order to get my script more readable, I decided to get into my CTE the subquery based on the top 5 five customers who have spent the most on Rockbuster. With that in mind, I used the WITH function and named my CTE, placing it at the beginning of my query, just before the main statement.

Then, I wrote my main statement, which has the ultimate names of my three categories: country, all_customer_count, and top_customer_count regarding the requested task ("Find out how many of the top 5 customers are based within each country".) Moreover, I used the LEFT JOIN because I wanted to combine my main query with my CTE.

Step 2: Compare the performance of your CTEs and subqueries

- Which approach do you think will perform better and why?

Initially, I would assume CTE performs better because it is better organized, at least for the human eye. The way computers read this code differs from how humans read it, so both readings may vary.

However, we must check each case in detail because there would be more variables to consider to perform our queries in time and cost through the EXPLAIN + ANALYZE command.

Compare the costs of all the queries by creating query plans for each one.

QUERY	COST	POP-UP WINDOW	PLANNING TIME	EXECUTION TIME
Step 1 (3.8.) – SQ	(cost=67.91..67.92 rows=1 width=32)	39 msc. 22 rows affected	2.600 ms	3.069 ms
Step 1 (3.9.) – CTE	(cost=67.91..67.92 rows=1 width=32)	34 msc. 22 rows affected	1.487 ms	1.076 ms
Step 2 (3.8.) – SQ	(cost=130.61..130.97 rows=36 width=25)	107 msc. 47 rows affected	2.511 ms	9.225 ms
Step 2 (3.9.) – CTE	(cost=130.61..130.97 rows=36 width=25)	77 msc. 47 rows affected	3.558	3.133

- Did the results surprise you? Write a few sentences to explain your answer.

At first, I would have thought that queries with CTEs would be more efficient than those containing subqueries. With the results in hand, it seems that this rule is fulfilled in general terms.

Regarding the cost required by each query with its characteristics obtained through the EXPLAIN command, we cannot observe differences between both pairs. Now, concerning the final execution times received automatically by PostgreSQL when a query is made, we find slight differences between both procedures. Since these are small data sets and short queries, the difference may be more significant as both the data to be analyzed and the complexity of the queries increase.

We can observe two absorbing matters if we analyze the data provided by the EXPLAIN ANALYZE command. For one thing, planning and execution times are higher when subqueries are used than when they are not. On the other hand, subqueries perform their task in a longer time than estimated, while CTEs do it in less time than estimated. Not only do they work better than expected, they even do it almost three times faster than the same query, but they include subqueries inside. In this sense, it is appropriate to ask, is it not necessary to pay more attention to the final time executed than the cost each query represents?

Step 3

Write 1 to 2 paragraphs on the challenges you faced when replacing your subqueries with CTEs.

At first, it was a bit difficult to get out of the abstraction and materialize the tables now with the new CET, but it is more than achievable if you think more carefully about the syntax. Once you know the syntax, it is easier to write your queries, and that rule applies here too.

On the other hand, it would have been a different situation if I had made the query from scratch since I had almost everything written thanks to exercise 3.8. In this sense, I had to configure the CET, join it to the main query, and execute it as general steps.