

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

3.8.1 Find the average amount paid by the top 5 customers

Query		Query History
1	WITH top_customers AS	
2	(SELECT c.customer_id,	
3	c.first_name,	
4	c.last_name,	
5	co.country,	
6	ci.city,	
7	SUM(p.amount) AS total_amount_paid	
8	FROM customer c	
9	JOIN address a ON c.address_id = a.address_id	
10	JOIN city ci ON a.city_id = ci.city_id	
11	JOIN country co ON ci.country_id = co.country_id	
12	JOIN payment p ON c.customer_id = p.customer_id	
13	WHERE ci.city IN ('Aurora',	
14	'Acua',	
15	'Citrus Heights',	
16	'Iwaki',	
17	'Ambattur',	
18	'Shanwei',	
19	'So Leopoldo',	
20	'Teboksary',	
21	'Tianjin',	
22	'Cianju')	
23	GROUP BY c.customer_id,	
24	c.first_name,	
25	c.last_name,	
26	co.country,	
27	ci.city	
28	ORDER BY total_amount_paid DESC	
29	LIMIT 5)	
30	SELECT AVG(total_amount_paid) AS average	
31	FROM top_customers;	
Data Output		Messages Notifications
	average	
	numeric	
1	105.5540000000000000	

The original query contained a subquery to find the top 5 customers based on certain criteria. The first step was to recognize this subquery and understand its purpose.

The next step involved transforming the subquery into a CTE. I named the CTE **top_customers** and placed the existing subquery logic inside it. This helped in making the code cleaner.

Then I modified the main query to reference the CTE instead of the subquery by replacing the subquery reference with the CTE name (**top_customers**) and adjusting the join conditions. To enhance readability, I used aliases for the table names in both the CTE and the main query.

Throughout the process, I ensured that the changes did not alter the logic of the original query. The CTE was correctly referenced in the main query, and the join conditions remained accurate.

3.8.2 Find out how many of the top 5 customers are based within each country

Query Query History

```
1 WITH top_customers AS
2 (SELECT c.customer_id,
3      c.first_name,
4      c.last_name,
5      co.country_id,
6      ci.city,
7      SUM(p.amount) AS total_amount_paid
8 FROM customer c
9 INNER JOIN address a ON c.address_id = a.address_id
10 INNER JOIN city ci ON a.city_id = ci.city_id
11 INNER JOIN country co ON ci.country_id = co.country_id
12 INNER JOIN payment p ON c.customer_id = p.customer_id
13 WHERE ci.city IN ('Aurora',
14                  'Acua',
15                  'Citrus Heights',
16                  'Iwaki',
17                  'Ambattur',
18                  'Shanwei',
19                  'So Leopoldo',
20                  'Teboksary',
21                  'Tianjin',
22                  'Cianju'))
23 GROUP BY c.customer_id,
24          c.first_name,
25          c.last_name,
26          co.country_id,
27          ci.city
28 ORDER BY total_amount_paid DESC
29 LIMIT 5)
30 SELECT co.country,
31        COUNT(DISTINCT c.customer_id) AS all_customer_count,
32        COUNT(DISTINCT tc.customer_id) AS top_customer_count
33 FROM customer c
34 INNER JOIN address a ON c.address_id = a.address_id
35 INNER JOIN city ci ON a.city_id = ci.city_id
36 INNER JOIN country co ON ci.country_id = co.country_id
37 LEFT JOIN top_customers tc ON c.customer_id = tc.customer_id
38 GROUP BY co.country
39 ORDER BY all_customer_count DESC
40 LIMIT 5;
```

Data Output Messages Notifications

	country character varying (50)	all_customer_count bigint	top_customer_count bigint
1	India	60	1
2	China	53	1
3	United States	36	1
4	Japan	31	1
5	Mexico	30	1

Total rows: 5 of 5 Query complete 00:00:00.059

To achieve our objective, I created a CTE named 'top_customers' to nest the subquery, retrieving necessary information from joined tables. In this subquery, I applied filters for specific cities, aggregated data for each customer, and ordered the results by the total amount paid in descending order, limiting to the top 5 customers.

In the main query, I employed INNER JOINS to establish connections between the 'customer', 'address', 'city', and 'country' tables based on their respective foreign keys. A LEFT JOIN was then performed with the top_customers' CTE on the 'customer_id' to bring in information about the top 5 customers. Using COUNT(DISTINCT), I calculated both the total number of customers ('all_customer_count') and the count of top customers ('top_customer_count') in each country. The results were grouped by country and ordered by 'all_customer_count' in descending order, with a LIMIT of 5 to retrieve the top 5 countries.

In terms of readability and consistency, I maintained clear aliasing for table names and columns throughout the query.








The logical structure ensures ease of comprehension. I tested the query in a database environment to confirm correctness and performance, especially with larger datasets. This step-by-step approach facilitates clarity, testing, and potential future modifications.

Step 2: Compare the performance of your CTEs and subqueries.

Average Amount Paid by the Top 5 Customers

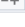
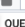
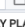
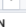
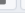

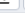
CTE

Query	Query History
1	EXPLAIN
2	WITH top_customers AS
3	(SELECT c.customer_id,
4	c.first_name,
5	c.last_name,
6	co.country,
7	ci.city,
8	SUM(p.amount) AS total_amount_paid
9	FROM customer c
10	JOIN address a ON c.address_id = a.address_id
11	JOIN city ci ON a.city_id = ci.city_id
12	JOIN country co ON ci.country_id = co.country_id
13	JOIN payment p ON c.customer_id = p.customer_id
14	WHERE ci.city IN ('Aurora',
15	'Acua',
16	'Citrus Heights',
17	'Iwaki',
18	'Ambattur',
19	'Shanwei',
20	'So Leopoldo',
21	'Teboksary',
22	'Tianjin',
23	'Cianju')
24	GROUP BY c.customer_id,
25	c.first_name,
26	c.last_name,
27	co.country,
28	ci.city
29	ORDER BY total_amount_paid DESC
30	LIMIT 5)
31	SELECT AVG(total_amount_paid) AS average
32	FROM top_customers;

Data Output	Messages	Notifications
      		
QUERY PLAN		
text		
1	Aggregate	(cost=64.49..64.50 rows=1 width=32)
2	-> Limit	(cost=64.41..64.43 rows=5 width=67)
3	-> Sort	(cost=64.41..65.02 rows=244 width=67)
4	Sort Key:	(sum(p.amount)) DESC
5	-> HashAggregate	(cost=57.31..60.36 rows=244 width=67)
Total rows: 22 of 22 Query complete 00:00:00.254		

Subquery

Query	Query History
1	EXPLAIN
2	SELECT AVG(total_amount_paid) AS average
3	FROM
4	(SELECT c.customer_id,
5	c.first_name,
6	c.last_name,
7	co.country,
8	ci.city,
9	SUM(p.amount) AS total_amount_paid
10	FROM customer c
11	JOIN address a ON c.address_id = a.address_id
12	JOIN city ci ON a.city_id = ci.city_id
13	JOIN country co ON ci.country_id = co.country_id
14	JOIN payment p ON c.customer_id = p.customer_id
15	WHERE ci.city IN ('Aurora',
16	'Acua',
17	'Citrus Heights',
18	'Iwaki',
19	'Ambattur',
20	'Shanwei',
21	'So Leopoldo',
22	'Teboksary',
23	'Tianjin',
24	'Cianju')
25	GROUP BY c.customer_id,
26	c.first_name,
27	c.last_name,
28	co.country,
29	ci.city
30	ORDER BY total_amount_paid DESC
31	LIMIT 5) AS top_customers;

Data Output	Messages	Notifications
      		
QUERY PLAN		
text		
1	Aggregate	(cost=64.49..64.50 rows=1 width=32)
2	-> Limit	(cost=64.41..64.43 rows=5 width=67)
3	-> Sort	(cost=64.41..65.02 rows=244 width=67)
4	Sort Key:	(sum(p.amount)) DESC
5	-> HashAggregate	(cost=57.31..60.36 rows=244 width=67)
Total rows: 22 of 22 Query complete 00:00:00.046		

Aggregate (cost=64.49..64.50 rows=1 width=32)

Top 5 Customers within Each Country

CTE

Query	Query History
1	EXPLAIN
2	WITH top_customers AS
3	(SELECT c.customer_id,
4	c.first_name,
5	c.last_name,
6	co.country_id,
7	ci.city,
8	SUM(p.amount) AS total_amount_paid
9	FROM customer c
10	INNER JOIN address a ON c.address_id = a.address_id
11	INNER JOIN city ci ON a.city_id = ci.city_id
12	INNER JOIN country co ON ci.country_id = co.country_id
13	INNER JOIN payment p ON c.customer_id = p.customer_id
14	WHERE ci.city IN ('Aurora',
15	'Acua',
16	'Citrus Heights',
17	'Iwak',
18	'Ambattur',
19	'Shanwei',
20	'So Leopoldo',
21	'Teboksary',
22	'Tianjin',
23	'Cianju')
24	GROUP BY c.customer_id,
25	c.first_name,
26	c.last_name,
27	co.country_id,
28	ci.city
29	ORDER BY total_amount_paid DESC
30	LIMIT 5)
31	SELECT co.country,
32	COUNT(DISTINCT c.customer_id) AS all_customer_count,
33	COUNT(DISTINCT tc.customer_id) AS top_customer_count
34	FROM customer c
35	INNER JOIN address a ON c.address_id = a.address_id
36	INNER JOIN city ci ON a.city_id = ci.city_id
37	INNER JOIN country co ON ci.country_id = co.country_id
38	LEFT JOIN top_customers tc ON c.customer_id = tc.customer_id
39	GROUP BY co.country
40	ORDER BY all_customer_count DESC
41	LIMIT 5;

Data Output	Messages	Notifications
QUERY PLAN text		
1	Limit (cost=166.66..166.68 rows=5 width=25)	
2	→ Sort (cost=166.66..166.94 rows=109 width=25)	
3	Sort Key: (count(DISTINCT c.customer_id)) DESC	
4	→ GroupAggregate (cost=157.77..164.85 rows=109 width=25)	
5	Group Key: co.country	
6	→ Sort (cost=157.77..159.27 rows=999 width=17)	
7	Sort Key: co.country, c.customer_id	
Total rows: 45 of 45 Query complete 00:00:00.072		

Subquery

Query	Query History
1	EXPLAIN
2	SELECT co.country,
3	COUNT(DISTINCT c.customer_id) AS all_customer_count,
4	COUNT(DISTINCT tc.customer_id) AS top_customer_count
5	FROM customer c
6	INNER JOIN address a ON c.address_id = a.address_id
7	INNER JOIN city ci ON a.city_id = ci.city_id
8	INNER JOIN country co ON ci.country_id = co.country_id
9	LEFT JOIN
10	(SELECT c.customer_id,
11	c.first_name,
12	c.last_name,
13	co.country_id,
14	ci.city
15	FROM customer c
16	INNER JOIN address a ON c.address_id = a.address_id
17	INNER JOIN city ci ON a.city_id = ci.city_id
18	INNER JOIN country co ON ci.country_id = co.country_id
19	INNER JOIN payment p ON c.customer_id = p.customer_id
20	WHERE ci.city IN ('Aurora',
21	'Acua',
22	'Citrus Heights',
23	'Iwak',
24	'Ambattur',
25	'Shanwei',
26	'So Leopoldo',
27	'Teboksary',
28	'Tianjin',
29	'Cianju')
30	GROUP BY c.customer_id,
31	c.first_name,
32	c.last_name,
33	co.country_id,
34	ci.city
35	ORDER BY SUM(p.amount) DESC
36	LIMIT 5) tc ON c.customer_id = tc.customer_id
37	GROUP BY co.country
38	ORDER BY all_customer_count DESC
39	LIMIT 5;

Data Output	Messages	Notifications
QUERY PLAN text		
1	Limit (cost=166.66..166.68 rows=5 width=25)	
2	→ Sort (cost=166.66..166.94 rows=109 width=25)	
3	Sort Key: (count(DISTINCT c.customer_id)) DESC	
4	→ GroupAggregate (cost=157.77..164.85 rows=109 width=25)	
5	Group Key: co.country	
6	→ Sort (cost=157.77..159.27 rows=999 width=17)	
7	Sort Key: co.country, c.customer_id	
Total rows: 45 of 45 Query complete 00:00:00.091		

"Limit (cost=166.66..166.68 rows=5 width=25)"

1. I initially thought that using CTEs might result in better performance compared to subqueries. However, in my testing, both approaches exhibited the same performance in terms of costs.
2. When comparing the costs of all the queries by examining the query plans, I found that the costs for CTEs and subqueries were similar. The query planner estimated comparable costs for both methods, indicating that there wasn't a significant difference in performance.
3. The results did surprise me to some extent. I anticipated that CTEs, being a more structured and readable way to organize complex queries, might have some performance benefits. However, the similarity in costs suggests that, in the specific context of these queries, the database optimizer treated CTEs and subqueries similarly in terms of execution plans and efficiency.

Step 3

Incorporating Common Table Expressions (CTEs) instead of subqueries posed some challenges. Initially, understanding the concept of CTEs required a bit of extra effort, as it was a new technique for me. Figuring out how to structure the CTE and connect it seamlessly with the main query was a bit tricky, but examples and practice helped to overcome this hurdle.

The transition from subqueries to CTEs involved adapting my mindset to a more modular approach. While subqueries are embedded within the main query, CTEs provide a cleaner and more organized way to structure complex queries. Once I grasped the idea of breaking down the task into smaller, manageable parts with a CTE, it became a valuable tool for improving code readability and maintainability. Overall, the challenges were more about adapting to a new methodology, but the benefits in terms of code structure and clarity made the learning process worthwhile.