# Data Analytics Immersion
## 3.9: Common Table Expressions
*Mehreen Becker*

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

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

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

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

Rockbuster/postgres@My Server

Query    Query History                                                                Scratch Pad ×

```sql
1    WITH top_countries AS (
2        SELECT
3            D.country
4        FROM customer A
5        INNER JOIN address B ON A.address_id = B.address_id
6        INNER JOIN city C ON B.city_id = C.city_id
7        INNER JOIN country D ON C.country_id = D.country_id
8        GROUP BY D.country
9        ORDER BY COUNT(A.customer_id) DESC
10       LIMIT 10
11   ),
12   top_cities AS (
13       SELECT
14           D.country,
15           C.city
16       FROM customer A
17       INNER JOIN address B ON A.address_id = B.address_id
18       INNER JOIN city C ON B.city_id = C.city_id
19       INNER JOIN country D ON C.country_id = D.country_id
20       WHERE D.country IN (SELECT country FROM top_countries)
21       GROUP BY D.country, C.city
22       ORDER BY COUNT(A.customer_id) DESC
23       LIMIT 10
```

Data Output    Messages    Explain ×    Notifications

Showing rows: 1 to 1    Page No: 1    of 1

|   | average_amount_paid numeric |
|---|---|
| 1 | 105.5540000000000000 |

Total rows: 1    Query complete 00:00:00.073                          LF    Ln 7, Col 11

---

Rockbuster/postgres@My Server

Query    Query History                                                                Scratch Pad ×

```sql
22       ORDER BY COUNT(A.customer_id) DESC
23       LIMIT 10
24   ),
25   customer_totals AS (
26       SELECT
27           A.customer_id,
28           A.first_name,
29           A.last_name,
30           D.country,
31           C.city,
32           SUM(P.amount) AS total_amount_paid
33       FROM customer A
34       INNER JOIN address B ON A.address_id = B.address_id
35       INNER JOIN city C ON B.city_id = C.city_id
36       INNER JOIN country D ON C.country_id = D.country_id
37       INNER JOIN payment P ON A.customer_id = P.customer_id
38       WHERE C.city IN (SELECT city FROM top_cities)
39       GROUP BY A.customer_id, A.first_name, A.last_name, D.country, C.city
40       ORDER BY total_amount_paid DESC
41       LIMIT 5
42   )
43   SELECT AVG(total_amount_paid) AS average_amount_paid
44   FROM customer_totals;
```

Data Output    Messages    Explain ×    Notifications

Showing rows: 1 to 1    Page No: 1    of 1

|   | average_amount_paid numeric |
|---|---|
| 1 | 105.5540000000000000 |

Total rows: 1    Query complete 00:00:00.073                          LF    Ln 17, Col 11

pgAdmin 4

Dependents ✕   Processes ✕   ▤ Testing/postgres… ✕   ▤ Testing/postgres… ✕   ▤ postgres/postgres… ✕   Preferences ✕   ▤ Rockbuster/postgres@My Server* ✕   ⌄   ⋮

Rockbuster/postgres@My Server

No limit

Query   Query History                                                                        Scratch Pad ✕

```sql
1   SELECT
2       d.country,
3       COUNT(DISTINCT a.customer_id) AS all_customer_count,
4       COUNT(DISTINCT CASE
5           WHEN (
6               SELECT SUM (p2.amount)
7               FROM payment p2
8               WHERE p2.customer_id = a.customer_id
9           ) > (
10              SELECT AVG(total_sum)
11              FROM (
12                  SELECT SUM(p.amount) AS total_sum
13                  FROM payment p
14                  GROUP BY p.customer_id
15              ) AS customer_totals
16              )
17          THEN a.customer_id
18          ELSE NULL
19          END) AS top_customer_count
20  FROM customer a
21  JOIN address b ON a.address_id = b.address_id
22  JOIN city c ON b.city_id = c.city_id
23  JOIN country d ON c.country_id = d.country_id
24  GROUP BY d.country
25  ORDER BY top_customer_count DESC
26  LIMIT 10;
```

Data Output   Messages   Explain ✕   Notifications

Showing rows: 1 to 10   Page No: 1   of 1

| country character varying (50) | all_customer_count bigint | top_customer_count bigint |
|---|---|---|

Total rows: 10   Query complete 00:00:00.300                                          LF   Ln 2, Col 12

---

pgAdmin 4

Dependents ✕   Processes ✕   ▤ Testing/postgres… ✕   ▤ Testing/postgres… ✕   ▤ postgres/postgres… ✕   Preferences ✕   ▤ Rockbuster/postgres@My Server* ✕   ⌄   ⋮

Rockbuster/postgres@My Server

No limit

Query   Query History                                                                        Scratch Pad ✕

```sql
14              GROUP BY p.customer_id
15          ) AS customer_totals
16          )
17          THEN a.customer_id
18          ELSE NULL
19          END) AS top_customer_count
20  FROM customer a
21  JOIN address b ON a.address_id = b.address_id
22  JOIN city c ON b.city_id = c.city_id
23  JOIN country d ON c.country_id = d.country_id
24  GROUP BY d.country
25  ORDER BY top_customer_count DESC
26  LIMIT 10;
```

Data Output   Messages   Explain ✕   Notifications

Showing rows: 1 to 10   Page No: 1   of 1

| | country character varying (50) | all_customer_count bigint | top_customer_count bigint |
|---|---|---|---|
| 1 | India | 60 | 26 |
| 2 | China | 53 | 25 |
| 3 | United States | 36 | 16 |
| 4 | Japan | 31 | 14 |
| 5 | Russian Federation | 28 | 13 |
| 6 | Brazil | 28 | 12 |
| 7 | Mexico | 30 | 11 |
| 8 | Philippines | 20 | 11 |
| 9 | Taiwan | 10 | 7 |
| 10 | Turkey | 15 | 7 |

Total rows: 10   Query complete 00:00:00.300                                          LF   Ln 2, Col 12

I first identified the top countries with the most customers by joining the customer, address, city and country tables and grouping them by country. Then, I narrowed the focus to the top 10 cities located within those countries, again by counting customers per city and ordering by that count. After that, I calculated information about the top 5 customers who paid the most using the payment table joined to the geographic data, and limited it to the customers from the top city. Finally I summarized the results in one query by find the average of the top 5 total payments, and in the other by comparing counts of all customers versus top customers per country.

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

1. **Which approach do you think will perform better and why?** I'm sure that both CTEs and subqueries produce the same, if not a similar, result. With my limited experience with CTEs, I feel like subqueries allow for better optimization potentially and I also felt like I had to scan fewer rows.

2. **Compare the costs of all the queries by creating query plans for each one.**
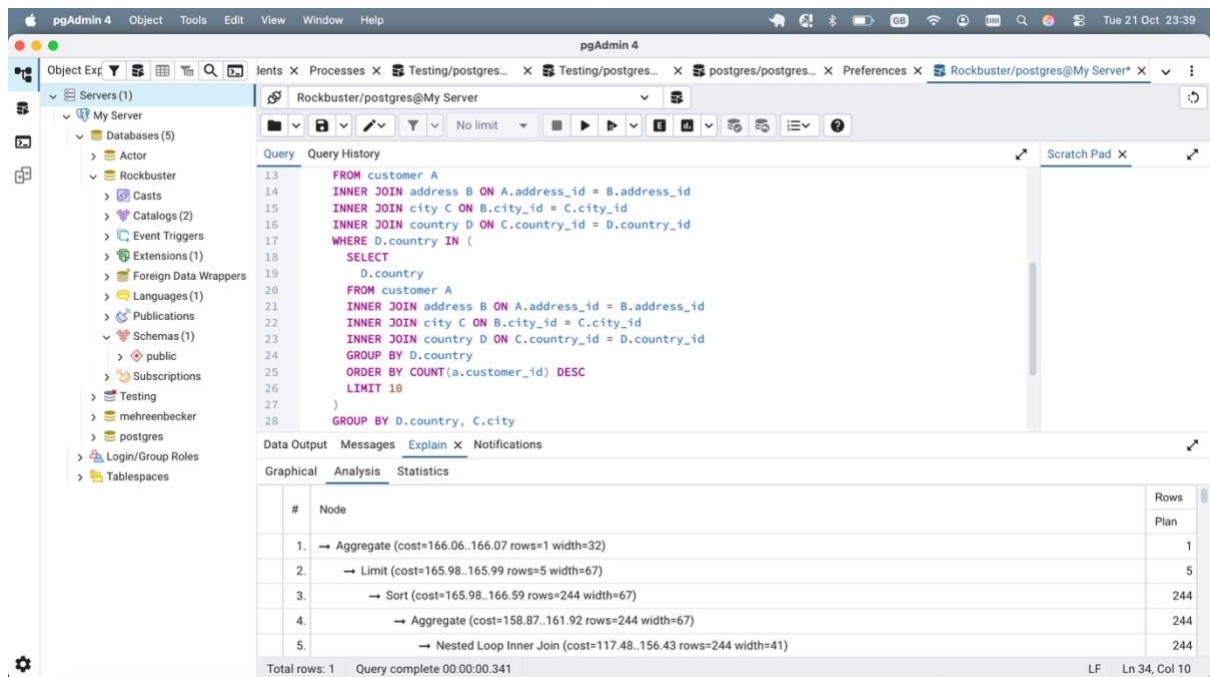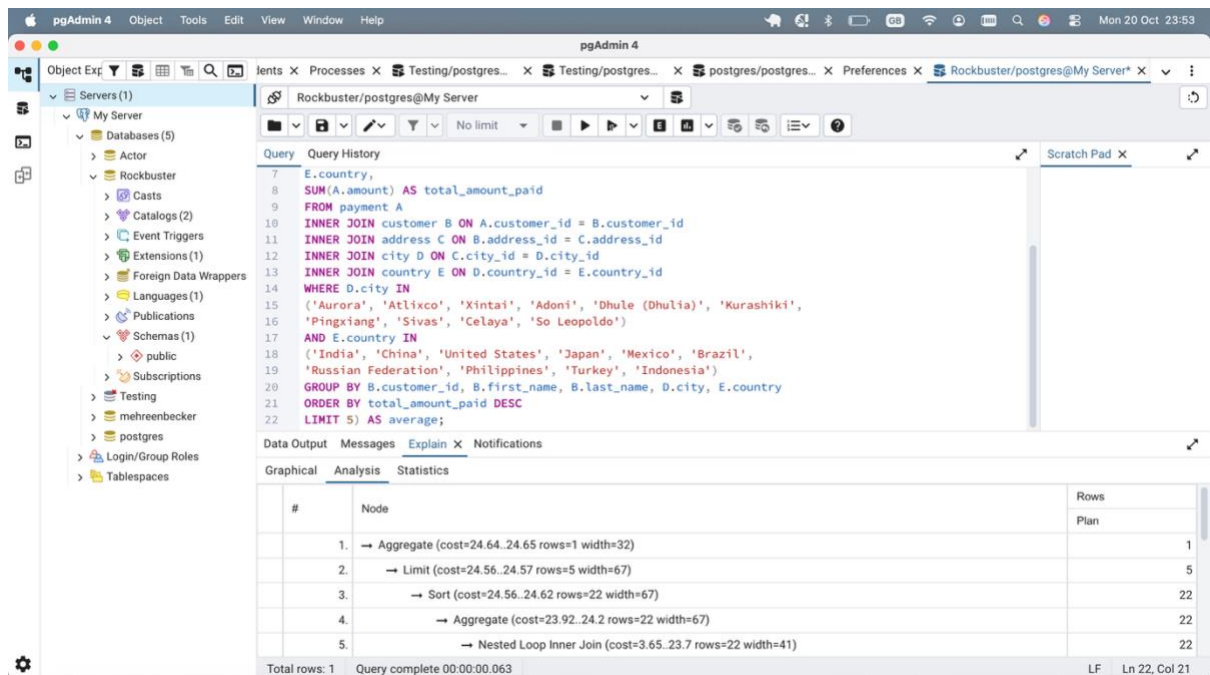
Difference in cost:



*Figure 1: Cost (CTE)*



*Figure 2: Cost (Subquery)*

3. The **EXPLAIN** command gives you an *estimated* cost. To find out the actual speed of your queries, run them in pgAdmin 4. After you've run each query, a popup window will display its speed in milliseconds.
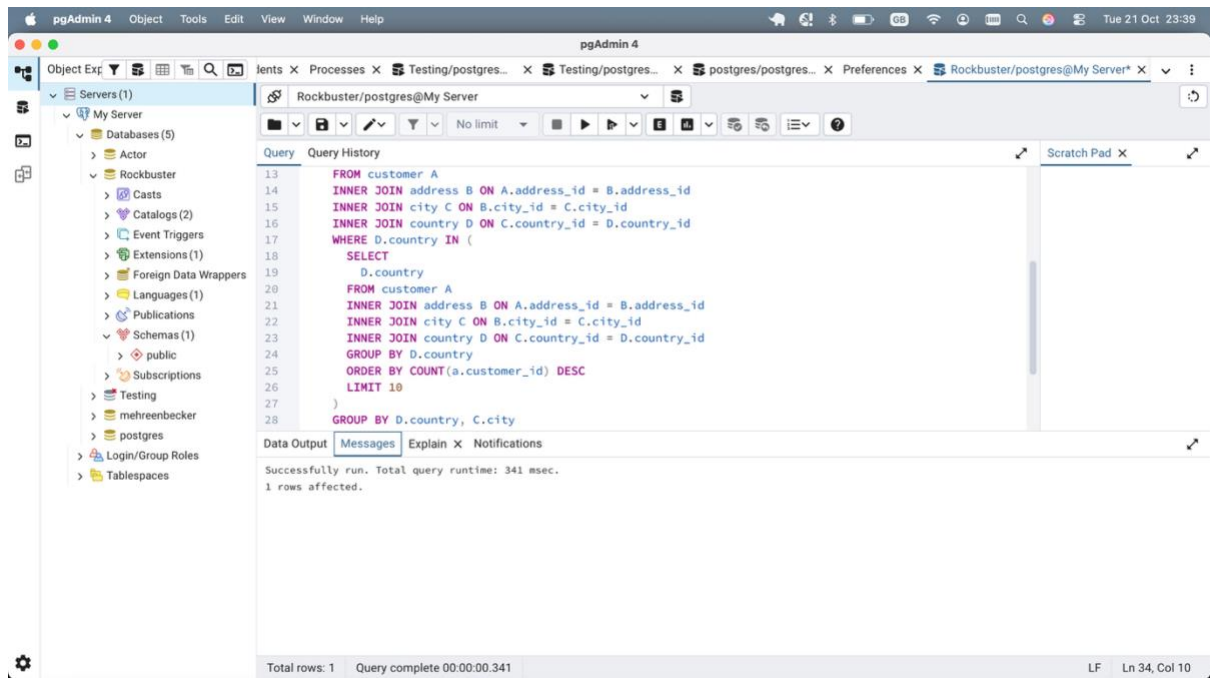
Difference in query run time:

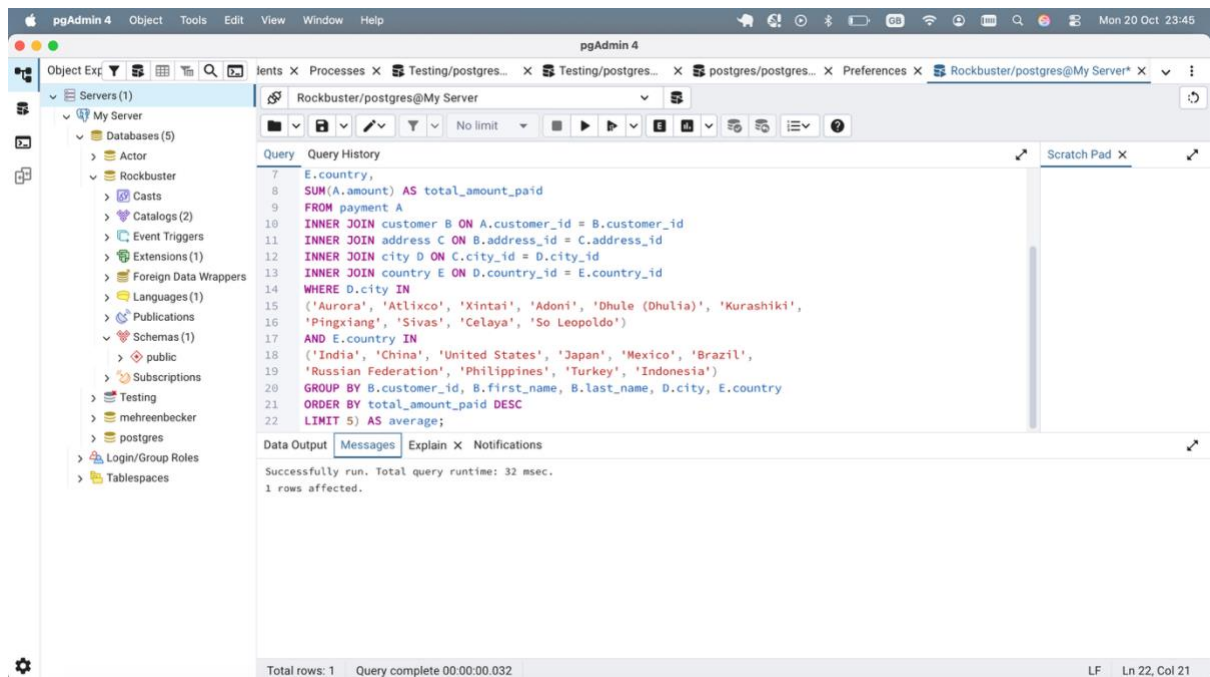*Figure 1: Total Query Runtime: 341 msec (CTE)*



*Figure 2: Total Query Runtime: 32 msec (Subquery)*

4. **Did the results surprise you?** Write a few sentences to explain your answer. To be honest, no. I did think that subqueries would run faster (not that I noticed the few milisecond difference) but more because they feel more straightforward and less complex.

**Step 3: Write 1 to 2 paragraphs on the challenges you faced when replacing your subqueries with CTEs.** Since we were asked to use CTEs and not subqueries, one of the main challenges for me was to make sure that everything made sense. I did find myself feeling a little lost and confused since it's easier with CTEs to misplace filters and to also join things together. I

found myself scanning rows multiple times to ensure that the string made sense and that everything was connected with each other and connected in a way that it should be.