

Advanced Query Techniques

DATA 503: Fundamentals of Data Engineering

Lucas P. Cordova, Ph.D.

2026-04-08

A concise version of advanced SQL patterns: restoring the practice database, subqueries in WHERE and FROM, IN and EXISTS, LATERAL for top-N per group, CTEs for readability, a short crosstab example, and CASE for classification. Expect checkpoints with hints, then solutions on the following slide. Based on Lecture 11-1 and Chapter 13 of Practical SQL, 2nd Edition.

Table of contents

1	Learning objectives	2
2	Restore the practice database	2
3	Subqueries: WHERE and FROM	3
4	IN and EXISTS	5
5	LATERAL joins	7
6	Review: Common Table Expressions (WITH)	9
7	Crosstab	10
8	CASE and classification	11
9	Summary	13

1 Learning objectives

1.1 What you should be able to do

1. Use scalar and derived-table subqueries for filtering and reshaping aggregates
2. Choose `IN` vs `EXISTS` / `NOT EXISTS` for membership tests and anti-joins
3. Write `LATERAL` joins for per-row computations and top-N patterns
4. Refactor multi-step logic with CTEs (`WITH`) and optional `CASE` classification

...

Course connection: These patterns show up in analytical SQL, incremental pipelines, and any time you need correct, readable transformations over relational data.

2 Restore the practice database

2.1 Create a Database

Create a database called `class_dbs` using the `createdb` command in a terminal OR using any method you usually use to create a new database.

```
1 createdb -h localhost -U postgres class_dbs
```

2.2 Download the files from Canvas

You will get **one** of these (not both required):

File	Format	Tool
<code>class_dbs.sql</code>	Plain SQL script	<code>psql</code> with <code>-f</code>
<code>class_dbs.dump</code>	Custom <code>pg_dump -Fc</code> archive	<code>pg_restore</code>

Save the file somewhere easy to find. The examples below assume it landed in your **Downloads** folder.

2.3 Local PostgreSQL: plain SQL (psql)

Open a terminal, go to the folder that contains `class_dbs.sql`, then run `psql` against the `class_dbs` database you created above.

```
1 cd ~/Downloads
2 psql -h localhost -U postgres -d class_dbs -f class_dbs.sql
```

2.4 Local PostgreSQL: custom archive (pg_restore)

The `class_dbs` database must exist before restore (same as the SQL path above).

```
1 cd ~/Downloads
2 pg_restore -h localhost -U postgres -d class_dbs --verbose class_dbs.dump
```

2.5 Railway: same files, database railway

Replace host and port with the values from your Railway Postgres service (Variables / Connect). The database name is often `railway`.

Plain SQL:

```
1 cd ~/Downloads
2 psql "postgresql://postgres:YOUR_PASSWORD@YOUR_HOST:YOUR_PORT/railway" -f class_dbs.sql
```

Custom dump:

```
1 cd ~/Downloads
2 pg_restore "postgresql://postgres:YOUR_PASSWORD@YOUR_HOST:YOUR_PORT/railway" --verbose class_dbs.dump
```

3 Subqueries: WHERE and FROM

3.1 Scalar subquery in WHERE

A **scalar** subquery returns one value. The database can evaluate it once, then use it like a constant in the outer query.

Run this. Counties at or above the national 90th percentile for `pop_est_2019`:

```

1 SELECT county_name, state_name, pop_est_2019
2 FROM us_counties_pop_est_2019
3 WHERE pop_est_2019 >= (
4     SELECT percentile_cont(0.9) WITHIN GROUP (ORDER BY pop_est_2019)
5     FROM us_counties_pop_est_2019
6 )
7 ORDER BY pop_est_2019 DESC;

```

This subquery is **uncorrelated**: it does not reference the outer row. It typically runs once, then the outer query filters.

3.2 Derived table in FROM

A subquery in FROM builds a **derived table**. PostgreSQL requires an alias (AS ...).

```

1 SELECT round(calcs.avg_pop, 0) AS avg_pop,
2         calcs.median_pop,
3         round(calcs.avg_pop - calcs.median_pop, 0) AS mean_minus_median
4 FROM (
5     SELECT avg(pop_est_2019) AS avg_pop,
6            percentile_cont(0.5)
7            WITHIN GROUP (ORDER BY pop_est_2019)::numeric AS median_pop
8     FROM us_counties_pop_est_2019
9 ) AS calcs;

```

Use this when you need to aggregate first, then treat the result as a table you can select from or join.

3.3 Checkpoint: above-average counties

Task: Return every county whose `pop_est_2019` is **strictly greater** than the national average.

Required columns: `county_name`, `state_name`, `pop_est_2019`

Requirements:

- Use a **scalar subquery** in WHERE that computes `avg(pop_est_2019)` over all counties
- Order by `pop_est_2019` descending

Expected shape (first rows, illustrative):

county_name	state_name	pop_est_2019
Los Angeles County	California	10039107
Cook County	Illinois	5150233
Harris County	Texas	4713320

Hint: Compare pop_est_2019 to (SELECT avg(pop_est_2019) FROM us_counties_pop_est_2019).

3.4 Solution: above-average counties

```

1 SELECT county_name, state_name, pop_est_2019
2 FROM us_counties_pop_est_2019
3 WHERE pop_est_2019 > (
4     SELECT avg(pop_est_2019)
5     FROM us_counties_pop_est_2019
6 )
7 ORDER BY pop_est_2019 DESC;
```

Why fewer than half the rows?

4 IN and EXISTS

4.1 IN (subquery)

IN tests membership against the set returned by the subquery.

```

1 SELECT first_name, last_name
2 FROM employees
3 WHERE emp_id IN (
4     SELECT id FROM retirees
5 )
6 ORDER BY emp_id;
```

Readable when the subquery is small and uncorrelated.

4.2 EXISTS and NOT EXISTS

EXISTS returns true if the subquery returns **at least one row**. The subquery can reference outer columns; that makes it **correlated** (conceptually: evaluated per outer row, though the planner may rewrite).

```
1 SELECT first_name, last_name
2 FROM employees AS e
3 WHERE EXISTS (
4     SELECT 1
5     FROM retirees AS r
6     WHERE r.id = e.emp_id
7 );
```

NOT EXISTS is the usual pattern for anti-joins: rows in A with no match in B. Prefer **NOT EXISTS** over **NOT IN** when the subquery can produce NULLs, because **NOT IN** with NULLs in the list behaves badly.

4.3 Checkpoint: active employees

Task: List employees who are **not** in retirees.

Required columns: `first_name`, `last_name`

Requirements: Use **NOT EXISTS** with a correlated subquery matching `retirees.id` to `employees.emp_id`. Order by `last_name`, `first_name`

Expected shape:

<code>first_name</code>	<code>last_name</code>
Samuel	Cole
Arthur	Pappas
...	...

...

Hint: `SELECT 1` inside **EXISTS** is enough; you only care whether a row appears.

4.4 Solution: active employees

```
1 SELECT first_name, last_name
2 FROM employees AS e
3 WHERE NOT EXISTS (
4     SELECT 1
5     FROM retirees AS r
6     WHERE r.id = e.emp_id
7 )
8 ORDER BY last_name, first_name;
```

5 LATERAL joins

5.1 Top-N per group

LATERAL lets a subquery in FROM reference columns from earlier tables. Classic use: **for each parent row**, run a query and keep a few rows.

...

For each teacher, return their two most recent lab accesses.

```
1 SELECT t.first_name, t.last_name, a.access_time, a.lab_name
2 FROM teachers AS t
3 LEFT JOIN LATERAL (
4     SELECT *
5     FROM teachers_lab_access
6     WHERE teacher_id = t.id
7     ORDER BY access_time DESC
8     LIMIT 2
9 ) AS a ON true
10 ORDER BY t.id;
```

...

LEFT JOIN LATERAL keeps teachers with zero visits (NULLs in a.*). ON true satisfies join syntax; filtering is inside the lateral subquery.

5.2 Checkpoint: one visit per teacher

Task: Same as the demo, but return only each teacher's **single most recent** lab access.

...

Required columns:

- first_name
- last_name
- access_time
- lab_name

...

Requirements:

- LEFT JOIN LATERAL with ORDER BY access_time DESC and LIMIT 1

...

Expected shape (illustrative):

first_name	last_name	access_time	lab_name
Janet	Smith	2022-12-21 ...	Chemistry A
Lee	Reynolds

...

Hint: Change one number from the demo query.

5.3 Solution: one visit per teacher

```
1 SELECT t.first_name, t.last_name, a.access_time, a.lab_name
2 FROM teachers AS t
3 LEFT JOIN LATERAL (
4     SELECT *
5     FROM teachers_lab_access
6     WHERE teacher_id = t.id
7     ORDER BY access_time DESC
8     LIMIT 1
9 ) AS a ON true
10 ORDER BY t.id;
```

6 Review: Common Table Expressions (WITH)

6.1 Why CTEs

A CTE names an intermediate result. Use them when you:

...

- chain several logical steps
- reuse the same expression instead of repeating a subquery
- want readers (including future you) to follow the story top to bottom

...

```
1 WITH big_counties AS (  
2     SELECT county_name, state_name, pop_est_2019  
3     FROM us_counties_pop_est_2019  
4     WHERE pop_est_2019 >= 100000  
5 )  
6 SELECT state_name, count(*) AS num_big_counties  
7 FROM big_counties  
8 GROUP BY state_name  
9 ORDER BY num_big_counties DESC;
```

6.2 Checkpoint: state totals with two CTEs

Task: Build a state-level summary (table: `us_counties_pop_est_2019`)

...

Required columns:

- `state_name`
- `total_pop` (sum of `pop_est_2019`)
- `num_counties` (count of rows per state)
- `avg_pop_per_county` (`total_pop / num_counties`, rounded to whole numbers is fine)

...

Requirements:

- **CTE 1:** aggregate population by `state_name`
- **CTE 2:** aggregate county counts by `state_name`
- Join the two CTEs on `state_name`
- Order by `avg_pop_per_county` descending

...

Expected shape (first rows):

state_name	total_pop	num_counties	avg_pop_per_county
California
Texas

...

Hint: Same GROUP BY state_name in both CTEs; join keys match exactly.

6.3 Solution: state totals with two CTEs

```
1 WITH state_pop AS (  
2     SELECT state_name, sum(pop_est_2019) AS total_pop  
3     FROM us_counties_pop_est_2019  
4     GROUP BY state_name  
5 ),  
6 state_counties AS (  
7     SELECT state_name, count(*) AS num_counties  
8     FROM us_counties_pop_est_2019  
9     GROUP BY state_name  
10 )  
11 SELECT sp.state_name,  
12        sp.total_pop,  
13        sc.num_counties,  
14        round(sp.total_pop::numeric / sc.num_counties, 0) AS avg_pop_per_county  
15 FROM state_pop AS sp  
16 JOIN state_counties AS sc USING (state_name)  
17 ORDER BY avg_pop_per_county DESC;
```

7 Crosstab

7.1 Enable and pivot

Reporting often needs **long to wide** layout. PostgreSQL provides `crosstab()` in the `tablefunc` extension.

...

```

1 CREATE EXTENSION IF NOT EXISTS tablefunc;
2
3 SELECT *
4 FROM crosstab(
5     $$
6     SELECT office, flavor, count(*) AS response_count
7     FROM ice_cream_survey
8     GROUP BY office, flavor
9     ORDER BY office
10    $$,
11    $$ SELECT flavor
12       FROM ice_cream_survey
13       GROUP BY flavor
14       ORDER BY flavor
15    $$
16 ) AS ct(
17     office text,
18     chocolate bigint,
19     strawberry bigint,
20     vanilla bigint
21 );

```

...

Contract: the first query must expose **three** columns: row id, category, value. The second query lists category labels. The outer AS (...) names output columns in that order.

8 CASE and classification

8.1 Pattern

CASE evaluates **first match wins**. Put the most specific conditions first when ranges overlap.

```

1 SELECT max_temp,
2     CASE
3         WHEN max_temp >= 90 THEN 'Hot'
4         WHEN max_temp >= 70 THEN 'Warm'
5         WHEN max_temp >= 50 THEN 'Pleasant'
6         ELSE 'Cool or missing'
7     END AS temp_band
8 FROM temperature_readings
9 LIMIT 5;

```

8.2 Checkpoint: population tiers with a CTE

Task: Classify each county, then summarize. (table: `us_counties_pop_est_2019`)

Required columns (final result):

- `tier`
- `num_counties`
- `total_pop`
- `avg_county_pop`

Requirements:

- Use a **CTE** that adds `tier` with `CASE`
- Outer query: `GROUP BY tier` with aggregates

Tier rules on `pop_est_2019`:

- Metro: ≥ 500000
- Urban: 100000 to 499999
- Suburban: 50000 to 99999
- Rural: below 50000

Expected shape:

<code>tier</code>	<code>num_counties</code>	<code>total_pop</code>	<code>avg_county_pop</code>
Rural
Metro

Hint: `CASE` should mirror the order above so a county falls into exactly one tier.

8.3 Solution: population tiers

```
1 WITH county_tiers AS (  
2     SELECT county_name,  
3            state_name,  
4            pop_est_2019,  
5            CASE  
6                WHEN pop_est_2019 >= 500000 THEN 'Metro'  
7                WHEN pop_est_2019 >= 100000 THEN 'Urban'  
8                WHEN pop_est_2019 >= 50000 THEN 'Suburban'  
9                ELSE 'Rural'  
10           END AS tier
```

```

11     FROM us_counties_pop_est_2019
12 )
13 SELECT tier,
14        count(*) AS num_counties,
15        sum(pop_est_2019) AS total_pop,
16        round(avg(pop_est_2019), 0) AS avg_county_pop
17 FROM county_tiers
18 GROUP BY tier
19 ORDER BY total_pop DESC;

```

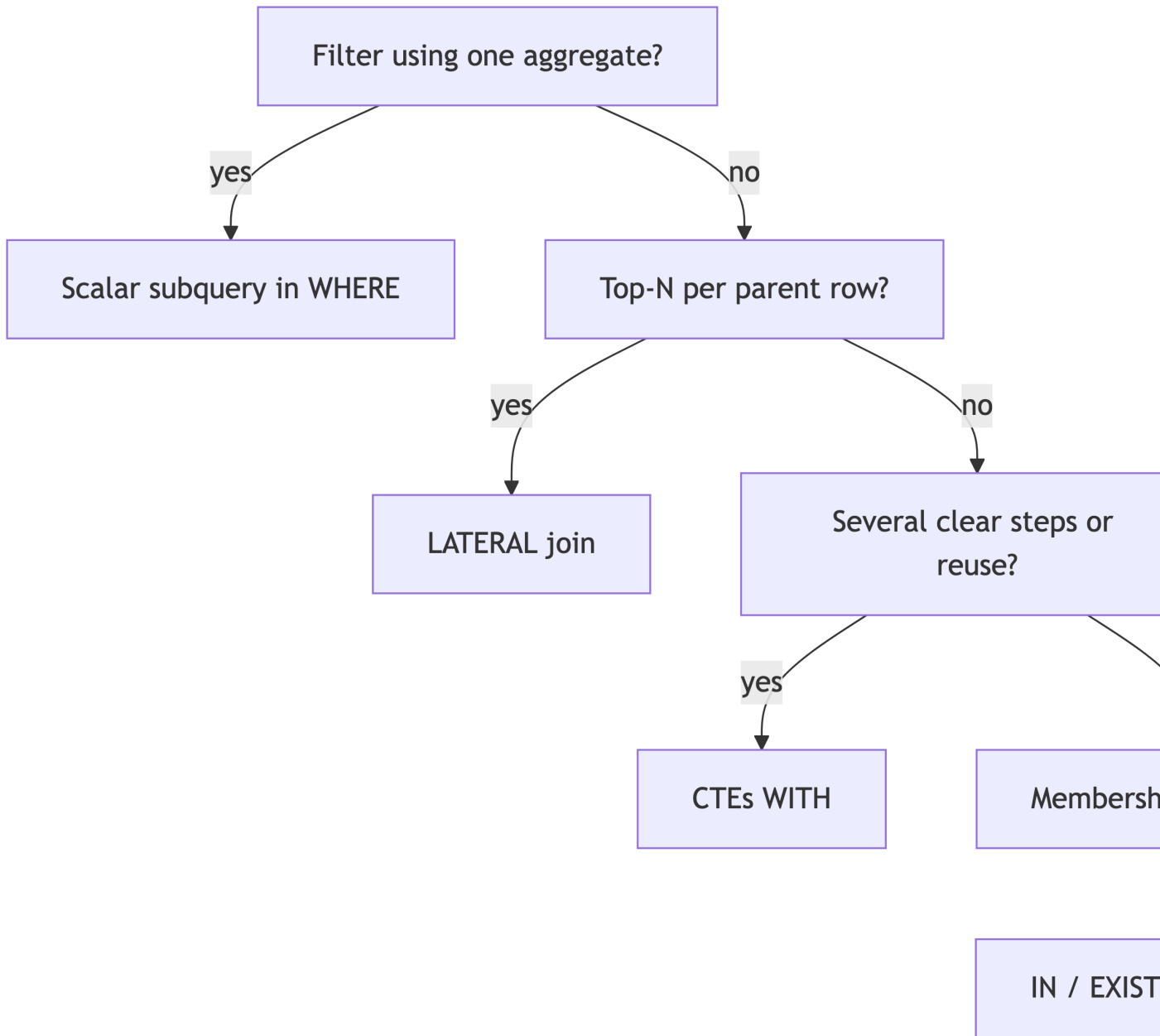
Typical story: many counties are rural, but a small number of metro counties hold a large share of total population.

9 Summary

9.1 Techniques in one place

Technique	Use it when
Scalar subquery in WHERE	Compare rows to one computed threshold
Derived table in FROM	Aggregate, then select or join the aggregate
IN	Simple membership against a set
EXISTS / NOT EXISTS	Correlated checks; anti-joins; nullable-safe negation
LATERAL	Per-row subquery, top-N per group
CTE (WITH)	Multiple named steps, less repetition
crosstab()	Pivot reporting shapes
CASE	Bucket values before or during aggregation

9.2 Choosing a pattern



9.3 References

1. DeBarros, A. (2022). *Practical SQL* (2nd ed.). No Starch Press. Chapter 13.

2. PostgreSQL: [WITH Queries \(CTEs\)](#)
3. PostgreSQL: [Subquery Expressions](#)
4. PostgreSQL: [LATERAL](#)
5. PostgreSQL: [tablefunc](#)