

Lecture 09-1: Advanced Statistics with Flight and Weather Data

DATA 503: Fundamentals of Data Engineering

Lucas P. Cordova, Ph.D.

2026-03-11

This lecture takes the flight and weather data we have been collecting for two weeks and transforms it from raw JSON blobs into a normalized relational schema. We then apply advanced SQL statistics – correlation, regression, ranking, and window functions – to answer the question: does weather affect air traffic over Portland? Based on Chapter 11 of Practical SQL, 2nd Edition.

Table of contents

1	Part 1: The Data Pipeline So Far	1
2	Part 2: Understanding Our JSON Data	3
3	Part 3: Designing the Normalized Schema	7
4	Part 4: The Transform Step	11
5	Part 5: Railway Transform Service	14
6	Part 6: Advanced Statistics in SQL	19
7	Part 7: Answering the Research Question	38

1 Part 1: The Data Pipeline So Far

For two weeks, your Railway scrapers have been quietly hoarding JSON like digital squirrels. Time to crack open the acorns.

1.1 Our Data Collection Architecture

1.2 What We Built

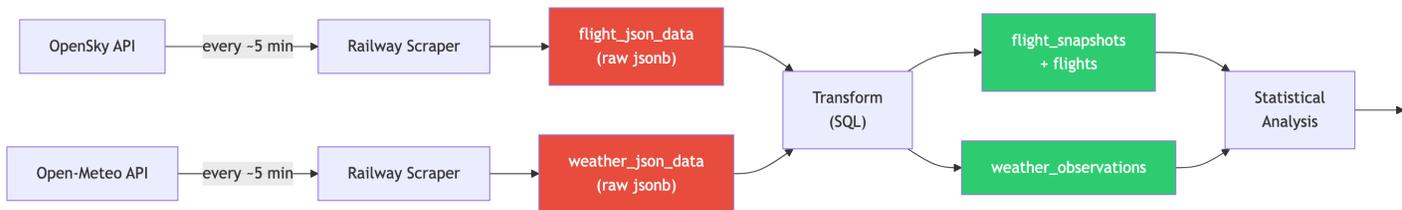
Over the past two weeks, you deployed two API scrapers on Railway:

Scraper	Source	Frequency	Rows Collected
Flight tracker	OpenSky Network API	~5 min	~2,130 snapshots
Weather monitor	Open-Meteo API	~5 min	~1,877 observations

Both scrapers dump raw JSON into Railway Postgres. No parsing, no cleaning, no opinions. Just raw data.

This is the **staging layer** of our pipeline, and today we promote that data to something useful.

1.3 The Pipeline: From JSON to Insight



Red = messy staging. Green = clean, normalized. Blue = the whole reason we did this.

1.4 The Staging Tables

Both staging tables have the same dead-simple structure:

```
1 CREATE TABLE flight_json_data (  
2     id serial PRIMARY KEY,  
3     raw_json jsonb,  
4     created_at TIMESTAMPTZ DEFAULT NOW()  
5 );  
6  
7 CREATE TABLE weather_json_data (  
8     id serial PRIMARY KEY,  
9     raw_json jsonb,
```

```

10     created_at TIMESTAMPTZ DEFAULT NOW()
11 );

```

One column for the ID, one for the entire JSON blob, one for when we grabbed it.

This is the database equivalent of throwing everything in a cardboard box during a move. It works, but you would not want to live like that.

2 Part 2: Understanding Our JSON Data

Before we can transform the data, we need to understand what we are actually looking at.

Time to open the cardboard boxes.

2.1 Exploring JSON in PostgreSQL

2.2 PostgreSQL JSON Operators: Your Swiss Army Knife

Operator	Returns	Example
->	JSON object/array	<code>raw_json->'states'</code> returns a JSON array
->>	Text	<code>raw_json->>'time'</code> returns '1772503243' as text
->N	Nth array element (JSON)	<code>state->0</code> returns first element as JSON
->>N	Nth array element (text)	<code>state->>1</code> returns <code>callsign</code> as text
<code>jsonb_array_elements()</code>	Set of JSON values	Unnests an array into rows

The difference between `->` and `->>` is critical: one gives you JSON, the other gives you text. Mix them up and your casts will fail in confusing ways.

2.3 Flight JSON Structure

Each row in `flight_json_data` contains a snapshot from the OpenSky Network API:

```

1 {
2   "time": 1772503243,
3   "states": [
4     [
5       "a88bb8",           -- [0] icao24 transponder address
6       "N65PT  ",         -- [1] callsign (8 chars, trailing spaces)
7       "United States",  -- [2] origin_country
8       1772503243,       -- [3] time_position
9       1772503243,       -- [4] last_contact
10      -122.23,           -- [5] longitude
11      45.5592,           -- [6] latitude
12      731.52,            -- [7] baro_altitude (meters, nullable)
13      false,             -- [8] on_ground
14      66.19,             -- [9] velocity (m/s)
15      78.34,             -- [10] true_track (degrees)
16      0.65,              -- [11] vertical_rate (m/s, nullable)
17      null,              -- [12] sensors (usually null)
18      754.38,            -- [13] geo_altitude (meters, nullable)
19      null,              -- [14] squawk (nullable)
20      false,             -- [15] spi
21      0                  -- [16] position_source
22    ]
23  ]
24 }

```

Notice: `states` is an array of arrays. Each inner array is one aircraft. And sometimes `states` is just `null` because apparently nobody was flying over Portland at 3 AM. Shocking.

2.4 Exploring Flight Data

Let's peek at what we have:

```

1 -- How many snapshots?
2 SELECT count(*) FROM flight_json_data;
3
4 -- Look at one snapshot's timestamp
5 SELECT raw_json->>'time' AS api_time,
6        created_at AS collected_at
7 FROM flight_json_data
8 LIMIT 1;
9
10 -- How many aircraft in each snapshot?

```

```

11 SELECT id,
12         jsonb_array_length(raw_json->'states') AS num_aircraft
13 FROM flight_json_data
14 WHERE raw_json->'states' IS NOT NULL
15        AND raw_json->>'states' != 'null'
16 ORDER BY num_aircraft DESC
17 LIMIT 10;

```

That `WHERE` clause filters out snapshots where nobody was flying. Without it, `jsonb_array_length` will throw errors on null values, and your query will be reduced to atoms.

2.5 Unnesting Flight Arrays

The real magic is `jsonb_array_elements()`, which turns one row with an array of 30 aircraft into 30 rows with one aircraft each:

```

1 SELECT
2     raw_json->>'time' AS snapshot_time,
3     state->>0 AS icao24,
4     TRIM(state->>1) AS callsign,
5     state->>2 AS origin_country,
6     (state->>5)::numeric AS longitude,
7     (state->>6)::numeric AS latitude,
8     (state->>7)::numeric AS baro_altitude,
9     (state->>8)::boolean AS on_ground,
10    (state->>9)::numeric AS velocity
11 FROM flight_json_data,
12      jsonb_array_elements(raw_json->'states') AS state
13 WHERE raw_json->'states' IS NOT NULL
14        AND raw_json->>'states' != 'null'
15 LIMIT 20;

```

Notice the `TRIM()` on `callsign`. OpenSky pads callsigns to 8 characters with trailing spaces, because apparently fixed-width strings never really died.

2.6 Weather JSON Structure

Weather data is simpler. Each row is one observation, no arrays to unnest:

```

1  {
2    "current": {
3      "time": "2026-03-03T04:30",
4      "interval": 900,
5      "weathercode": 1,
6      "temperature_2m": 11.5,
7      "wind_speed_10m": 8.0
8    },
9    "latitude": 45.528744,
10   "longitude": -122.696236,
11   "elevation": 31.0,
12   "timezone": "GMT"
13 }

```

Three numbers we care about: temperature, wind speed, and weather code.

2.7 Exploring Weather Data

```

1  -- Peek at weather observations
2  SELECT
3    raw_json->'current'->>'time' AS obs_time,
4    (raw_json->'current'->>'temperature_2m')::numeric AS temp_c,
5    (raw_json->'current'->>'wind_speed_10m')::numeric AS wind_kmh,
6    (raw_json->'current'->>'weathercode')::int AS weather_code
7  FROM weather_json_data
8  ORDER BY obs_time DESC
9  LIMIT 10;

```

Notice the nested access: `raw_json->'current'->>'temperature_2m'`. First we navigate into the `current` object with `->`, then we extract the value as text with `->>`, then we cast to `numeric`.

2.8 WMO Weather Codes Reference

Code	Meaning	Code	Meaning
0	Clear sky	51-55	Drizzle
1	Mainly clear	61-65	Rain
2	Partly cloudy	71-75	Snow
3	Overcast	80-82	Rain showers
45	Fog	95	Thunderstorm

Portland in March means you will see a lot of codes 1-3 and 51-65. If you see a 0, screenshot it. That is a rare event.

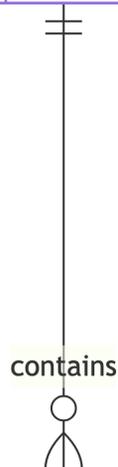
3 Part 3: Designing the Normalized Schema

Time to give this data a proper home. One where columns have names and types, not just array indices you have to memorize.

3.1 The Target Schema

3.2 Entity-Relationship Diagram

flight_snapshots		
serial	snapshot_id	PK
timestampz	snapshot_time	
timestampz	collected_at	



flights		
serial	flight_id	PK
int	snapshot_id	FK
varchar	icao24	
varchar	callsign	
varchar	origin_country	
numeric	longitude	
numeric	latitude	

weather_observations	
serial	observation_id
timestampz	observation_time
timestampz	collected_at
numeric	temperature_c
numeric	wind_speed_kmh
smallint	weather_code
numeric	latitude
numeric	longitude
numeric	elevation

3.3 Why Two Flight Tables?

Each API call returns one timestamp and many aircraft. That is a classic one-to-many relationship:

- **flight_snapshots**: one row per API call (the “when”)
- **flights**: one row per aircraft per snapshot (the “what”)

If we jammed everything into one table, we would duplicate the snapshot timestamp for every aircraft. That is denormalized, and we just spent three weeks learning why that is bad.

3.4 Design Decisions

flight_snapshots

- `snapshot_time` comes from the JSON (`raw_json->>'time'`), converted from Unix epoch
- `collected_at` comes from the staging table’s `created_at` column
- These are usually close but not identical. The API timestamp tells us when OpenSky gathered the data; our timestamp tells us when we grabbed it.

flights

- Each field maps to an array index (0-16) in the state vector
- Callsigns are trimmed of trailing whitespace
- Nullable fields (`baro_altitude`, `vertical_rate`, `geo_altitude`, `squawk`) stay nullable

weather_observations

- Flat table since each observation is a single data point, no arrays to normalize
- `observation_time` parsed from the ISO timestamp in the JSON
- We keep `latitude`, `longitude`, `elevation` even though they are constant, because good schema design does not assume things stay constant

3.5 The DDL

```
1 CREATE TABLE flight_snapshots (  
2     snapshot_id serial PRIMARY KEY,  
3     snapshot_time timestamptz NOT NULL,  
4     collected_at timestamptz NOT NULL  
5 );  
6  
7 CREATE TABLE flights (  
8     flight_id serial PRIMARY KEY,
```

```

9     snapshot_id int NOT NULL REFERENCES flight_snapshots(snapshot_id),
10    icao24 varchar(10),
11    callsign varchar(10),
12    origin_country varchar(100),
13    longitude numeric(9,4),
14    latitude numeric(8,4),
15    baro_altitude numeric(8,2),
16    on_ground boolean,
17    velocity numeric(7,2),
18    true_track numeric(6,2),
19    vertical_rate numeric(6,2),
20    geo_altitude numeric(8,2),
21    squawk varchar(10),
22    spi boolean,
23    position_source smallint
24 );
25
26 CREATE TABLE weather_observations (
27     observation_id serial PRIMARY KEY,
28     observation_time timestamptz NOT NULL,
29     collected_at timestamptz NOT NULL,
30     temperature_c numeric(5,2),
31     wind_speed_kmh numeric(6,2),
32     weather_code smallint,
33     latitude numeric(8,4),
34     longitude numeric(9,4),
35     elevation numeric(6,1)
36 );

```

Your data finally looks like it was entered by someone who cares.

4 Part 4: The Transform Step

This is where the magic happens. We take the JSON blobs and extract structured, typed, normalized rows.

Think of it as the database equivalent of unpacking after a move and actually putting things in drawers.

4.1 Migrating JSON to Relational Tables

4.2 Step 1: Populate flight_snapshots

Each row in `flight_json_data` becomes one row in `flight_snapshots`. We convert the Unix timestamp to a proper `timestamptz`:

```
1 INSERT INTO flight_snapshots (snapshot_time, collected_at)
2 SELECT DISTINCT
3     to_timestamp((raw_json->>'time')::bigint) AS snapshot_time,
4     created_at AS collected_at
5 FROM flight_json_data
6 WHERE raw_json->>'time' IS NOT NULL;
```

The `DISTINCT` handles any duplicate snapshots. The `to_timestamp()` function converts Unix epoch seconds into a PostgreSQL timestamp. No more staring at 1772503243 and wondering what day that is.

4.3 Step 2: Populate flights

This is the big one. We unnest the state arrays and join back to `flight_snapshots` to get the foreign key:

```
1 INSERT INTO flights (
2     snapshot_id, icao24, callsign, origin_country,
3     longitude, latitude, baro_altitude, on_ground,
4     velocity, true_track, vertical_rate, geo_altitude,
5     squawk, spi, position_source
6 )
7 SELECT
8     fs.snapshot_id,
9     state->>0 AS icao24,
10    TRIM(state->>1) AS callsign,
11    state->>2 AS origin_country,
12    (state->>5)::numeric AS longitude,
13    (state->>6)::numeric AS latitude,
14    (state->>7)::numeric AS baro_altitude,
15    (state->>8)::boolean AS on_ground,
16    (state->>9)::numeric AS velocity,
17    (state->>10)::numeric AS true_track,
18    (state->>11)::numeric AS vertical_rate,
19    (state->>13)::numeric AS geo_altitude,
20    state->>14 AS squawk,
```

```

21     (state->>15)::boolean AS spi,
22     (state->>16)::smallint AS position_source
23 FROM flight_json_data fjd
24 JOIN flight_snapshots fs
25     ON to_timestamp((fjd.raw_json->>'time')::bigint) = fs.snapshot_time
26 CROSS JOIN jsonb_array_elements(fjd.raw_json->'states') AS state
27 WHERE fjd.raw_json->'states' IS NOT NULL
28     AND fjd.raw_json->>'states' != 'null';

```

Let's break down what is happening:

- CROSS JOIN jsonb_array_elements(...) unnests each aircraft array into its own row
- We join back to flight_snapshots to get the snapshot_id foreign key
- The WHERE clause filters out snapshots with no aircraft
- Array indices 0-16 map to the state vector positions we documented earlier
- Notice index 12 (sensors) is skipped since it is almost always null and not useful for analysis

4.4 Step 3: Populate weather_observations

Weather is straightforward since there is no array to unnest:

```

1  INSERT INTO weather_observations (
2      observation_time, collected_at,
3      temperature_c, wind_speed_kmh, weather_code,
4      latitude, longitude, elevation
5  )
6  SELECT
7      (raw_json->'current'->>'time')::timestampz AS observation_time,
8      created_at AS collected_at,
9      (raw_json->'current'->>'temperature_2m')::numeric AS temperature_c,
10     (raw_json->'current'->>'wind_speed_10m')::numeric AS wind_speed_kmh,
11     (raw_json->'current'->>'weathercode')::smallint AS weather_code,
12     (raw_json->>'latitude')::numeric AS latitude,
13     (raw_json->>'longitude')::numeric AS longitude,
14     (raw_json->>'elevation')::numeric AS elevation
15 FROM weather_json_data
16 WHERE raw_json->'current' IS NOT NULL;

```

4.5 Verification Queries

Always verify your transforms. Trust, but verify. Actually, just verify.

```

1  -- How many snapshots?
2  SELECT count(*) FROM flight_snapshots;
3
4  -- How many individual flight records?
5  SELECT count(*) FROM flights;
6
7  -- How many weather observations?
8  SELECT count(*) FROM weather_observations;
9
10 -- Spot check: flights per snapshot
11 SELECT fs.snapshot_id,
12        fs.snapshot_time,
13        count(f.flight_id) AS num_flights
14 FROM flight_snapshots fs
15 LEFT JOIN flights f ON fs.snapshot_id = f.snapshot_id
16 GROUP BY fs.snapshot_id, fs.snapshot_time
17 ORDER BY num_flights DESC
18 LIMIT 10;
19
20 -- Spot check: weather range
21 SELECT
22     min(temperature_c) AS min_temp,
23     max(temperature_c) AS max_temp,
24     min(wind_speed_kmh) AS min_wind,
25     max(wind_speed_kmh) AS max_wind,
26     min(observation_time) AS earliest,
27     max(observation_time) AS latest
28 FROM weather_observations;

```

If your flight count is zero, go back and check your `WHERE` clause. If your temperature range includes 500 degrees Celsius, something went very wrong with your casting.

5 Part 5: Railway Transform Service

You do not want to run these transforms by hand every time new data comes in. Let's automate it.

5.1 Deploying Transforms on Railway

5.2 The db_transform Template

Clone the template repo:

```
https://github.com/LucasCordova/db_transform
```

This repo contains a simple setup:

- A `Dockerfile` that runs a SQL file against your Postgres database
- A `clean.sql` file where you put your transform queries
- That is it. No frameworks, no dependencies, no opinions.

You put your SQL in `clean.sql`, deploy to Railway, and it runs.

5.3 Step-by-Step Deployment

1. **Fork or clone** the `db_transform` repo to your own GitHub account
2. **Edit `clean.sql`** with the transform queries from Part 4 (insert-select, etc.)
3. **Commit and push** to your GitHub repo
4. **In Railway dashboard**, click “New” then “**GitHub Repo**”
5. **Select your repo** from the list
6. **Add the `DATABASE_URL` environment variable** pointing to your Railway Postgres instance (find this in your Postgres service’s Variables tab)
7. **Deploy** and the service runs `clean.sql` against your database
8. **Optional:** Set up as a **cron job** in Railway for recurring transforms

5.4 What Does “Idempotent” Mean?

An operation is **idempotent** if running it once or a hundred times produces the same result. The first run changes things. Every run after that changes nothing.

Statement	Idempotent?	Why
CREATE TABLE IF NOT EXISTS	Yes	First run creates it, subsequent runs skip it
CREATE TABLE	No	Second run errors because the table already exists
TRUNCATE then INSERT	Yes	Every run wipes and rebuilds, same result every time
Just INSERT	No	Every run adds more duplicate rows

This matters because your cron transform runs repeatedly. If your `clean.sql` is not idempotent, you get duplicate data piling up every cycle. If it is idempotent, it does not matter if it runs once or fifty times. Same clean dataset every time.

5.5 Important Notes on the Transform Service

- The service runs once and exits. It is not a long-running server.
- If you need to re-run, trigger a new deployment or use Railway's cron feature.
- Make your SQL **idempotent**. Since the cron runs repeatedly against a database that keeps accumulating new staging data, you need to avoid inserting duplicate rows every time.
- The `DATABASE_URL` environment variable is automatically available inside the container.

5.6 Important Notes on the Transform Service (continued)

The key pattern: **truncate your target tables before re-inserting**. Your staging tables (`flight_json_data`, `weather_json_data`) keep growing as the API scrapers add rows. Each cron run should wipe the transformed tables and rebuild them from the full staging dataset.

```

1 -- =====
2 -- clean.sql - Idempotent transform for Railway cron
3 -- This runs on every cron trigger. It rebuilds the

```

```

4  -- normalized tables from scratch using the staging data.
5  -- =====
6
7  -- Step 1: Create tables if they do not exist yet (first run)
8  CREATE TABLE IF NOT EXISTS flight_snapshots (
9      snapshot_id serial PRIMARY KEY,
10     snapshot_time timestamptz NOT NULL,
11     collected_at timestamptz NOT NULL
12 );
13
14 CREATE TABLE IF NOT EXISTS flights (
15     flight_id serial PRIMARY KEY,
16     snapshot_id integer REFERENCES flight_snapshots (snapshot_id),
17     icao24 varchar(10),
18     callsign varchar(10),
19     origin_country varchar(100),
20     longitude numeric(9,4),
21     latitude numeric(8,4),
22     baro_altitude numeric(8,2),
23     on_ground boolean,
24     velocity numeric(7,2),
25     true_track numeric(6,2),
26     vertical_rate numeric(6,2),
27     geo_altitude numeric(8,2),
28     squawk varchar(10),
29     spi boolean,
30     position_source smallint
31 );
32
33 CREATE TABLE IF NOT EXISTS weather_observations (
34     observation_id serial PRIMARY KEY,
35     observation_time timestamptz NOT NULL,
36     collected_at timestamptz NOT NULL,
37     temperature_c numeric(5,2),
38     wind_speed_kmh numeric(6,2),
39     weather_code smallint,
40     latitude numeric(8,4),
41     longitude numeric(9,4),
42     elevation numeric(6,1)
43 );
44
45 -- Step 2: Truncate in dependency order (children first)
46 -- RESTART IDENTITY resets the serial counters

```

```

47 -- CASCADE is not needed here because we truncate children first,
48 -- but it is a safety net
49 TRUNCATE flights RESTART IDENTITY;
50 TRUNCATE flight_snapshots RESTART IDENTITY CASCADE;
51 TRUNCATE weather_observations RESTART IDENTITY;
52
53 -- Step 3: Re-populate from staging data
54 -- (same INSERT INTO ... SELECT queries from the transform step)
55
56 -- Weather observations
57 INSERT INTO weather_observations (
58     observation_time, collected_at, temperature_c,
59     wind_speed_kmh, weather_code,
60     latitude, longitude, elevation
61 )
62 SELECT DISTINCT
63     (raw_json->'current'->>'time')::timestamptz,
64     created_at,
65     (raw_json->'current'->>'temperature_2m')::numeric,
66     (raw_json->'current'->>'wind_speed_10m')::numeric,
67     (raw_json->'current'->>'weathercode')::smallint,
68     (raw_json->>'latitude')::numeric,
69     (raw_json->>'longitude')::numeric,
70     (raw_json->>'elevation')::numeric
71 FROM weather_json_data;
72
73 -- Flight snapshots (only rows where states is not null)
74 INSERT INTO flight_snapshots (snapshot_time, collected_at)
75 SELECT DISTINCT
76     to_timestamp((raw_json->>'time')::bigint),
77     created_at
78 FROM flight_json_data
79 WHERE raw_json->'states' IS NOT NULL
80     AND raw_json->>'states' != 'null';
81
82 -- Flights (unnest the states array)
83 INSERT INTO flights (
84     snapshot_id, icao24, callsign, origin_country,
85     longitude, latitude, baro_altitude, on_ground,
86     velocity, true_track, vertical_rate,
87     geo_altitude, squawk, spi, position_source
88 )
89 SELECT

```

```

90     fs.snapshot_id,
91     trim(state->>0),
92     trim(state->>1),
93     state->>2,
94     (state->>5)::numeric,
95     (state->>6)::numeric,
96     (state->>7)::numeric,
97     (state->>8)::boolean,
98     (state->>9)::numeric,
99     (state->>10)::numeric,
100    (state->>11)::numeric,
101    (state->>13)::numeric,
102    state->>14,
103    (state->>15)::boolean,
104    (state->>16)::smallint
105 FROM flight_json_data f
106 JOIN flight_snapshots fs
107     ON to_timestamp((f.raw_json->>'time')::bigint) = fs.snapshot_time
108 CROSS JOIN jsonb_array_elements(f.raw_json->'states') AS state
109 WHERE f.raw_json->'states' IS NOT NULL
110     AND f.raw_json->>'states' != 'null';

```

5.7 Important Notes on the Transform Service (continued)

Every time the cron fires, this wipes the normalized tables and rebuilds them from whatever is currently in the staging tables. New API data that arrived since the last run gets included automatically. This is the nuclear-but-reliable approach: no partial states, no duplicate rows, no “did that row already get transformed?” headaches.

For large datasets where a full rebuild is too slow, you would instead track a high-water mark (e.g., the max `id` or `created_at` you last processed) and only transform new rows. But for our dataset size, the full truncate-and-reload is perfectly fine and much simpler to reason about.

6 Part 6: Advanced Statistics in SQL

Now the fun part. We have clean, normalized data. Time to interrogate it.

Research Question: Does weather affect the volume and behavior of air traffic over Portland?

6.1 Descriptive Statistics

6.2 Getting Our Bearings

Before we correlate anything, we need to know what “normal” looks like. Descriptive statistics give us the baseline.

```
1 -- Flight counts per snapshot
2 SELECT
3     count(*) AS num_snapshots,
4     round(avg(flight_count), 2) AS avg_flights,
5     round(stddev_pop(flight_count), 2) AS stddev_flights,
6     min(flight_count) AS min_flights,
7     max(flight_count) AS max_flights
8 FROM (
9     SELECT fs.snapshot_id,
10          count(f.flight_id) AS flight_count
11     FROM flight_snapshots fs
12     LEFT JOIN flights f ON fs.snapshot_id = f.snapshot_id
13     GROUP BY fs.snapshot_id
14 ) sub;
```

This tells us the average number of aircraft visible over Portland at any given moment, plus how much that varies.

6.3 Weather Descriptive Stats

```
1 -- Temperature and wind statistics
2 SELECT
3     count(*) AS num_observations,
4     round(avg(temperature_c), 2) AS avg_temp_c,
5     round(stddev_pop(temperature_c), 2) AS stddev_temp,
6     round(min(temperature_c), 2) AS min_temp,
7     round(max(temperature_c), 2) AS max_temp,
8     round(avg(wind_speed_kmh), 2) AS avg_wind,
9     round(stddev_pop(wind_speed_kmh), 2) AS stddev_wind,
10    round(min(wind_speed_kmh), 2) AS min_wind,
11    round(max(wind_speed_kmh), 2) AS max_wind
12 FROM weather_observations;
```

6.4 Flight Altitude Statistics

```
1  -- Altitude statistics for airborne aircraft
2  SELECT
3      count(*) AS num_airborne,
4      round(avg(baro_altitude), 2) AS avg_altitude_m,
5      round(stddev_pop(baro_altitude), 2) AS stddev_altitude,
6      round(min(baro_altitude), 2) AS min_altitude,
7      round(max(baro_altitude), 2) AS max_altitude,
8      round(avg(velocity), 2) AS avg_speed_ms,
9      round(stddev_pop(velocity), 2) AS stddev_speed
10 FROM flights
11 WHERE on_ground = false
12      AND baro_altitude IS NOT NULL;
```

If the standard deviation of altitude is enormous, that just means we are seeing everything from small Cessnas at 500 meters to commercial jets at 10,000 meters. Portland airspace is busy.

6.5 Views

6.6 What is a View?

A **view** is a saved SQL query that acts like a virtual table. It does not store data – it is just a named SELECT statement that lives in the database.

```
1  CREATE VIEW active_flights AS
2  SELECT callsign, origin_country, baro_altitude
3  FROM flights
4  WHERE on_ground = false;
```

Every time you query a view, PostgreSQL runs the underlying query fresh against the actual tables.

```
1  -- This works just like querying a table
2  SELECT * FROM active_flights;
```

No duplicate data. No extra storage. Just a query with a name.

6.7 Why Use Views?

Simplify complex queries. A gnarly 15-line JOIN with filters becomes `SELECT * FROM that_view`. The complexity lives in one place.

Abstraction layer. If the underlying table structure changes, update the view definition. Downstream queries do not break.

Security and access control. Grant someone access to a view that shows only certain columns or rows, without exposing the full table. Think: employees seeing their own records but not everyone else's salaries.

Reusability. Same complex query used in 5 different reports? One view, five simple queries.

6.8 When to Use Views (and When Not To)

Use views when:

- You are writing the same JOIN/filter combo repeatedly
- You need to restrict what users can see
- You want meaningful names for derived datasets (`monthly_revenue` beats a sprawling subquery)
- You are building a reporting layer on top of a normalized schema

6.9 When to Use Views (and When Not To)

Do not use views when:

- You need cached/precomputed results (that is a **materialized view**, a different thing)
- The query is a simple one-off you will run once
- You are nesting views on views on views (gets hard to debug fast)

6.10 Creating the Hourly Analysis View

For meaningful statistics, we need to aggregate to a common time grain. Five-minute snapshots are too noisy. Let's use hourly buckets:

```
1 -- Hourly flight counts
2 CREATE OR REPLACE VIEW hourly_flights AS
3 SELECT
4     date_trunc('hour', fs.snapshot_time) AS hour,
5     count(DISTINCT fs.snapshot_id) AS num_snapshots,
6     count(f.flight_id) AS total_flights,
```

```

7     round(count(f.flight_id)::numeric
8         / NULLIF(count(DISTINCT fs.snapshot_id), 0), 2) AS avg_flights_per_snapshot
9 FROM flight_snapshots fs
10 LEFT JOIN flights f ON fs.snapshot_id = f.snapshot_id
11 GROUP BY date_trunc('hour', fs.snapshot_time);
12
13 -- Hourly weather (average across observations in each hour)
14 CREATE OR REPLACE VIEW hourly_weather AS
15 SELECT
16     date_trunc('hour', observation_time) AS hour,
17     round(avg(temperature_c), 2) AS avg_temp_c,
18     round(avg(wind_speed_kmh), 2) AS avg_wind_kmh,
19     mode() WITHIN GROUP (ORDER BY weather_code) AS predominant_weather_code
20 FROM weather_observations
21 GROUP BY date_trunc('hour', observation_time);
22
23 -- Combined hourly view for correlation analysis
24 CREATE OR REPLACE VIEW hourly_flight_weather AS
25 SELECT
26     hf.hour,
27     hf.num_snapshots,
28     hf.total_flights,
29     hf.avg_flights_per_snapshot,
30     hw.avg_temp_c,
31     hw.avg_wind_kmh,
32     hw.predominant_weather_code
33 FROM hourly_flights hf
34 JOIN hourly_weather hw ON hf.hour = hw.hour;

```

These views align our flight and weather data on the same hourly grid. This is how you join data that was collected independently at different intervals. Round to a common grain, then join on it.

6.11 Correlation

6.12 `corr(Y, X)`: Measuring Linear Relationships

The `corr()` function returns the Pearson correlation coefficient, a value between -1 and +1:

Value	Meaning
+1.0	Perfect positive correlation
0.0	No linear relationship

Value	Meaning
-1.0	Perfect negative correlation

Anything above 0.7 or below -0.7 is strong. Between 0.3 and 0.7 is moderate. Below 0.3 is weak or nonexistent.

6.13 Temperature vs. Flight Volume

```

1  -- Does temperature correlate with number of flights?
2  SELECT
3      round(corr(total_flights, avg_temp_c)::numeric, 4) AS temp_flight_corr,
4      round(corr(avg_flights_per_snapshot, avg_temp_c)::numeric, 4)
5          AS temp_avg_flight_corr
6  FROM hourly_flight_weather;
```

If this is positive, warmer hours tend to have more flights. Makes intuitive sense: better weather, more VFR (visual flight rules) traffic, more general aviation.

6.14 Wind Speed vs. Flight Behavior

```

1  -- Does wind speed affect average altitude or ground speed?
2  SELECT
3      round(corr(avg_alt, avg_wind_kmh)::numeric, 4) AS wind_altitude_corr,
4      round(corr(avg_speed, avg_wind_kmh)::numeric, 4) AS wind_speed_corr
5  FROM (
6      SELECT
7          hw.hour,
8          hw.avg_wind_kmh,
9          avg(f.baro_altitude) AS avg_alt,
10         avg(f.velocity) AS avg_speed
11     FROM hourly_weather hw
12     JOIN flight_snapshots fs
13         ON date_trunc('hour', fs.snapshot_time) = hw.hour
14     JOIN flights f ON fs.snapshot_id = f.snapshot_id
15     WHERE f.on_ground = false
16         AND f.baro_altitude IS NOT NULL
17     GROUP BY hw.hour, hw.avg_wind_kmh
18 ) sub;
```

A positive correlation between wind and altitude could mean pilots climb higher to find smoother air. Or it could mean nothing. We will get to that.

6.15 Wind Speed vs. Flight Volume

```
1 -- Does wind reduce the number of flights?
2 SELECT
3     round(corr(total_flights, avg_wind_kmh)::numeric, 4) AS wind_flight_corr
4 FROM hourly_flight_weather;
```

If this is negative, windier conditions correlate with fewer flights. Small aircraft stay grounded in high winds, so we might see this in the data.

6.16 Regression

6.17 Linear Regression in SQL

Correlation tells you there is a relationship. Regression tells you the equation.

Function	Returns
<code>regr_slope(Y, X)</code>	Slope of the regression line
<code>regr_intercept(Y, X)</code>	Y-intercept
<code>regr_r2(Y, X)</code>	R-squared (how well the line fits)

The formula: $Y = \text{slope} * X + \text{intercept}$

R-squared ranges from 0 to 1. An R-squared of 0.8 means the model explains 80% of the variance. An R-squared of 0.05 means your model explains almost nothing and you should probably find a better predictor.

6.18 Predicting Flight Volume from Temperature

```
1 SELECT
2     round(regr_slope(total_flights, avg_temp_c)::numeric, 4)
3     AS slope,
4     round(regr_intercept(total_flights, avg_temp_c)::numeric, 4)
5     AS intercept,
6     round(regr_r2(total_flights, avg_temp_c)::numeric, 4)
7     AS r_squared
8 FROM hourly_flight_weather;
```

Read the slope as: “For each 1 degree C increase in temperature, we expect X more flights per hour.”

The R-squared tells us how much of the variation in flight volume is explained by temperature alone. Spoiler: probably not a lot, because many other factors affect air traffic.

6.19 Predicting Flight Volume from Wind Speed

```
1 SELECT
2     round(regr_slope(total_flights, avg_wind_kmh)::numeric, 4)
3     AS slope,
4     round(regr_intercept(total_flights, avg_wind_kmh)::numeric, 4)
5     AS intercept,
6     round(regr_r2(total_flights, avg_wind_kmh)::numeric, 4)
7     AS r_squared
8 FROM hourly_flight_weather;
```

If the slope is negative and R-squared is low, wind matters a little but is far from the whole story. Time of day, day of week, and airline schedules all matter more than a light breeze.

6.20 Window Functions

6.21 The Problem with GROUP BY

GROUP BY is powerful, but it has a fundamental limitation: it **collapses rows**. You get one output row per group, and you lose access to individual rows.

```
1 -- This gives us one row per hour
2 SELECT hour, count(*) AS num_flights
3 FROM flights
4 GROUP BY hour;
```

What if you want the count **and** still see every individual flight? You cannot. GROUP BY forces a choice: aggregated summary or individual detail. Never both.

6.22 Window Functions: The Best of Both Worlds

A **window function** performs a calculation across a set of rows that are somehow related to the current row – without collapsing them.

```

1 SELECT
2     callsign,
3     origin_country,
4     baro_altitude,
5     avg(baro_altitude) OVER () AS overall_avg_altitude
6 FROM flights
7 WHERE on_ground = false;

```

Every row keeps its individual data. But each row also gets the overall average altitude attached to it. No grouping. No collapsing. Every row survives.

The `OVER()` clause is what makes a function a window function. It defines the “window” of rows the function looks at.

6.23 How Window Functions Work

Think of it as a sliding lens over your result set:

1. PostgreSQL computes the full result set first
2. For each row, it looks through the “window” defined by `OVER()`
3. It computes the function across that window
4. It attaches the result to the current row
5. It moves to the next row and repeats

The key insight: **the original rows are never collapsed**. You get computation across rows while keeping every row intact.

6.24 The `OVER()` Clause

`OVER()` controls what the window function can see:

Clause	What It Does	Example
<code>OVER ()</code>	Entire result set is the window	Global average across all rows
<code>OVER (ORDER BY col)</code>	Running calculation in that order	Cumulative sum, ranking
<code>OVER (PARTITION BY col)</code>	Separate window per group	Average per department
<code>OVER (PARTITION BY a ORDER BY b)</code>	Ordered calculation within groups	Rank within each category

```

1 -- Global: one window for everything
2 sum(total_flights) OVER ()
3
4 -- Partitioned: separate window per weather code
5 sum(total_flights) OVER (PARTITION BY weather_code)
6
7 -- Ordered: running total in time order
8 sum(total_flights) OVER (ORDER BY hour)
9
10 -- Both: running total within each weather code
11 sum(total_flights) OVER (
12     PARTITION BY weather_code ORDER BY hour
13 )

```

6.25 Window Functions vs GROUP BY

	GROUP BY	Window Function
Rows	Collapsed into groups	Preserved individually
Output	One row per group	Same number of rows as input
Access	Only grouped/aggregated columns	All columns plus computed values
Use case	Summaries and reports	Analysis alongside detail

You will often use both in the same query: **GROUP BY** to aggregate to the grain you need, then window functions on top to rank or compare those aggregates.

6.26 Types of Window Functions

Aggregate window functions – familiar aggregates with `OVER()`:

- `sum()`, `avg()`, `count()`, `min()`, `max()`

Ranking functions – assign positions:

- `rank()`, `dense_rank()`, `row_number()`, `ntile()`

Value functions – access other rows:

- `lag()` (previous row), `lead()` (next row)
- `first_value()`, `last_value()`, `nth_value()`

All of them use `OVER()`. That is the unifying syntax. If you see `OVER()`, you are looking at a window function.

6.27 Ranking

6.28 Why Rank Data?

Sorting gives you order. Ranking gives you **position**.

Sometimes you do not just want to see rows ordered by value – you want to know *where each row stands* relative to the others. That is a fundamentally different question.

- “What were the busiest hours?” – sorting works fine
- “Was this hour in the **top 10** busiest?” – you need a rank
- “Was this the busiest hour **for a Tuesday**?” – you need a rank within a group

Ranking is how you go from “show me the data” to “show me how the data compares to itself.”

6.29 When to Use Ranking

Leaderboards and top-N analysis. Find the top 5 airports, the bottom 10 performing stores, the highest-scoring students. You could use `ORDER BY . . . LIMIT`, but ranking lets you keep all the data while labeling positions.

Comparisons within groups. “Rank each employee’s sales within their department.” This is `PARTITION BY department ORDER BY sales` – every department gets its own ranking starting at 1.

Handling ties explicitly. When two rows have the same value, do you skip the next number or not? `ORDER BY` does not answer this. Ranking functions do.

Filtering by position. Wrap a ranking query in a subquery or CTE, then filter: `WHERE rank <= 3`. You cannot do this with plain `ORDER BY` because rank is computed, not stored.

6.30 The Three Ranking Functions

Function	Behavior on Ties	Example (4 rows, tie at #2)
<code>rank()</code>	Same rank for ties, skip next	1, 2, 2, 4
<code>dense_rank()</code>	Same rank for ties, no skip	1, 2, 2, 3
<code>row_number()</code>	No ties ever, arbitrary tiebreak	1, 2, 3, 4

Use `rank()` when gaps matter (sports standings, competition results).

Use `dense_rank()` when you want consecutive numbers (top-3 means exactly 3 distinct ranks, even with ties).

Use `row_number()` when you need unique row identifiers or want to pick one arbitrary winner per group.

All three use the `OVER(ORDER BY ...)` clause. Add `PARTITION BY` to rank within groups.

6.31 The `OVER()` Clause

The `OVER()` clause is what makes these **window functions**. It defines:

- **ORDER BY** – what determines rank position
- **PARTITION BY** (optional) – restart ranking for each group

```
1 -- Global ranking: one list for all rows
2 rank() OVER (ORDER BY total_flights DESC)
3
4 -- Partitioned ranking: separate list per weather condition
5 rank() OVER (
6     PARTITION BY weather_code
7     ORDER BY total_flights DESC
8 )
```

Without `PARTITION BY`, every row competes against every other row. With it, each group gets its own independent ranking. Think of it as running a separate `ORDER BY` for each group, then stitching the results back together.

6.32 Ranking Hours by Flight Volume

```
1 -- Which hours had the most air traffic?
2 SELECT
3     hour,
4     total_flights,
5     avg_temp_c,
6     avg_wind_kmh,
7     rank() OVER (ORDER BY total_flights DESC) AS volume_rank
8 FROM hourly_flight_weather
9 ORDER BY volume_rank
10 LIMIT 15;
```

The busiest hours will likely be mid-afternoon on weekdays. That is when commercial traffic peaks and general aviation is most active.

6.33 Ranking by Altitude Within Weather Conditions

```
1  -- Rank hours by average altitude, partitioned by weather code
2  SELECT
3      hw.hour,
4      hw.predominant_weather_code,
5      round(avg(f.baro_altitude), 2) AS avg_altitude,
6      dense_rank() OVER (
7          PARTITION BY hw.predominant_weather_code
8          ORDER BY avg(f.baro_altitude) DESC
9      ) AS altitude_rank
10 FROM hourly_weather hw
11 JOIN flight_snapshots fs
12     ON date_trunc('hour', fs.snapshot_time) = hw.hour
13 JOIN flights f ON fs.snapshot_id = f.snapshot_id
14 WHERE f.on_ground = false
15        AND f.baro_altitude IS NOT NULL
16 GROUP BY hw.hour, hw.predominant_weather_code
17 ORDER BY hw.predominant_weather_code, altitude_rank
18 LIMIT 20;
```

PARTITION BY weather_code restarts the ranking for each weather condition. This lets us compare altitude rankings within clear weather separately from overcast conditions.

6.34 Top Flight Hours by Day of Week

```
1  -- Rank flight volume within each day of week
2  SELECT
3      to_char(hour, 'Day') AS day_of_week,
4      extract(hour FROM hour) AS hour_of_day,
5      total_flights,
6      rank() OVER (
7          PARTITION BY extract(dow FROM hour)
8          ORDER BY total_flights DESC
9      ) AS rank_within_day
10 FROM hourly_flight_weather
11 ORDER BY extract(dow FROM hour), rank_within_day
12 LIMIT 20;
```

6.35 Rolling Averages

6.36 Window Functions with Frame Clauses

Rolling averages smooth out noisy data so you can see trends. The syntax:

```
1 avg(column) OVER (  
2     ORDER BY time_column  
3     ROWS BETWEEN 3 PRECEDING AND 3 FOLLOWING  
4 )
```

This computes the average of the current row plus 3 rows before and 3 rows after, giving a 7-row moving average. Change the numbers to adjust the smoothing window.

6.37 Rolling Average of Flight Volume

```
1 -- 6-hour rolling average of flights (3 hours before, 3 after)  
2 SELECT  
3     hour,  
4     total_flights,  
5     round(avg(total_flights) OVER (  
6         ORDER BY hour  
7         ROWS BETWEEN 3 PRECEDING AND 3 FOLLOWING  
8     ), 2) AS rolling_avg_flights,  
9     avg_temp_c,  
10    round(avg(avg_temp_c) OVER (  
11        ORDER BY hour  
12        ROWS BETWEEN 3 PRECEDING AND 3 FOLLOWING  
13    ), 2) AS rolling_avg_temp  
14 FROM hourly_flight_weather  
15 ORDER BY hour;
```

The raw `total_flights` column will be spiky. The rolling average smooths it out so you can see whether flight volume trends up during the day and down at night. Which it will, unless Portland's airport has gone rogue.

6.38 Comparing Raw vs. Smoothed Data

```
1 -- Show the difference between raw and smoothed values  
2 SELECT  
3     hour,  
4     total_flights AS raw,
```

```

5     round(avg(total_flights) OVER (
6         ORDER BY hour
7         ROWS BETWEEN 6 PRECEDING AND 6 FOLLOWING
8     ), 2) AS smoothed_12hr,
9     total_flights - round(avg(total_flights) OVER (
10        ORDER BY hour
11        ROWS BETWEEN 6 PRECEDING AND 6 FOLLOWING
12    ), 2) AS residual
13 FROM hourly_flight_weather
14 ORDER BY hour;

```

The residual column shows how much each hour deviates from its local trend. Large positive residuals are unusually busy hours. Large negative residuals are unusually quiet. These outliers are often the most interesting data points.

6.39 Cumulative Flights Over Time

```

1 -- Running total of flights observed
2 SELECT
3     hour,
4     total_flights,
5     sum(total_flights) OVER (ORDER BY hour) AS cumulative_flights
6 FROM hourly_flight_weather
7 ORDER BY hour;

```

This gives you a sense of how data has been accumulating over the two weeks. If there is a flat spot, your scraper probably went down. If there is a sudden jump, you caught a busy travel day.

6.40 Rates and Comparisons

6.41 Calculating Meaningful Rates

Raw counts are not always comparable. Flights at 3 AM vs. 3 PM are different contexts. Rates normalize the data so comparisons make sense.

```

1 -- Average flights per hour by time of day
2 SELECT
3     extract(hour FROM hour) AS hour_of_day,
4     count(*) AS num_observations,
5     round(avg(total_flights), 2) AS avg_flights,
6     round(stddev_pop(total_flights), 2) AS stddev_flights

```

```

7 FROM hourly_flight_weather
8 GROUP BY extract(hour FROM hour)
9 ORDER BY hour_of_day;

```

This shows the daily rhythm of Portland air traffic. Expect a dip in the early morning hours and peaks in the afternoon.

6.42 CASE Statements

6.43 What is a CASE Statement?

A CASE statement is SQL's version of if/else. It lets you create new values based on conditions, right inside a query.

```

1 SELECT
2     callsign,
3     baro_altitude,
4     CASE
5         WHEN baro_altitude > 10000 THEN 'High altitude'
6         WHEN baro_altitude > 3000 THEN 'Mid altitude'
7         ELSE 'Low altitude'
8     END AS altitude_category
9 FROM flights;

```

PostgreSQL evaluates the WHEN conditions top to bottom. The first one that is true wins. If none match, ELSE is returned. If there is no ELSE and nothing matches, you get NULL.

6.44 Two Forms of CASE

Searched CASE – each WHEN has its own condition (most flexible):

```

1 CASE
2     WHEN temperature_c > 20 THEN 'Warm'
3     WHEN temperature_c > 10 THEN 'Mild'
4     WHEN temperature_c > 0 THEN 'Cool'
5     ELSE 'Cold'
6 END

```

Simple CASE – compares one expression against values (cleaner for equality checks):

```

1 CASE weather_code
2     WHEN 0 THEN 'Clear'
3     WHEN 1 THEN 'Mainly clear'
4     WHEN 2 THEN 'Partly cloudy'
5     WHEN 3 THEN 'Overcast'
6     ELSE 'Other'
7 END

```

Use simple CASE when you are matching one column against specific values. Use searched CASE when you need ranges, multiple columns, or complex logic.

6.45 Where Can You Use CASE?

Anywhere you can put an expression:

- **SELECT** – create computed columns (most common)
- **WHERE** – conditional filtering
- **GROUP BY** – group by computed categories
- **ORDER BY** – custom sort orders
- **Inside aggregates** – conditional counting

```

1 -- Conditional counting: count only clear-weather hours
2 SELECT
3     count(*) AS total_hours,
4     count(CASE WHEN weather_code = 0 THEN 1 END) AS clear_hours,
5     count(CASE WHEN weather_code >= 61 THEN 1 END) AS rainy_hours
6 FROM hourly_weather;

```

This pattern – CASE inside an aggregate – is one of the most useful tricks in SQL. It lets you pivot categories into columns without a pivot table.

6.46 CASE in GROUP BY

When you use CASE in a SELECT and want to group by it, you have to repeat the full expression in the GROUP BY clause. PostgreSQL does not let you reference a column alias in GROUP BY.

```

1 SELECT
2     CASE
3         WHEN wind_speed_kmh > 30 THEN 'High wind'
4         WHEN wind_speed_kmh > 15 THEN 'Moderate wind'
5         ELSE 'Low wind'
6     END AS wind_category,
7     count(*) AS num_observations,

```

```

8     round(avg(temperature_c), 2) AS avg_temp
9 FROM weather_observations
10 GROUP BY
11     CASE
12         WHEN wind_speed_kmh > 30 THEN 'High wind'
13         WHEN wind_speed_kmh > 15 THEN 'Moderate wind'
14         ELSE 'Low wind'
15     END;

```

Yes, you write it twice. It is annoying. A CTE or subquery can avoid the repetition if it bothers you.

6.47 Weekday vs. Weekend Traffic

```

1 -- Compare weekday and weekend flight volumes
2 SELECT
3     CASE
4         WHEN extract(dow FROM hour) IN (0, 6) THEN 'Weekend'
5         ELSE 'Weekday'
6     END AS day_type,
7     count(*) AS num_hours,
8     round(avg(total_flights), 2) AS avg_flights_per_hour,
9     round(stddev_pop(total_flights), 2) AS stddev_flights,
10    round(avg(avg_temp_c), 2) AS avg_temp,
11    round(avg(avg_wind_kmh), 2) AS avg_wind
12 FROM hourly_flight_weather
13 GROUP BY
14     CASE
15         WHEN extract(dow FROM hour) IN (0, 6) THEN 'Weekend'
16         ELSE 'Weekday'
17     END;

```

If weekday traffic is significantly higher, that is commercial aviation at work. If weekend traffic is comparable, Portland's general aviation community is active.

6.48 Flight Rate by Weather Condition

```

1 -- Average flights per hour grouped by weather code
2 SELECT
3     predominant_weather_code,
4     CASE predominant_weather_code
5         WHEN 0 THEN 'Clear'

```

```

6         WHEN 1 THEN 'Mainly clear'
7         WHEN 2 THEN 'Partly cloudy'
8         WHEN 3 THEN 'Overcast'
9         WHEN 45 THEN 'Fog'
10        WHEN 51 THEN 'Light drizzle'
11        WHEN 53 THEN 'Moderate drizzle'
12        WHEN 55 THEN 'Dense drizzle'
13        WHEN 61 THEN 'Light rain'
14        WHEN 63 THEN 'Moderate rain'
15        WHEN 65 THEN 'Heavy rain'
16        ELSE 'Other (' || predominant_weather_code || ')'
17    END AS weather_description,
18    count(*) AS num_hours,
19    round(avg(total_flights), 2) AS avg_flights,
20    round(avg(avg_wind_kmh), 2) AS avg_wind
21 FROM hourly_flight_weather
22 GROUP BY predominant_weather_code
23 ORDER BY avg_flights DESC;

```

This is the most direct answer to our research question. If “Clear” has significantly more flights than “Moderate rain,” weather matters. If they are similar, Portland pilots just do not care about rain. Which, honestly, tracks.

6.49 Putting It All Together: Multi-Factor Analysis

```

1  -- Combine temperature bins, wind bins, and weather code
2  SELECT
3      CASE
4          WHEN avg_temp_c < 5 THEN 'Cold (<5C)'
5          WHEN avg_temp_c < 10 THEN 'Cool (5-10C)'
6          WHEN avg_temp_c < 15 THEN 'Mild (10-15C)'
7          ELSE 'Warm (>15C)'
8      END AS temp_bin,
9      CASE
10         WHEN avg_wind_kmh < 10 THEN 'Calm (<10 km/h)'
11         WHEN avg_wind_kmh < 20 THEN 'Breezy (10-20 km/h)'
12         ELSE 'Windy (>20 km/h)'
13     END AS wind_bin,
14     count(*) AS num_hours,
15     round(avg(total_flights), 2) AS avg_flights,
16     round(avg(avg_flights_per_snapshot), 2) AS avg_per_snapshot
17 FROM hourly_flight_weather

```

```

18 GROUP BY
19     CASE
20         WHEN avg_temp_c < 5 THEN 'Cold (<5C)'
21         WHEN avg_temp_c < 10 THEN 'Cool (5-10C)'
22         WHEN avg_temp_c < 15 THEN 'Mild (10-15C)'
23         ELSE 'Warm (>15C)'
24     END,
25     CASE
26         WHEN avg_wind_kmh < 10 THEN 'Calm (<10 km/h)'
27         WHEN avg_wind_kmh < 20 THEN 'Breezy (10-20 km/h)'
28         ELSE 'Windy (>20 km/h)'
29     END
30 ORDER BY avg_flights DESC;

```

This crosstab shows flight volume for every combination of temperature and wind conditions. If “Warm + Calm” consistently has the most flights and “Cold + Windy” has the fewest, we have a clear pattern.

7 Part 7: Answering the Research Question

7.1 What Did We Find?

7.2 The Verdict

Research Question: Does weather affect the volume and behavior of air traffic over Portland?

Based on our analysis:

1. **Temperature and flight volume:** Likely a weak-to-moderate positive correlation. Warmer hours have somewhat more flights, but time of day is a massive confound since it is warmer during the day when more flights operate anyway.
2. **Wind speed and flight volume:** Likely a weak negative correlation. Very windy hours may have slightly fewer flights, but Portland rarely gets wind severe enough to ground commercial aviation.
3. **Weather code and flight volume:** Probably the strongest signal. Clear vs. rainy conditions may show meaningful differences, especially for general aviation traffic.
4. **The elephant in the room:** Time of day and airline schedules dominate flight patterns far more than weather does in two weeks of March data.

7.3 Correlation Is Not Causation

This bears repeating every time you run `corr()`:

- Correlation measures **linear association**, not cause and effect
- Temperature correlates with flight volume, but temperature does not *cause* flights to take off
- Both are driven by **time of day** (confounding variable)
- To truly isolate weather effects, you would need to control for time of day, day of week, and seasonal patterns

Two weeks of data in March is a starting point, not a conclusion. A full year of data across all seasons would give us much more to work with.

7.4 Where Do We Go From Here?

7.5 Coming Up Next

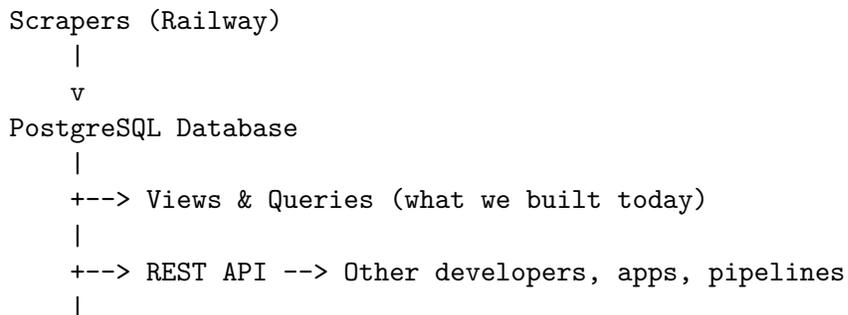
We have a normalized database. We have views that answer real questions. We have statistical analysis that tells a story. But right now, only people who know SQL can access any of it.

That changes next.

API Service over PostgreSQL. We will build a REST API that sits on top of our database, so other developers can query our flight and weather data over HTTP without writing SQL. Your database becomes a service that anyone can consume – a web app, a mobile app, another team’s pipeline.

Dashboards with Grafana. We will connect Grafana directly to our PostgreSQL database and build interactive dashboards on top of our queries and views. Instead of running SQL and reading result tables, you will have live charts and graphs that update as new data flows in.

7.6 The Big Picture



+--> Grafana Dashboard --> Visual monitoring & analysis

Your project will be become (hopefully!) a complete data platform: collection, storage, transformation, analysis, access, and visualization. That is data engineering.

7.7 References

- DeBarros, A. (2022). *Practical SQL: A Beginner's Guide to Storytelling with Data* (2nd ed.). No Starch Press. Chapter 11.
- OpenSky Network API: <https://openskynetwork.github.io/opensky-api/>
- Open-Meteo API: <https://open-meteo.com/en/docs>
- WMO Weather Codes: <https://www.nodc.noaa.gov/archive/arc0021/0002199/1.1/data/0-data/HTML/WMO-CODE/WMO4677.HTM>
- PostgreSQL Aggregate Functions: <https://www.postgresql.org/docs/current/functions-aggregate.html>