

Lecture 10-2: Working with Dates and Times

DATA 503: Fundamentals of Data Engineering

Lucas P. Cordova, Ph.D.

2026-03-16

This lecture covers working with dates, times, and time zones in PostgreSQL. We explore date/time data types, extraction functions, time zone handling, date math, and interval calculations. We apply these concepts to NYC yellow taxi trip data and Amtrak train schedules, then connect them to the pipeline work you have been doing all semester. Based on Chapter 12 of Practical SQL, 2nd Edition.

Table of contents

1 Part 1: Setting Up	2
2 Part 2: Date and Time Data Types	3
3 Part 3: Extracting Date and Time Components	5
4 Part 4: Creating Dates and Times	8
5 Part 5: Time Zones	9
6 Part 6: Date and Time Math	12
7 Part 7: NYC Taxi Trip Analysis	14
8 Part 8: Amtrak Train Trips	18
9 Part 9: CTEs and Useful PostgreSQL Features	21
10 Part 10: Putting It All Together (Exercises)	27
11 Part 11: What We Learned	29

1 Part 1: Setting Up

Before we start manipulating time itself, we need to load some data.

1.1 Loading the Chapter Database

1.2 Step 1: Create the Database

Create a fresh database for this lecture:

```
1 createdb -U postgres dates_and_times;
```

Or use Beekeeper Studio / pgAdmin to create it through the GUI.

1.3 Step 2: Create and Load the Taxi Table

The NYC taxi trips table contains 368,774 metered trips from June 1, 2016:

```
1 CREATE TABLE nyc_yellow_taxi_trips (  
2     trip_id bigint GENERATED ALWAYS AS IDENTITY PRIMARY KEY,  
3     vendor_id text NOT NULL,  
4     tpep_pickup_datetime timestamptz NOT NULL,  
5     tpep_dropoff_datetime timestamptz NOT NULL,  
6     passenger_count integer NOT NULL,  
7     trip_distance numeric(8,2) NOT NULL,  
8     pickup_longitude numeric(18,15) NOT NULL,  
9     pickup_latitude numeric(18,15) NOT NULL,  
10    rate_code_id text NOT NULL,  
11    store_and_fwd_flag text NOT NULL,  
12    dropoff_longitude numeric(18,15) NOT NULL,  
13    dropoff_latitude numeric(18,15) NOT NULL,  
14    payment_type text NOT NULL,  
15    fare_amount numeric(9,2) NOT NULL,  
16    extra numeric(9,2) NOT NULL,  
17    mta_tax numeric(5,2) NOT NULL,  
18    tip_amount numeric(9,2) NOT NULL,  
19    tolls_amount numeric(9,2) NOT NULL,  
20    improvement_surcharge numeric(9,2) NOT NULL,  
21    total_amount numeric(9,2) NOT NULL  
22 );
```

1.4 Step 3: Import the CSV

```
1 \copy nyc_yellow_taxi_trips (vendor_id, tpep_pickup_datetime,  
2     tpep_dropoff_datetime, passenger_count, trip_distance,  
3     pickup_longitude, pickup_latitude, rate_code_id,  
4     store_and_fwd_flag, dropoff_longitude, dropoff_latitude,  
5     payment_type, fare_amount, extra, mta_tax, tip_amount,  
6     tolls_amount, improvement_surcharge, total_amount)  
7 FROM '/path/to/nyc_yellow_taxi_trips.csv'  
8 WITH (FORMAT CSV, HEADER);
```

Replace /path/to/ with the actual path to the CSV file on your machine.

1.5 Step 4: Create an Index and Verify

```
1 CREATE INDEX tpep_pickup_idx  
2 ON nyc_yellow_taxi_trips (tpep_pickup_datetime);  
3  
4 -- Verify: you should see 368,774 rows  
5 SELECT count(*) FROM nyc_yellow_taxi_trips;
```

1.6 Step 5: Set the Time Zone

The taxi data is from New York City:

```
1 SET TIME ZONE 'US/Eastern';
```

This is a **session-level** setting. It changes how timestamps are *displayed*, not how they are *stored*. PostgreSQL always stores `timestamptz` as UTC internally.

2 Part 2: Date and Time Data Types

Time to talk about time.

2.1 The Four Types

2.2 PostgreSQL Date/Time Types

PostgreSQL gives us four data types for temporal data:

Type	Stores	Example
<code>timestampz</code>	Date + time + time zone	2022-12-01 18:37:12 EST
<code>date</code>	Date only	2022-12-01
<code>time</code>	Time only	18:37:12
<code>interval</code>	A duration	2 days 3 hours

The first three are **datetime types**. The fourth, `interval`, represents a length of time rather than a specific moment.

2.3 Which One Should You Use?

Almost always: `timestampz`.

A timestamp without a time zone is like a coordinate without a datum. Sure, “45.5, -122.6” means something, but in which reference system? NAD27? WGS84? The difference matters.

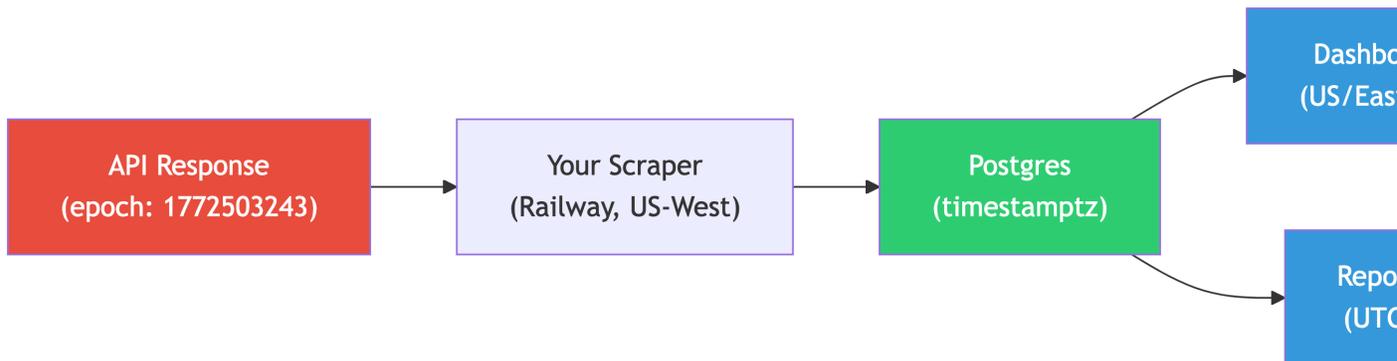
`date` is fine when you genuinely only care about the calendar date. `time` alone is almost never useful because a time without a date or zone is ambiguous.

`interval` shows up as the *result* of subtracting two timestamps. Think of it as the answer to “how long?” rather than “when?”

2.4 Why This Matters for Data Engineering

Think about your Railway scrapers. Each row has a `created_at` timestamp. If that column were `timestamp` instead of `timestampz`, and your Railway instance moved to a different region, your time-series data would silently break. Every record would look like it shifted by hours.

This is not hypothetical. It is one of the most common data pipeline bugs in production systems.



One moment in time, stored once, displayed correctly everywhere. That is the entire point of `timestampz`.

2.5 Try It: Pick the Type

What data type would you use for each of these?

1. When a scraper pulled data from an API
2. A person's date of birth
3. How long an ETL job took to run
4. The daily closing time of a store

...

Answers:

1. `timestampz` (you need to know which time zone the server was in)
2. `date` (nobody was born at 3:47 PM UTC)
3. `interval` (e.g., '2 hours 23 minutes')
4. Trick question. `time` *seems* right, but without a date, you cannot handle daylight saving changes. In practice, store this as text or use application logic.

3 Part 3: Extracting Date and Time Components

Sometimes you do not need the whole timestamp. You just want the year, or the hour, or the day of the week.

3.1 `date_part()` and `extract()`

3.2 Pulling Apart a Timestamp

The `date_part()` function extracts a single component:

```
1 SELECT
2     date_part('year',    '2022-12-01 18:37:12 EST'::timestampz) AS year,
3     date_part('month',  '2022-12-01 18:37:12 EST'::timestampz) AS month,
4     date_part('day',    '2022-12-01 18:37:12 EST'::timestampz) AS day,
5     date_part('hour',   '2022-12-01 18:37:12 EST'::timestampz) AS hour,
6     date_part('minute', '2022-12-01 18:37:12 EST'::timestampz) AS minute,
7     date_part('seconds', '2022-12-01 18:37:12 EST'::timestampz) AS seconds;
```

Run this now. Your hour value might differ from your neighbor's. Why?

...

Because PostgreSQL converts the timestamp to *your* session time zone. EST is UTC-5. If your server is set to Pacific time (UTC-8), the hour shows as 15 instead of 18. The underlying moment is identical.

3.3 More Components

```
1 SELECT
2     date_part('timezone_hour', '2022-12-01 18:37:12 EST'::timestampz) AS tz,
3     date_part('week',          '2022-12-01 18:37:12 EST'::timestampz) AS week,
4     date_part('quarter',       '2022-12-01 18:37:12 EST'::timestampz) AS quarter,
5     date_part('epoch',         '2022-12-01 18:37:12 EST'::timestampz) AS epoch;
```

Component	What It Returns
timezone_hour	UTC offset in hours (e.g., -5 for EST)
week	ISO 8601 week number (1-53)
quarter	Quarter of the year (1-4)
epoch	Seconds since January 1, 1970 00:00:00 UTC

Epoch is how computers think about time. Every timestamp is just a big number counting seconds from 1970. Remember the OpenSky API? The `time` field in your flight data was an epoch value. Now you know how to work with them.

3.4 The SQL Standard Alternative: `extract()`

`extract()` does the same thing with slightly different syntax:

```
1 SELECT extract(year FROM '2022-12-01 18:37:12 EST'::timestampz) AS year;
```

No quotes around the component name, and `FROM` instead of a comma. Both functions return the same result. Use whichever your team prefers.

3.5 Try It: Your Turn

Write a query that extracts the **day of the week** (dow) and the **hour** from the current timestamp.

Hint: `date_part('dow', ...)` returns 0 for Sunday through 6 for Saturday.

```
1 -- Your query here
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

3.6 date_trunc(): Rounding Timestamps

3.7 Truncating to a Precision

While `date_part()` pulls out a component, `date_trunc()` rounds a timestamp *down* to a specified precision:

```
1 SELECT
2     date_trunc('year',    '2022-12-01 18:37:12 EST'::timestampz) AS year,
3     date_trunc('month',  '2022-12-01 18:37:12 EST'::timestampz) AS month,
4     date_trunc('day',    '2022-12-01 18:37:12 EST'::timestampz) AS day,
5     date_trunc('hour',   '2022-12-01 18:37:12 EST'::timestampz) AS hour;
```

Precision	Result
year	2022-01-01 00:00:00-05
month	2022-12-01 00:00:00-05
day	2022-12-01 00:00:00-05
hour	2022-12-01 18:00:00-05

This is essential for time-series aggregation. Want daily totals? `GROUP BY date_trunc('day', timestamp_col)`. Hourly rollups? `GROUP BY date_trunc('hour', timestamp_col)`.

In data pipelines, `date_trunc` is how you build time-bucketed summary tables and partition data by date ranges.

4 Part 4: Creating Dates and Times

Sometimes your data arrives in pieces. PostgreSQL can reassemble them.

4.1 Making Datetimes

4.2 Building from Components

```
1 -- Make a date
2 SELECT make_date(2022, 2, 22);
3
4 -- Make a time (no time zone)
5 SELECT make_time(18, 4, 30.3);
6
7 -- Make a full timestamp with time zone
8 SELECT make_timestamptz(2022, 2, 22, 18, 4, 30.3, 'Europe/Lisbon');
```

Run the `make_timestamptz` query. If your session is in US/Eastern, you should see 2022-02-22 13:04:30.3-05. Lisbon is UTC+0, Eastern in February is UTC-5, so 6:04 PM Lisbon = 1:04 PM New York.

4.3 Getting the Current Date and Time

```
1 SELECT
2     current_timestamp,    -- timestamp with time zone (SQL standard)
3     localtimestamp,      -- timestamp WITHOUT time zone (avoid this)
4     current_date,        -- just the date
5     now();               -- PostgreSQL shorthand for current_timestamp
```

All of these return the time at the *start* of the transaction. If your query takes 10 seconds, every row gets the *same* timestamp.

4.4 `current_timestamp` vs. `clock_timestamp()`

What if you want the *actual* clock time as each row is processed?

```
1 CREATE TABLE current_time_example (
2     time_id integer GENERATED ALWAYS AS IDENTITY,
3     current_timestamp_col timestamptz,
4     clock_timestamp_col timestamptz
5 );
```

```

6
7 INSERT INTO current_time_example
8     (current_timestamp_col, clock_timestamp_col)
9     (SELECT current_timestamp,
10          clock_timestamp()
11     FROM generate_series(1,1000));
12
13 SELECT * FROM current_time_example;

```

The `current_timestamp_col` is identical for all 1,000 rows. The `clock_timestamp_col` increases slightly with each row.

When would this matter? Imagine logging the exact processing time of each row during a large ETL load. `current_timestamp` would give you one timestamp for the entire batch. `clock_timestamp()` gives you per-row precision. This is how you measure row-level throughput in a pipeline.

5 Part 5: Time Zones

This is where dates and times get *spicy*.

5.1 Understanding Time Zones

5.2 Why Time Zones Matter

A server in Oregon logs an event at 2026-03-15 14:00:00. A server in New York logs an event at 2026-03-15 17:00:00. Which happened first?

...

They happened at the **exact same moment**. Oregon is UTC-7, New York is UTC-4. Both are 21:00 UTC.

Without time zone information, you would think the Oregon event happened 3 hours earlier. In a distributed pipeline where data flows through servers in multiple regions, this creates silent data corruption.

5.3 Checking Your Time Zone

```
1 SHOW timezone;
```

If you followed setup, it should say US/Eastern.

5.4 Browsing Available Time Zones

```
1  -- Abbreviations (short list)
2  SELECT * FROM pg_timezone_abbrevs ORDER BY abbrev;
3
4  -- Full named zones (long list)
5  SELECT * FROM pg_timezone_names ORDER BY name;
6
7  -- Filter to a region
8  SELECT * FROM pg_timezone_names
9  WHERE name LIKE 'US%'
10 ORDER BY name;
```

5.5 Try It: Find a Time Zone

Write a query to find all time zones in Asia. How many are there?

```
1  -- Your query here
. . .
1  SELECT count(*) FROM pg_timezone_names
2  WHERE name LIKE 'Asia/%';
```

5.6 Changing and Converting Time Zones

5.7 SET TIME ZONE

Change your session's display time zone:

```
1  SET TIME ZONE 'US/Pacific';
```

This changes how timestamps are *displayed*, not how they are *stored*. Internally, PostgreSQL always stores `timestampz` values as UTC.

5.8 The Great Demonstration

Watch the same moment change its display:

```
1 CREATE TABLE time_zone_test (  
2     test_date timestamptz  
3 );  
4  
5 INSERT INTO time_zone_test VALUES ('2023-01-01 4:00');  
6  
7 SET TIME ZONE 'US/Pacific';  
8 SELECT test_date FROM time_zone_test;  
9 -- Shows: 2023-01-01 04:00:00-08  
10  
11 SET TIME ZONE 'US/Eastern';  
12 SELECT test_date FROM time_zone_test;  
13 -- Shows: 2023-01-01 07:00:00-05
```

Same row, same data on disk, different display. This is why `timestamptz` is the right default.

5.9 AT TIME ZONE: One-Off Conversions

Peek at a value in a different time zone without changing your session:

```
1 SELECT test_date AT TIME ZONE 'Asia/Seoul'  
2 FROM time_zone_test;  
3 -- Shows: 2023-01-01 21:00:00
```

4 AM Pacific on January 1st is 9 PM the same day in Seoul.

Quirk alert: When you use `AT TIME ZONE` on a `timestamptz`, the output becomes a `timestamp` *without* time zone. PostgreSQL strips the zone because you already specified which one you want.

5.10 Try It: World Clocks

It is midnight on New Year's Day 2030 in New York ('2030-01-01 00:00:00 US/Eastern').

Write a query showing what time it is at that moment in London, Tokyo, and Sydney.

```
1 -- Your query here  
...  
...
```

```

1 SELECT
2     '2030-01-01 00:00:00 US/Eastern'::timestamptz
3     AT TIME ZONE 'Europe/London' AS london,
4     '2030-01-01 00:00:00 US/Eastern'::timestamptz
5     AT TIME ZONE 'Asia/Tokyo' AS tokyo,
6     '2030-01-01 00:00:00 US/Eastern'::timestamptz
7     AT TIME ZONE 'Australia/Sydney' AS sydney;

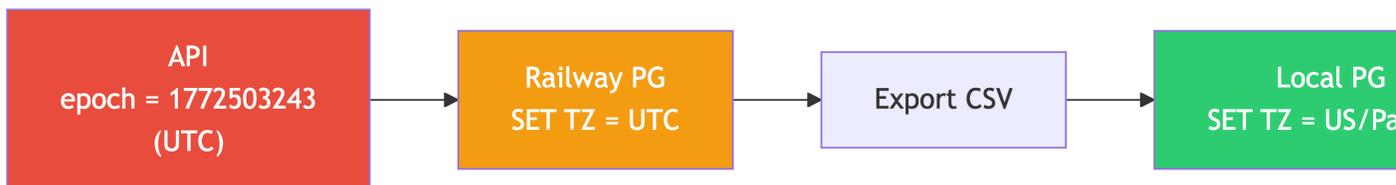
```

London: 5 AM. Tokyo: 2 PM. Sydney: 4 PM.

5.11 Time Zones in Data Pipelines

5.12 The Pipeline Time Zone Problem

Consider your Railway scrapers. The API returns Unix epoch timestamps (UTC). Your Railway Postgres might be in US/West. Your local Postgres for analysis might be in US/Pacific.



If every column uses `timestamptz`, this chain works seamlessly. PostgreSQL handles the conversions. If any column uses plain `timestamp`, you get silent drift.

5.13 Best Practices for Timestamps in Pipelines

1. Always use `timestamptz` for storage
2. Store raw epoch values alongside human-readable timestamps when ingesting API data
3. Set session time zone explicitly at the start of every script
4. Convert at the display layer, not at the storage layer
5. Document your convention in your pipeline README

6 Part 6: Date and Time Math

You can do arithmetic with dates and times, and it actually makes sense.

6.1 Calculations with Dates

6.2 Subtracting Dates

Subtract one date from another and you get an integer (days between them):

```
1 SELECT '1929-09-30'::date - '1929-09-27'::date;
```

Result: 3. Three days between September 27 and September 30, 1929.

6.3 Adding Intervals

Add an interval to a date:

```
1 SELECT '1929-09-30'::date + '5 years'::interval;
```

Result: 1934-09-30 00:00:00. PostgreSQL handles leap years, month lengths, and calendar quirks automatically.

6.4 Subtracting Timestamps

Subtract two timestamps and you get an interval:

```
1 SELECT
2     '2026-03-18 17:00:00 US/Eastern'::timestampz -
3     '2026-03-18 09:00:00 US/Eastern'::timestampz
4     AS work_day_length;
```

Result: 08:00:00.

6.5 Getting Duration in Specific Units

When you need duration as a number (not an interval), extract the epoch:

```
1 SELECT
2     EXTRACT(EPOCH FROM (
3         '2026-03-18 17:00:00 US/Eastern'::timestampz -
4         '2026-03-18 09:00:00 US/Eastern'::timestampz
5     )) / 3600 AS hours;
```

Result: 8. This pattern – subtract timestamps, extract epoch, divide – is how you compute durations in minutes, hours, or days as plain numbers. Essential for aggregations and comparisons.

6.6 Try It: How Old Is PostgreSQL?

PostgreSQL's first release was on July 8, 1996. Write a query calculating how many days ago that was.

```
1 -- Your query here
. . .
1 SELECT current_date - '1996-07-08'::date AS days_since_postgres;
```

7 Part 7: NYC Taxi Trip Analysis

Time to put all of this to work on real data. We have 368,774 taxi rides from June 1, 2016, in New York City.

7.1 Patterns in Pickup Times

7.2 The Busiest Hours

When do New Yorkers take the most taxi rides?

```
1 SET TIME ZONE 'US/Eastern';
2
3 SELECT
4     date_part('hour', tpep_pickup_datetime) AS trip_hour,
5     count(*)
6 FROM nyc_yellow_taxi_trips
7 GROUP BY trip_hour
8 ORDER BY trip_hour;
```

Run this now. You should get 24 rows.

. . .

The busiest hours are in the evening (6 PM to 10 PM). The slowest are early morning (2 AM to 5 AM). The evening peak reflects commutes and nightlife on a summer Wednesday.

7.3 Try It: Busiest Hour

Write a query to find the **single busiest hour** and its trip count.

```
1 -- Your query here
. . .
1 SELECT
2     date_part('hour', tpep_pickup_datetime) AS trip_hour,
3     count(*) AS num_trips
4 FROM nyc_yellow_taxi_trips
5 GROUP BY trip_hour
6 ORDER BY num_trips DESC
7 LIMIT 1;
```

7.4 Trip Duration Analysis

7.5 Median Trip Time by Hour

How long do taxi rides take at different times of day?

```
1 SELECT
2     date_part('hour', tpep_pickup_datetime) AS trip_hour,
3     percentile_cont(.5)
4         WITHIN GROUP (ORDER BY
5             tpep_dropoff_datetime - tpep_pickup_datetime) AS median_trip
6 FROM nyc_yellow_taxi_trips
7 GROUP BY trip_hour
8 ORDER BY trip_hour;
```

The longest median trips are around midday (~15 minutes at 1 PM). The shortest are early morning (~8 minutes at 5 AM). Traffic matters.

7.6 Try It: The Longest Rides

Find the 10 longest taxi rides by duration. Include pickup time, dropoff time, calculated duration, and trip distance.

Do any results look suspicious?

```
1 -- Your query here
. . .
```

```

1 SELECT
2     tpep_pickup_datetime,
3     tpep_dropoff_datetime,
4     tpep_dropoff_datetime - tpep_pickup_datetime AS duration,
5     trip_distance
6 FROM nyc_yellow_taxi_trips
7 ORDER BY duration DESC
8 LIMIT 10;

```

You will see rides lasting 20+ hours with zero distance. Data quality issues. The meter was left running, or timestamps were recorded incorrectly.

This is why data engineers never trust raw data. Always profile. Always validate.

7.7 Try It: Average Fare by Hour

Write a query showing the average `total_amount` for each hour. Which hour has the highest average fare?

```

1 -- Your query here
2
3 ...
4
5 SELECT
6     date_part('hour', tpep_pickup_datetime) AS trip_hour,
7     round(avg(total_amount), 2) AS avg_fare
8 FROM nyc_yellow_taxi_trips
9 GROUP BY trip_hour
10 ORDER BY avg_fare DESC;

```

Early morning hours (4-5 AM) have the highest average fares. Fewer but longer rides – airport runs and long distances with no traffic.

7.8 Time-Series Aggregation with `date_trunc`

7.9 Hourly Rollups

In data engineering, you frequently need to aggregate time-series data into fixed buckets. This is the basis of every dashboard and monitoring system.

```

1 SELECT
2     date_trunc('hour', tpep_pickup_datetime) AS pickup_hour,
3     count(*) AS trips,
4     round(avg(total_amount), 2) AS avg_fare
5 FROM nyc_yellow_taxi_trips
6 GROUP BY pickup_hour
7 ORDER BY pickup_hour;

```

Unlike `date_part('hour', ...)` which gives you just the number 0-23, `date_trunc` preserves the full timestamp truncated to the hour. This matters when you have multi-day data – hour 14 on Monday is different from hour 14 on Tuesday.

7.10 Try It: Daily Revenue

Write a query using `date_trunc` to calculate total daily revenue (`total_amount`). Since this dataset is one day, you should get a single row.

```

1 -- Your query here
2
3 ...
4
5 SELECT
6     date_trunc('day', tpep_pickup_datetime) AS trip_date,
7     round(sum(total_amount), 2) AS daily_revenue,
8     count(*) AS total_trips
9 FROM nyc_yellow_taxi_trips
10 GROUP BY trip_date;

```

7.11 Generating Time Series with `generate_series`

7.12 Filling in Missing Time Slots

Real-world time-series data has gaps. If no trips happened at 3:17 AM, that minute is simply absent from your results. Dashboards hate gaps. Pipeline consumers hate gaps.

`generate_series` creates a continuous sequence of timestamps:

```

1 SELECT generate_series(
2     '2016-06-01 00:00:00 US/Eastern'::timestampz,
3     '2016-06-01 23:00:00 US/Eastern'::timestampz,
4     '1 hour'::interval
5 ) AS hour_slot;

```

This produces 24 rows, one for each hour, regardless of whether any trips happened.

7.13 Joining with Generated Series

To get a complete hourly report (with zeros for empty hours):

```
1 WITH hours AS (  
2     SELECT generate_series(  
3         '2016-06-01 00:00:00 US/Eastern'::timestampz,  
4         '2016-06-01 23:00:00 US/Eastern'::timestampz,  
5         '1 hour'::interval  
6     ) AS hour_slot  
7 )  
8 SELECT  
9     h.hour_slot,  
10    coalesce(count(t.trip_id), 0) AS trip_count  
11 FROM hours h  
12 LEFT JOIN nyc_yellow_taxi_trips t  
13     ON date_trunc('hour', t.tpep_pickup_datetime) = h.hour_slot  
14 GROUP BY h.hour_slot  
15 ORDER BY h.hour_slot;
```

This pattern – generate a time spine, left join your data – is fundamental to building reliable time-series reports. You will see it in every analytics platform.

8 Part 8: Amtrak Train Trips

Let us switch from taxis to trains and watch timestamps cross time zones.

8.1 Building the Train Data

8.2 Creating and Loading the Table

```
1 CREATE TABLE train_rides (  
2     trip_id bigint GENERATED ALWAYS AS IDENTITY PRIMARY KEY,  
3     segment text NOT NULL,  
4     departure timestampz NOT NULL,  
5     arrival timestampz NOT NULL  
6 );  
7  
8 INSERT INTO train_rides (segment, departure, arrival)  
9 VALUES  
10    ('Chicago to New York',
```

```

11         '2020-11-13 21:30 CST', '2020-11-14 18:23 EST'),
12     ('New York to New Orleans',
13         '2020-11-15 14:15 EST', '2020-11-16 19:32 CST'),
14     ('New Orleans to Los Angeles',
15         '2020-11-17 13:45 CST', '2020-11-18 9:00 PST'),
16     ('Los Angeles to San Francisco',
17         '2020-11-19 10:10 PST', '2020-11-19 21:24 PST'),
18     ('San Francisco to Denver',
19         '2020-11-20 9:10 PST', '2020-11-21 18:38 MST'),
20     ('Denver to Chicago',
21         '2020-11-22 19:10 MST', '2020-11-23 14:50 CST');

```

Notice how each timestamp includes the time zone of the city. The subtraction handles cross-zone arithmetic automatically because we used `timestamptz`.

8.3 Calculating Trip Durations

8.4 Segment Duration

```

1 SELECT segment,
2         to_char(departure, 'YYYY-MM-DD HH12:MI a.m. TZ') AS departure,
3         arrival - departure AS segment_duration
4 FROM train_rides;

```

Two things to notice:

1. `to_char()` formats timestamps into human-readable strings
2. **Subtracting timestamps** gives intervals. PostgreSQL correctly handles the time zone differences across segments

8.5 Try It: Longest Segment

Which segment takes the longest?

```

1 -- Your query here
2
3 ...
4
5 SELECT segment,
6         arrival - departure AS segment_duration
7 FROM train_rides
8 ORDER BY segment_duration DESC
9 LIMIT 1;

```

San Francisco to Denver: over 32 hours.

8.6 Cumulative Trip Duration

8.7 Running Total with Window Functions

How long has the entire trip taken after each segment?

```
1 SELECT segment,
2     arrival - departure AS segment_duration,
3     sum(arrival - departure) OVER (ORDER BY trip_id) AS cume_duration
4 FROM train_rides;
```

The output shows things like 2 days 85:47:00. PostgreSQL sums the days and hours separately without rolling up.

8.8 Fixing with justify_interval()

```
1 SELECT segment,
2     arrival - departure AS segment_duration,
3     justify_interval(sum(arrival - departure)
4                     OVER (ORDER BY trip_id)) AS cume_duration
5 FROM train_rides;
```

Now the final row shows 5 days 13:47:00. Five and a half days coast to coast by train.

8.9 Try It: Format the Output

Rewrite the query with the departure formatted as Mon DD, YYYY HH12:MI AM TZ using `to_char()`.

```
1 -- Your query here
2
3 ...
4
5 SET TIME ZONE 'US/Central';
6
7 SELECT segment,
8     to_char(departure, 'Mon DD, YYYY HH12:MI AM TZ') AS departure,
9     arrival - departure AS segment_duration,
10    justify_interval(sum(arrival - departure)
11                    OVER (ORDER BY trip_id)) AS cume_duration
12 FROM train_rides;
```

9 Part 9: CTEs and Useful PostgreSQL Features

Complex temporal queries get messy fast. CTEs and a few other PostgreSQL features make them manageable.

9.1 Common Table Expressions (CTEs)

9.2 What Is a CTE?

A **Common Table Expression** (CTE) is a temporary, named result set defined with `WITH` that exists only for the duration of a single query. Think of it as a named intermediate table.

```
1 WITH hourly_stats AS (  
2     SELECT  
3         date_part('hour', tpep_pickup_datetime) AS trip_hour,  
4         count(*) AS num_trips,  
5         round(avg(total_amount), 2) AS avg_fare  
6     FROM nyc_yellow_taxi_trips  
7     GROUP BY trip_hour  
8 )  
9 SELECT  
10    trip_hour,  
11    num_trips,  
12    avg_fare  
13 FROM hourly_stats  
14 WHERE num_trips > 10000  
15 ORDER BY trip_hour;
```

The `WITH` clause defines `hourly_stats`. The main query references it like a table. When the query finishes, `hourly_stats` disappears.

9.3 CTEs vs. Subqueries

The equivalent subquery version:

```
1 SELECT trip_hour, num_trips, avg_fare  
2 FROM (  
3     SELECT  
4         date_part('hour', tpep_pickup_datetime) AS trip_hour,  
5         count(*) AS num_trips,  
6         round(avg(total_amount), 2) AS avg_fare  
7     FROM nyc_yellow_taxi_trips
```

```

8     GROUP BY trip_hour
9 ) AS hourly_stats
10 WHERE num_trips > 10000
11 ORDER BY trip_hour;

```

Same result. But CTEs are:

- **More readable.** Logic flows top-to-bottom, not inside-out.
- **Reusable.** Reference the same CTE multiple times in one query.
- **Chainable.** Multiple CTEs can build on each other.
- **Self-documenting.** Good CTE names describe the data at each stage.

In data engineering, CTEs mirror the stages of a transformation pipeline. Each CTE is a logical step: extract, clean, aggregate, enrich.

9.4 Chaining CTEs: The Pipeline Pattern

You can define multiple CTEs separated by commas. Each one can reference the ones before it:

```

1 WITH hourly AS (
2     -- Stage 1: Aggregate by hour
3     SELECT
4         date_part('hour', tpep_pickup_datetime) AS trip_hour,
5         count(*) AS num_trips,
6         round(avg(total_amount), 2) AS avg_fare,
7         round(avg(trip_distance), 2) AS avg_distance
8     FROM nyc_yellow_taxi_trips
9     GROUP BY trip_hour
10 ),
11 ranked AS (
12     -- Stage 2: Rank hours by popularity
13     SELECT
14         trip_hour,
15         num_trips,
16         avg_fare,
17         avg_distance,
18         rank() OVER (ORDER BY num_trips DESC) AS popularity_rank
19     FROM hourly
20 )
21 -- Stage 3: Filter to top 5
22 SELECT *
23 FROM ranked

```

```
24 WHERE popularity_rank <= 5
25 ORDER BY popularity_rank;
```

Each CTE is simple. The chain is complex. This is the same principle behind well-designed ETL: small, testable, composable transformations.

9.5 Try It: Busiest vs. Quietest Hours

Using a CTE, find the 3 busiest and 3 quietest hours by trip count. Combine them into one result with UNION ALL.

```
1 -- Your query here
. . .
1 WITH hourly AS (
2     SELECT
3         date_part('hour', tpep_pickup_datetime) AS trip_hour,
4         count(*) AS num_trips
5     FROM nyc_yellow_taxi_trips
6     GROUP BY trip_hour
7 )
8 (SELECT trip_hour, num_trips, 'Busiest' AS category
9  FROM hourly ORDER BY num_trips DESC LIMIT 3)
10 UNION ALL
11 (SELECT trip_hour, num_trips, 'Quietest' AS category
12  FROM hourly ORDER BY num_trips ASC LIMIT 3)
13 ORDER BY category, num_trips DESC;
```

The CTE is defined once, used twice. Without it, the GROUP BY aggregation would be duplicated in both halves of the UNION ALL.

9.6 Useful PostgreSQL Features

9.7 COALESCE: Handling NULLs

COALESCE returns the first non-NULL value from its arguments:

```
1 SELECT coalesce(passenger_count, 0) AS passengers
2 FROM nyc_yellow_taxi_trips
3 LIMIT 5;
```

This is critical when combining `generate_series` with `LEFT JOIN`. Hours with no trips produce `NULL` counts, and `COALESCE(count, 0)` fills the gaps with zeros. In pipeline code, you will use `COALESCE` constantly – any time you join tables that might not match every row.

9.8 HAVING: Filtering After Aggregation

`WHERE` filters rows *before* grouping. `HAVING` filters groups *after* aggregation:

```
1 SELECT
2     date_part('hour', tpep_pickup_datetime) AS trip_hour,
3     count(*) AS num_trips,
4     round(avg(total_amount), 2) AS avg_fare
5 FROM nyc_yellow_taxi_trips
6 GROUP BY trip_hour
7 HAVING count(*) > 15000
8 ORDER BY trip_hour;
```

This returns only hours with more than 15,000 trips. The count does not exist until after `GROUP BY` runs, so `WHERE` cannot filter on it.

9.9 FILTER: Conditional Aggregation (PostgreSQL-Specific)

The `FILTER` clause lets you apply a condition to a specific aggregate without affecting others in the same query:

```
1 SELECT
2     date_part('hour', tpep_pickup_datetime) AS trip_hour,
3     count(*) AS total_trips,
4     count(*) FILTER (WHERE payment_type = '1') AS credit_trips,
5     count(*) FILTER (WHERE payment_type = '2') AS cash_trips,
6     round(
7         count(*) FILTER (WHERE payment_type = '1')::numeric
8         / count(*)::numeric * 100, 1
9     ) AS credit_pct
10 FROM nyc_yellow_taxi_trips
11 GROUP BY trip_hour
12 ORDER BY trip_hour;
```

Without `FILTER`, you would need `SUM(CASE WHEN ... THEN 1 ELSE 0 END)` for each conditional count. `FILTER` is cleaner. It is PostgreSQL-specific (not in the SQL standard), but most modern databases have equivalents.

9.10 Try It: Tipping Patterns by Hour

Using `FILTER`, write a query that shows for each hour: the total number of trips, the number of trips that included a tip (`tip_amount > 0`), and the tipping percentage. Which hours have the highest tipping rates?

```
1 -- Your query here
2
3 ...
4
5 SELECT
6     date_part('hour', tpep_pickup_datetime) AS trip_hour,
7     count(*) AS total_trips,
8     count(*) FILTER (WHERE tip_amount > 0) AS tipped_trips,
9     round(
10        count(*) FILTER (WHERE tip_amount > 0)::numeric
11        / count(*)::numeric * 100, 1
12    ) AS tip_pct
13 FROM nyc_yellow_taxi_trips
14 GROUP BY trip_hour
15 ORDER BY tip_pct DESC;
```

9.11 CTE + Date Functions: A Real Pattern

Combining CTEs with temporal analysis is where things get powerful. Here we compute each hour's trip count, then compare it to the overall average:

```
1 WITH hourly_counts AS (
2     SELECT
3         date_part('hour', tpep_pickup_datetime) AS trip_hour,
4         count(*) AS num_trips
5     FROM nyc_yellow_taxi_trips
6     GROUP BY trip_hour
7 ),
8 overall AS (
9     SELECT round(avg(num_trips), 0) AS avg_hourly_trips
10    FROM hourly_counts
11 )
12 SELECT
13     h.trip_hour,
14     h.num_trips,
15     o.avg_hourly_trips,
16     CASE
```

```

17         WHEN h.num_trips > o.avg_hourly_trips THEN 'Above Average'
18         WHEN h.num_trips < o.avg_hourly_trips THEN 'Below Average'
19         ELSE 'Average'
20     END AS comparison
21 FROM hourly_counts h
22 CROSS JOIN overall o
23 ORDER BY h.trip_hour;

```

CTE 1 aggregates by hour. CTE 2 computes the average across all hours. The main query joins them with `CROSS JOIN` (since `overall` is a single row) and classifies each hour.

This pattern – aggregate, compute a benchmark, compare – is the backbone of anomaly detection in data pipelines. Is this hour’s traffic normal? Is this day’s revenue within expected range? Same structure, different data.

9.12 Try It: Revenue Anomaly Detection

Using two CTEs, find the hours where total revenue is more than one standard deviation above or below the hourly average. Show the hour, its revenue, the average, the standard deviation, and label each as ‘High’, ‘Low’, or ‘Normal’.

Hint: use `stddev()` alongside `avg()` in the second CTE.

```

1  -- Your query here
2
3  ...
4
5  WITH hourly_revenue AS (
6      SELECT
7          date_part('hour', tpep_pickup_datetime) AS trip_hour,
8          round(sum(total_amount), 2) AS hour_revenue
9      FROM nyc_yellow_taxi_trips
10     GROUP BY trip_hour
11 ),
12 stats AS (
13     SELECT
14         round(avg(hour_revenue), 2) AS avg_revenue,
15         round(stddev(hour_revenue), 2) AS stddev_revenue
16     FROM hourly_revenue
17 )
18 SELECT
19     h.trip_hour,
20     h.hour_revenue,
21     s.avg_revenue,

```

```

18     s.stddev_revenue,
19     CASE
20         WHEN h.hour_revenue > s.avg_revenue + s.stddev_revenue THEN 'High'
21         WHEN h.hour_revenue < s.avg_revenue - s.stddev_revenue THEN 'Low'
22         ELSE 'Normal'
23     END AS anomaly_label
24 FROM hourly_revenue h
25 CROSS JOIN stats s
26 ORDER BY h.trip_hour;

```

10 Part 10: Putting It All Together (Exercises)

Let's combine several concepts from today into more complex queries.

10.1 Combined Exercises

10.2 Try It: Taxi Trips by Day of Week

Count trips by day of the week using a CASE statement for day names:

```

1  -- Your query here
. . .
1  SET TIME ZONE 'US/Eastern';
2
3  SELECT
4      CASE date_part('dow', tpep_pickup_datetime)
5          WHEN 0 THEN 'Sunday'
6          WHEN 1 THEN 'Monday'
7          WHEN 2 THEN 'Tuesday'
8          WHEN 3 THEN 'Wednesday'
9          WHEN 4 THEN 'Thursday'
10         WHEN 5 THEN 'Friday'
11         WHEN 6 THEN 'Saturday'
12     END AS day_name,
13     count(*) AS num_trips
14 FROM nyc_yellow_taxi_trips
15 GROUP BY day_name
16 ORDER BY num_trips DESC;

```

Since this dataset is from June 1, 2016 (a Wednesday), only one row appears. But the pattern works on multi-day datasets.

10.3 Try It: Time Zone Conversion

Show the first 5 taxi pickup times converted to Asia/Tokyo time.

```
1 -- Your query here
. . .
1 SELECT
2     tpep_pickup_datetime AS nyc_time,
3     tpep_pickup_datetime AT TIME ZONE 'Asia/Tokyo' AS tokyo_time
4 FROM nyc_yellow_taxi_trips
5 ORDER BY tpep_pickup_datetime
6 LIMIT 5;
```

Tokyo is 13 hours ahead of New York in summer. A 9 PM Wednesday pickup in NYC is 10 AM Thursday in Tokyo.

10.4 Try It: Trip Duration Percentiles

Calculate the 25th, 50th (median), and 75th percentile of trip duration:

```
1 -- Your query here
. . .
1 SELECT
2     percentile_cont(0.25)
3         WITHIN GROUP (ORDER BY tpep_dropoff_datetime - tpep_pickup_datetime)
4         AS p25,
5     percentile_cont(0.50)
6         WITHIN GROUP (ORDER BY tpep_dropoff_datetime - tpep_pickup_datetime)
7         AS median,
8     percentile_cont(0.75)
9         WITHIN GROUP (ORDER BY tpep_dropoff_datetime - tpep_pickup_datetime)
10        AS p75
11 FROM nyc_yellow_taxi_trips;
```

The IQR (P75 - P25) tells you how spread out typical trips are.

10.5 Try It: Correlation Challenge

Calculate the correlation coefficient between trip duration (in seconds) and total_amount. Then do the same for trip_distance and total_amount. Limit to rides of 3 hours or less.

Which is a better predictor of fare: time or distance?

```
1 -- Your query here
2
3
4
5
6
7
8
9
10
11
12
```

```
SELECT
  round(
    corr(
      date_part('epoch', tpep_dropoff_datetime - tpep_pickup_datetime),
      total_amount
    )::numeric, 4
  ) AS duration_fare_corr,
  round(
    corr(trip_distance, total_amount)::numeric, 4
  ) AS distance_fare_corr
FROM nyc_yellow_taxi_trips
WHERE tpep_dropoff_datetime - tpep_pickup_datetime <= '3 hours'::interval;
```

Distance correlates more strongly with fare. NYC taxi meters are primarily distance-based.

11 Part 11: What We Learned

11.1 Summary

11.2 Key Concepts

Concept	What It Does
timestampz	The go-to type for storing moments in time
date_part() / extract()	Pull year, month, hour, etc. from a timestamp
date_trunc()	Round timestamps down to a precision (essential for time-series)
make_timestampz()	Build timestamps from components
current_timestamp / now()	Time when the query started
clock_timestamp()	Actual clock time (changes per row)
SET TIME ZONE	Change display for your session
AT TIME ZONE	One-off conversion to another zone

Concept	What It Does
<code>to_char()</code>	Format timestamps as readable strings
<code>justify_interval()</code>	Clean up messy interval output
<code>generate_series()</code>	Create continuous time spines for gap-free reporting
<code>WITH (CTE)</code>	Define named temporary result sets; mirrors pipeline stages
<code>COALESCE</code>	Return the first non-NULL value (essential for outer joins)
<code>HAVING</code>	Filter groups after aggregation
<code>FILTER (WHERE ...)</code>	Apply conditions to individual aggregates (PostgreSQL-specific)

11.3 The Big Ideas

1. **Always use `timestampz`.** A timestamp without a zone is a bug waiting to happen.
2. **Time zones change the display, not the data.** PostgreSQL stores everything as UTC.
3. **Subtracting timestamps gives you intervals.** That is how you calculate durations.
4. **CTEs mirror pipeline stages.** Each CTE is a named, testable transformation step.
5. **`date_trunc` is your best friend** for time-series aggregation in pipelines.
6. **Generate time spines** with `generate_series` to fill gaps in your data.
7. **Use `FILTER` for conditional aggregation.** Cleaner than `SUM(CASE ...)` when you need multiple conditional counts or sums.
8. **Real-world data is messy.** Always sanity-check your results.

11.4 References

1. DeBarros, A. (2022). *Practical SQL: A Beginner's Guide to Storytelling with Data* (2nd ed.). No Starch Press. Chapter 12.
2. PostgreSQL Date/Time Functions: <https://www.postgresql.org/docs/current/functions-datetime.html>
3. PostgreSQL Formatting Functions: <https://www.postgresql.org/docs/current/functions-formatting.html>
4. NYC TLC Trip Record Data: <https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>
5. ISO 8601 Date and Time Standard: https://en.wikipedia.org/wiki/ISO_8601