

# Data Mining with Text

## DATA 503: Fundamentals of Data Engineering

Lucas P. Cordova, Ph.D.

2026-04-08

A concise version of text-oriented data work in PostgreSQL: string cleanup, POSIX regex for filters and extraction, parsing semi-structured crime narratives into columns, and full-text search with `tsvector` / `tsquery`, `@@`, GIN indexes, plus `ts_headline` and `ts_rank`. Checkpoints with hints, then solutions on the following slide. Based on Lecture 13-1 and Chapter 14 of Practical SQL, 2nd Edition.

## Table of contents

<b>1</b>	<b>Learning objectives</b>	<b>1</b>
<b>2</b>	<b>Load the practice data</b>	<b>2</b>
<b>3</b>	<b>String functions</b>	<b>2</b>
<b>4</b>	<b>Regular expressions</b>	<b>3</b>
<b>5</b>	<b>From blob to columns</b>	<b>5</b>
<b>6</b>	<b>Full-text search</b>	<b>8</b>
<b>7</b>	<b>Summary</b>	<b>10</b>

## 1 Learning objectives

### 1.1 What you should be able to do

1. Use core string functions for cleanup and slicing before heavier parsing
2. Apply `~` / `~*` and `regexp_match` / `regexp_matches` with capture groups to pull fields from text

3. Explain `tsvector` vs `tsquery`, filter with `@@`, and why GIN indexes matter for search at scale
4. Produce snippets with `ts_headline` and sort by relevance with `ts_rank` (including length normalization)

...

**Course connection:** Logs, tickets, and exports often arrive as prose. Regex and full-text search let you normalize and search near the database when you do not yet have a dedicated search service.

## 2 Load the practice data

### 2.1 Tables and sanity checks

Table	Role
<code>county_regex_demo</code>	Regex operator practice
<code>crime_reports</code>	Semi-structured police narratives
<code>president_speeches</code>	Full-text search (State of the Union)

Check this:

```
1 SELECT count(*) AS crime_reports FROM crime_reports;
2 SELECT count(*) AS speeches FROM president_speeches;
3 SELECT count(*) AS counties FROM county_regex_demo;
```

Expect 5, 79, and 8 respectively.

## 3 String functions

### 3.1 Why this layer exists

JSON and arrays get conference talks. Strings pay the rent: addresses, log lines, pasted emails, and CSV columns that are secretly paragraphs. PostgreSQL documents these under [String Functions](#).

### 3.2 Toolkit (high level)

Need	Examples
Case	upper, lower, initcap
Trim / measure	trim, char_length, position
Slice / replace	left, right, replace

### 3.3 Demo: cleanup

Run this.

```

1 SELECT trim(' Pat ') AS trimmed,
2        char_length(trim(' Pat ')) AS len,
3        left('703-555-1212', 3) AS area;

```

replace is greedy and left-to-right: `replace('aaa', 'aa', 'b')` yields `ba`, which surprises everyone exactly once.

### 3.4 Checkpoint: initcap and replace

**Task:** From the literal `' data engineering '`, produce a single row with:

- `title_phrase`: `initcap` after `trim` (one column)
- `demo_replace`: `replace('batch-etl-batch', 'batch', 'stream')` (shows non-overlapping replacement)

**Hint:** Nest `initcap(trim(...))`.

### 3.5 Solution: initcap and replace

```

1 SELECT initcap(trim(' data engineering ')) AS title_phrase,
2        replace('batch-etl-batch', 'batch', 'stream') AS demo_replace;

```

## 4 Regular expressions

### 4.1 What regex is for

A **regular expression** describes text shape. PostgreSQL uses **POSIX** patterns ([Pattern Matching](#)). Mind backslashes: you often write `'\\d'` in SQL so the pattern sees one `\`.

## 4.2 Notation (compressed)

Piece	Meaning
\d	Digit
{n,m}	Repeat between n and m
()	Group; also <b>capture</b> for <code>regexp_match</code>
\	Alternation
^ \$	Start and end of string (for whole-field tests)

## 4.3 Operators

Operator	Meaning
~	Matches (case sensitive)
~*	Matches (case insensitive)
!~ / !~*	Does not match

## 4.4 Demo: filter counties

Run this.

```
1 SELECT county_name
2 FROM county_regex_demo
3 WHERE county_name ~* 'ash' AND county_name !~ 'Wash'
4 ORDER BY county_name;
```

WHERE col ~ 'pattern' does not match NULL rows.

## 4.5 Checkpoint: alternation

**Task:** Return rows from `county_regex_demo` where `county_name` matches `(lade|lare)` case-insensitively. Select `county_name` only, ordered by name.

**Expected:** At least **Clare County**.

**Hint:** One WHERE with `~*` and parentheses.

## 4.6 Solution: alternation

```
1 SELECT county_name
2 FROM county_regex_demo
3 WHERE county_name ~* '(lade|lare)'
4 ORDER BY county_name;
```

## 5 From blob to columns

### 5.1 The crime narrative shape

crime\_reports stores original\_text blobs. Dates, streets, offense labels, and case numbers sit inside the text, not in typed columns yet.

**Run this.**

```
1 SELECT original_text
2 FROM crime_reports
3 ORDER BY crime_id;
```

### 5.2 regexp\_match vs regexp\_matches

- **regexp\_match**: first match as text[], or NULL
- **regexp\_matches(..., 'g')**: one row per non-overlapping match (set-returning)

**Run both** on crime\_id order and compare row 1 (two dates in the narrative).

```
1 SELECT crime_id,
2         regexp_match(original_text, '\d{1,2}/\d{1,2}/\d{2}')
3 FROM crime_reports
4 ORDER BY crime_id;
5
6 SELECT crime_id,
7         regexp_matches(original_text, '\d{1,2}/\d{1,2}/\d{2}', 'g')
8 FROM crime_reports
9 ORDER BY crime_id;
```

## 5.3 Capture groups and [1]

Parentheses **capture** substrings. Index with [1], [2], and so on when you want scalars.

```
1 SELECT crime_id,  
2     (regexp_match(original_text, '(?:C0|S0)[0-9]+')) [1] AS case_number  
3 FROM crime_reports  
4 ORDER BY crime_id;
```

(?: ... ) is a **non-capturing** group: it groups for precedence and repetition but does **not** add a numbered capture. That is why the case-number pattern uses (?:C0|S0) so the whole token still lands in [1].

If the match fails, the array is NULL and so is [1].

## 5.4 Capturing vs non-capturing (same URL)

Run these and compare the returned arrays.

```
1 -- Capturing: (s) is a numbered capture  
2 SELECT regexp_match('https://site.com', 'http(s)?://');  
3  
4 -- Non-capturing: (? :s) groups the optional s without a capture  
5 SELECT regexp_match('https://site.com', 'http(? :s)?://');
```

With **capturing** (s)?, PostgreSQL puts **only the captured substrings** in the result array (not the full match first). Here [1] is 's' when the URL uses **https**, because that is the first (and only) capture.

With **non-capturing** (? :s)?, there are no numbered captures left, so `regexp_match` returns a **one-element** array: the whole match, 'https://'.

## 5.5 How about this?

Run this for plain http:

```
1 SELECT regexp_match('http://site.com', 'http(s)?://');    -- first capture empty: NULL  
2 SELECT regexp_match('http://site.com', 'http(? :s)?://'); -- whole match: http://
```

## 5.6 Checkpoint: second date

**Task:** For each `crime_id`, return the **second** date in the narrative when it exists. Use `regexp_matches` with the `g` flag and pick the second element of the returned array, or use a lateral pattern you trust.

**Stretch (optional):** Only show rows where two dates exist.

**Hint:** `regexp_matches` with `g` returns one row per match; aggregate or filter on `ordinality` if you use `WITH ORDINALITY` in a lateral join, or take the second row in a subquery.

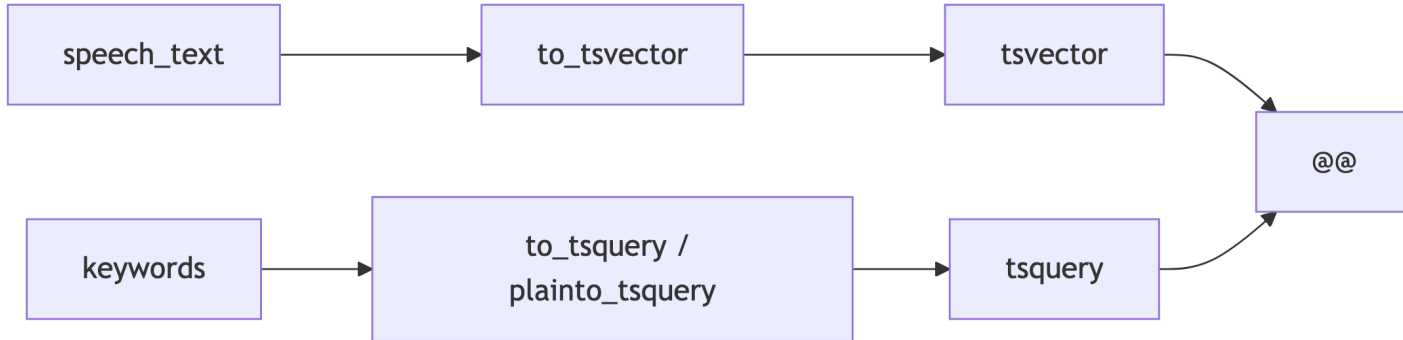
## 5.7 Solution: second date

```
1 SELECT c.crime_id,
2         m.match[1] AS second_date
3 FROM crime_reports AS c
4 JOIN LATERAL (
5     SELECT match
6     FROM regexp_matches(
7         c.original_text,
8         '\d{1,2}/\d{1,2}/\d{2}',
9         'g'
10    ) WITH ORDINALITY AS t(match, ord)
11 WHERE ord = 2
12 ) AS m ON true
13 ORDER BY c.crime_id;
```

Rows with only one date in the narrative do not appear in this inner join result. `LEFT JOIN LATERAL` plus `NULL` would keep those rows if you need them.

## 6 Full-text search

### 6.1 Pipeline



Store or compute a **tsvector** for documents; build a **tsquery** for the question; @@ is the match predicate. **ts\_headline** formats snippets; **ts\_rank** orders by relevance. A **GIN** index on the **tsvector** column keeps @@ from scanning every row.

### 6.2 Types and functions (short)

Piece	Role
<code>to_tsvector('english', text)</code>	Tokenize, drop stop words, stem
<code>to_tsquery / plainto_tsquery</code>	Build a query (&, \ , !, <->, <N> in <code>to_tsquery</code> )
<code>@@</code>	<code>tsvector @@ tsquery</code>

`plainto_tsquery` is handy for simple keyword boxes; `to_tsquery` shows operators explicitly.

### 6.3 GIN in one minute

**GIN** is an inverted index: lexeme to row lists. It fits @@ on **tsvector**. B-trees do not. Writes cost more than raw text; bulk load then create index is a common pattern.

The setup script already adds `search_speech_text` and a GIN index on `president_speeches`.

## 6.4 Demo: match and snippet

Run this.

```
1 SELECT president,
2     speech_date,
3     ts_headline(
4         speech_text,
5         to_tsquery('english', 'tax'),
6         'StartSel = <, StopSel = >, MinWords=5, MaxWords=7, MaxFragments=1'
7     )
8 FROM president_speeches
9 WHERE search_speech_text @@ to_tsquery('english', 'tax')
10 ORDER BY speech_date
11 LIMIT 5;
```

Filter with @@ first; ts\_headline is presentation on matching rows.

## 6.5 Demo: rank with normalization

Run both and compare top five ordering.

```
1 SELECT president, speech_date,
2     ts_rank(search_speech_text, to_tsquery('english', 'war & security')) AS r0
3 FROM president_speeches
4 WHERE search_speech_text @@ to_tsquery('english', 'war & security')
5 ORDER BY r0 DESC
6 LIMIT 5;
7
8 SELECT president, speech_date,
9     ts_rank(search_speech_text, to_tsquery('english', 'war & security'), 2)::numeric AS r2
10 FROM president_speeches
11 WHERE search_speech_text @@ to_tsquery('english', 'war & security')
12 ORDER BY r2 DESC
13 LIMIT 5;
```

Flag 2 divides by document length so long speeches do not always win.

## 6.6 Checkpoint: economy and jobs

**Task:** Find speeches where **economy and jobs** both match (use to\_tsquery with &). Return president, speech\_date, and a ts\_headline on speech\_text for the query. Sort by ts\_rank with normalization flag 2 descending; limit 10 rows.

**Hint:** Single `to_tsquery('english', 'economy & jobs')` for filter, headline, and rank.

## 6.7 Solution: economy and jobs

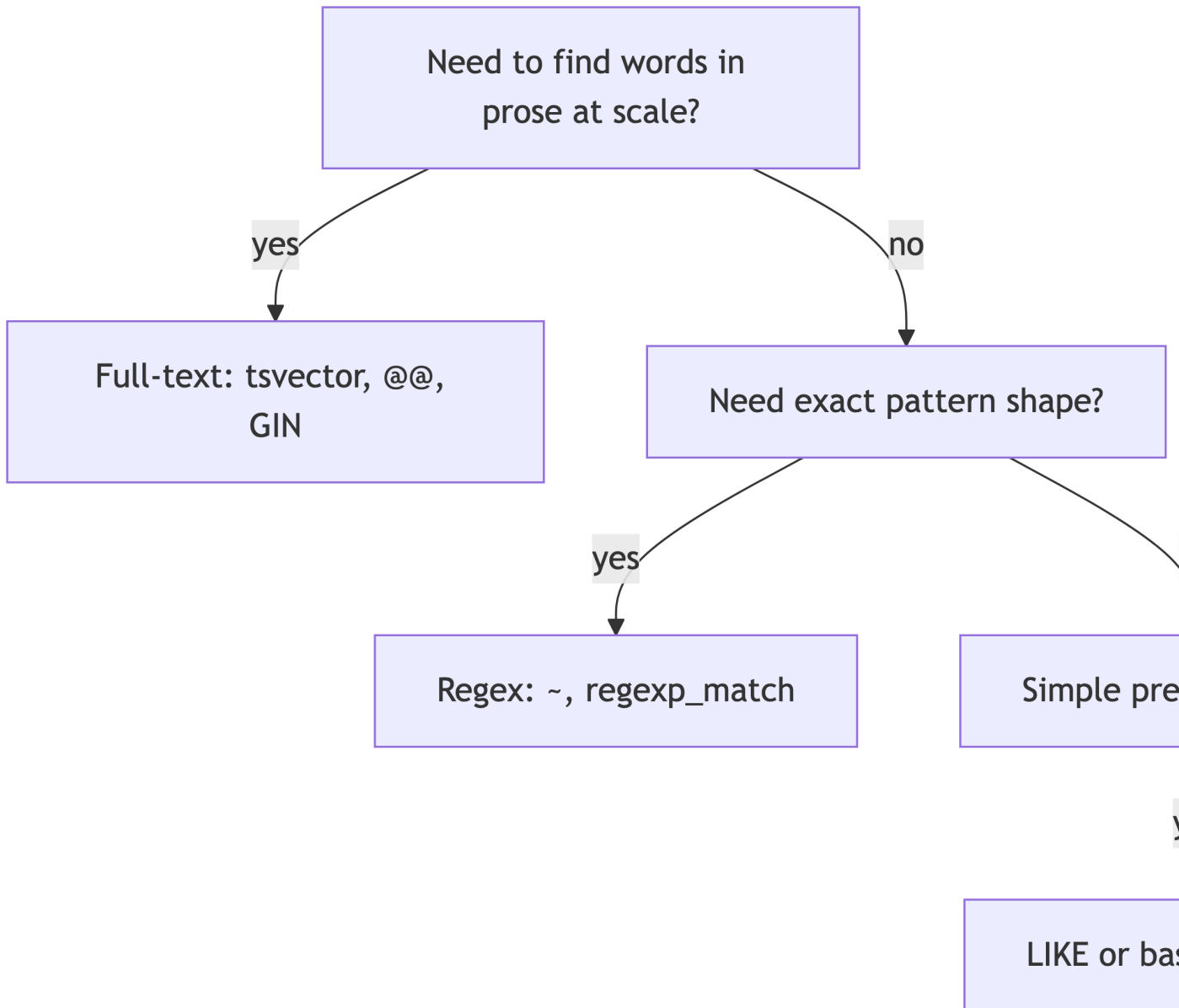
```
1 SELECT president,  
2     speech_date,  
3     ts_headline(  
4         speech_text,  
5         to_tsquery('english', 'economy & jobs'),  
6         'StartSel = <, StopSel = >, MinWords=5, MaxWords=10, MaxFragments=1'  
7     ),  
8     ts_rank(  
9         search_speech_text,  
10        to_tsquery('english', 'economy & jobs'),  
11        2  
12    )::numeric AS relevance  
13 FROM president_speeches  
14 WHERE search_speech_text @@ to_tsquery('english', 'economy & jobs')  
15 ORDER BY relevance DESC  
16 LIMIT 10;
```

## 7 Summary

### 7.1 Techniques in one place

Technique	Use it when
<code>trim, left, replace</code>	Normalize display and prep for parsing
<code>~ / ~*, regexp_match</code>	Filters and first structured pull from text
<code>regexp_matches(..., 'g')</code>	Every occurrence; drive lateral or counting
<code>to_tsvector + GIN</code>	Searchable document column
<code>@@, ts_headline, ts_rank</code>	Filter, snippet, sort

## 7.2 Choosing a tool



## 7.3 References

1. DeBarros, A. (2022). *Practical SQL* (2nd ed.). No Starch Press. Chapter 14.
2. PostgreSQL: [Pattern Matching](#)
3. PostgreSQL: [String Functions](#)

4. PostgreSQL: [Full Text Search](#)
5. PostgreSQL: [ts\\_rank](#)