

SQL: The Sequel

Christopher R. Genovese

Department of Statistics & Data Science

Tue 11 Nov 2025
Session #20

Plan

Recap

Plan

Recap

SQL Short Exercises

Plan

Recap

SQL Short Exercises

Schema Design Principles

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Announcements

- In **documents**: lecture notes, sql files, lec20-ex.org
- Download postgresql
- **Reading:**
 - Databases
 - pgexercises.com
 - sqlbolt.com
 - sql-practice.com – Structured Practice
 - mystery.knightlab.com – An SQL Murder Mystery
 - selectstarql.com
 - Fun with FP
 - cube composer
- **Homework:**
 - **sym-spell** assignment due Mon 17 Nov.

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

SQL Commands

CRUD Operations

- **CREATE TABLE**

```
create table NAME (attribute1 type1, attribute2 type2, ...);
```

Creates a table by name, giving names and attributes of all columns, and specifying any constraints.

See [documentation](#).

- **ALTER TABLE**

```
alter table [ if exists ] [ only ] <NAME> [ * ] <ACTION> [, ... ]
```

Change table features: adding and removing columns, renaming attributes, changing constraints or attribute types, and setting column defaults.

See [documentation](#).

- **DROP TABLE**

```
drop table [if exists] <NAME> [cascade | restrict];
```

Permanently eliminate a table, and optionally those that depend on it (cascade).

Use with care.

See [documentation](#).

SQL Commands (cont'd)

- **INSERT**

```
insert into <NAME> (<COLUMN1>, ..., <COLUMNK>)  
    values (<VALUE1>, ..., <VALUEK>);
```

Insert values into a table. A `returning` clause makes this return a table of results, such as the ids of inserted rows.

See [documentation](#).

- **\copy**

```
\COPY <NAME> from 'FILE.csv' with delimiter ',';  
select setval('NAME_id_seq', 1001, false);
```

A meta-command for copying a CSV file into a table.

See also [copy command](#).

- **UPDATE**

```
update <TABLE>  
    set <COL1> = expression1,  
        <COL2> = expression2,  
        ...  
    where <CONDITION>;
```

Sets the value of columns where condition is satisfied. **Do not omit the `where clause`.**

See [documentation](#).

SQL Commands (cont'd)

- SELECT

```
select <COLS> from <TABLE> [where <CONDITION>] [...]
```

Retrieves rows from a table or view. The workhorse query! Lots of variations. Here * is a shorthand for all columns. Column spec can include table prefixes, as clauses and more. See [documentation](#).

Examples:

```
select * from events;
select * from events where id > 20 and id < 40;
select score, element from events
    where persona = 1202
    order by element, score;

select 1 as one;
select ceiling(10 * random()) as r;
select 1 as ones from generate_series(1,10);

select min(r), avg(r) as mean, max(r) from
    (select random() as r from generate_series(1, 10000)) as _;
select '2015-01-22 08:00:00'::timestamp + random() * '64 days'::interval
    as w from generate_series(1, 10);
```

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Short Exercises

See associated file.

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Database Schema Design Principles

The key design principle for database schema is to keep the design DRY – that is, **eliminate data redundancy**. The process of making a design DRY is called **normalization**, and a DRY database is said to be in "normal form."

The basic modeling process:

- ① Identify and model the entities in your problem
- ② Model the relationships between entities
- ③ Include relevant attributes
- ④ Normalize by the steps below

Example

Consider a database to manage songs:

Album	Artist	Label	Songs
Talking Book	Stevie Wonder	Motown	You are the sunshine of my life, Maybe your baby, Superstition, ...
Miles Smiles	Miles Davis Quintet	Columbia	Orbits, Circle, ...
Speak No Evil	Wayne Shorter	Blue Note	Witch Hunt, Fee-Fi-Fo-Fum, ...
Headhunters	Herbie Hancock	Columbia	Chameleon, Watermelon Man, ...
Maiden Voyage	Herbie Hancock	Blue Note	Maiden Voyage
American Fool	John Cougar	Riva	Hurts so good, Jack & Diane, ...
...			

This seems fine at first, but why might this format be problematic or inconvenient?

Example (cont'd)

- Difficult to get songs from a long list in one column
- Same artist has multiple albums
- "Best Of" albums
- A few thoughts:
 - What happens if an artist changes names partway through his or her career (e.g., John Cougar)?
 - Suppose we want mis-spelled "Herbie Hancock" and wanted to update it. We would have to change every row corresponding to a Herbie Hancock album.
 - Suppose we want to search for albums with a particular song; we have to search specially within the list for each album.

Step 1. Give Each Entity a Unique Identifier

The schema here can be represented as

Album (artist, name, record_label, song_list)

where Album is the *entity* and the labels in parens are its *attributes*.

To normalize this design, we will add new entities and define their attributes so **each piece of data has a single authoritative copy**.

This will be its primary key, and we will call it id here.

Key features of a primary key are that it is unique, non-null, and it never changes for the lifetime of the entity.

Step 2. Give Each Attribute a Single (Atomic) Value

What does each attribute describe? What attributes are repeated in `Albums`, either implicitly or explicitly?

Consider the relationship between albums and songs. An album can have one or more songs; in other words, the attribute `song_list` is non-atomic (it is composed of other types, in this case a list of text strings). The attribute describes a collection of another entity – `Song`.

So, we now have two entities, `Album` and `Song`. How do we express these entities in our design? It depends on our model. Let's look at two ways this could play out.

Step 2. Give Each Attribute a Single (Atomic) Value

- ① Assume (at least hypothetically) that each song can only appear on *one* album. Then Album and Song would have a **one-to-many** relationship.

- `Album(id, title, label, artist)`
- `Song(id, name, duration, album_id)`

Question: What do our CREATE TABLE commands look like under this model?

- ② Alternatively, suppose our model recognizes that while an album can have one or more songs, a song can also appear on one or more albums (e.g., a greatest hits album). Then, these two entities have a **many-to-many** relationship.

This gives us two entities that look like:

- `Album(id, title, label, artist)`
- `Song(id, name, duration)`

This is fine, but it doesn't seem to capture that many-to-many relationship. How should we capture that?

- An answer: This model actually describes a *new* entity – Track.

The schema looks like:

- `Album(id, title, label, artist)`
- `Song(id, name, duration)`
- `Track(id, song_id, album_id, index)`

Step 3. Make All Non-Key Attributes Dependent Only on the Primary Key

This step is satisfied if each non-key column in the table serves to *describe* what the primary key *identifies*.

Any attributes that do not satisfy this condition should be moved to another table.

In our schema of the last step (and in the example table), both the artist and label field contain data that describes something else. We should move these to new tables, which leads to two new entities:

- Artist(id, name)
- RecordLabel (id, name, street_address, city, state_name, state_abbrev, zip)

Each of these may have additional attributes. For instance, producer in the latter case, and in the former, we may have additional entities describing members in the band.

Step 4. Make All Non-Key Attributes Independent of Other Non-Key Attributes

Consider RecordLabel. The state_name, state_abbrev, and zip code are all non-key fields that depend on each other. (If you know the zip code, you know the state name and thus the abbreviation.)

This suggests to another entity State, with name and abbreviation as attributes. And so on.

Exercise

Convert this normalized schema into a series of CREATE TABLE commands.

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Activity

Building a database with your data.

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Aggregate Functions and Grouping

Aggregate functions operate on one or more attributes to produce a summary value.

Examples: count, max, min, sum.

For example, let's calculate the average score obtained by students:

```
select avg(score) from events;
```

This produces a single row: the average score. Any aggregate function takes many rows and reduces them to a single row. This is why you **can't** write this:

```
select persona, avg(score) from events;
```

Try it; why does Postgres complain?

Aggregate Functions and Grouping

We often want to apply aggregate functions not just to whole columns but to **groups of rows** within columns. This is the province of the GROUP BY clause. It groups the data according to a specific value, and aggregate functions then produce a single result **per group**.

For example, if I wanted the average score for each separate user, I could write:

```
select persona, avg(score) as mean_score  
from events  
group by persona  
order by mean_score desc;
```

You can apply conditions on grouped queries. Instead of WHERE for those conditions, you use HAVING, with otherwise the same syntax. Short version: WHERE selects rows, and HAVING selects groups.

```
select persona, avg(score) as mean_score  
from events  
where moment > '2020-10-01 11:00:00'::timestamp  
group by persona  
having avg(score) > 50  
order by mean_score desc;
```

Common Table Expressions

Common Table Expressions, introduced via `with`, introduces temporary tables that can be used for a query.

Simple example:

```
with a as (
    select x from generate_series(1, 10) as x
), b as (
    select y from generate_series(11, 20) as y
)
select * from a, b;
```

Common Table Expressions (cont'd)

Subsequent terms in with can refer to earlier terms. What does this do?

```
WITH regional_sales AS (
    SELECT region, SUM(amount) AS total_sales
    FROM orders
    GROUP BY region
), top_regions AS (
    SELECT region
    FROM regional_sales
    WHERE total_sales > (SELECT SUM(total_sales)/10 FROM regional_sales)
)
SELECT region,
       product,
       SUM(quantity) AS product_units,
       SUM(amount) AS product_sales
FROM orders
WHERE region IN (SELECT region FROM top_regions)
GROUP BY region, product;
```

Common Table Expressions (cont'd)

CTEs can be *recursive*.

Try this:

```
WITH RECURSIVE t(n) AS (
    VALUES (1)          -- Base case
    UNION ALL
    SELECT n+1 FROM t WHERE n < 100  -- Inductive Case
)
SELECT sum(n) FROM t;
```

Common Table Expressions (cont'd)

CTEs can be fairly powerful. Challenge CTE for Binary Tree traversal (DFS, BFS).

Common Table Expressions (cont'd)

CTEs can be fairly powerful. Challenge CTE for Binary Tree traversal (DFS, BFS).

```
create table tree (
    id integer,
    data integer,
    left_child integer,
    right_child integer
);
insert into tree (id,,data, left_child, right_child)
values (1, 10, 2, 3),
       (2, 20, 4, 5),
       (3, 30, 6, 7),
       (4, 40, 8, NULL),
       (5, 50, NULL, 9),
       (6, 60, NULL, NULL),
       (7, 70, NULL, NULL),
       (8, 80, NULL, NULL),
       (9, 90, NULL, NULL);
```

Common Table Expressions (cont'd)

CTEs can be fairly powerful. Challenge CTE for Binary Tree traversal (DFS, BFS).

```
WITH RECURSIVE search_tree(id, data, left_child, right_child, path) AS (
    SELECT t.id, t.data, t.left_child, t.right_child, ARRAY[t.id]
    FROM tree t
    WHERE t.id = 1 -- Root is base case
  UNION ALL
    SELECT t.id, t.data, t.left_child, t.right_child, path || t.id
    FROM tree t, search_tree st
    WHERE t.id = st.left_child OR t.id = st.right_child
)
SELECT * FROM search_tree ORDER BY path;
```

Common Table Expressions (cont'd)

CTEs can be fairly powerful. Challenge CTE for Binary Tree traversal (DFS, BFS).

```
WITH RECURSIVE search_tree(id, data, left_child, right_child, path, depth) AS (
    SELECT t.id, t.data, t.left_child, t.right_child, ARRAY[t.id], 0
    FROM tree t
    WHERE t.id = 1 -- Root is base case
  UNION ALL
    SELECT t.id, t.data, t.left_child, t.right_child, path || t.id, depth + 1
    FROM tree t, search_tree st
    WHERE t.id = st.left_child OR t.id = st.right_child
)
SELECT * FROM search_tree ORDER BY depth;
```

Plan

Recap

SQL Short Exercises

Schema Design Principles

SQL Activity

Variations

Issues on SQL in Code

Storing Your Password

The code last time stores your password right in the source file. This is a **bad idea**.

If the code is ever shared with anyone, posted online, or otherwise revealed, anyone who sees it now has your database username and password and can view or modify any of your data.

If you commit the file to Git, your password is now in your Git history **forever**.

Fortunately, there are ways to work around this.

Storing Your Password: R

Run this R code:

```
file.edit(file.path("~", ".Rprofile"))
```

This will create a file called `~/.Rprofile` in your home directory and open it for editing.
In this file, write something like

```
DB_USER <- "yourusername"  
DB_PASSWORD <- "yourpassword"
```

Save and close the file. Start a new R session. The `DB_USER` and `DB_PASSWORD` variables will be defined in *any R script you run*, so you can use them in your code.

And since the `.Rprofile` is not in your assignments repository, you won't accidentally commit it to your Git history.

Storing Your Password: Python

Python doesn't have something like `~/.Rprofile`. Instead, when you do an assignment that requires SQL access, create a separate file `credentials.py` defining your username and password in variables.

You can import `credentials` and then use `credentials.DB_USER` and `credentials.DB_PASSWORD` in your other code files.

To avoid accidentally committing `credentials.py`, create (or modify, if it exists) a file called `.gitignore` in the root of your assignments repository.

Add the line `credentials.py` to it. This will make Git ignore any files called `credentials.py`, so you don't accidentally commit them.

You can use the `cryptography` or `bcrypt` packages to encrypt your stored password data as well, though the keys will need to be accessible in most cases.

SQL in Unit Tests

Question: If you have secret SQL credentials, but your unit tests need to check functions that access the database, how can your tests run on our "public" server? How can anyone else run your tests?

Answer: They can't. This is fine for our class; in a business environment, you would use a secret-management system to ensure code can always access the passwords it needs, or set up special test databases for development. But we're keeping things simple here.

So instead you should do the following:

- ① Ensure SQL access is limited to the functions that really need it. Most of your code doesn't need to access SQL – it just needs to be provided the right data in a convenient form. Only a few functions need to run SQL queries and put the data into the right form for the rest of your code to use. If you follow this design principle, most of your functions can still be tested.
- ② Set your tests to be skipped if the credentials are not available.

For point 2, most unit testing frameworks provide ways of marking tests as being skipped in certain circumstances.

SQL in Unit Tests

For example, in Python:

```
import pytest
import psycopg2

try:
    import credentials
except ImportError as e:
    no_database = True
else:
    no_database = False

@pytest.mark.skipif(no_database, reason="No database credentials available")
def test_some_database_thingy():
    conn = psycopg2.connect(...)

    # do stuff
```

SQL in Unit Tests

In R:

```
library(testthat)

test_that("database access works", {
  skip_if_not(exists("DB_USER"))

  # do some database stuff
})
```

Practicing Safe SQL

Suppose you've loaded some data from an external source – a CSV file, input from a user, from a website, another database, wherever. You need to use some of this data to do a SQL query.

```
result <- dbSendQuery(paste0("SELECT * FROM users ",  
                           "WHERE username = '", username, "' ",  
                           "AND password = '", password, "'"))
```

Now suppose `username` is the string `''; DROP TABLE users;--'`. What does the query look like before we send it to Postgres?

Practicing Safe SQL

Suppose you've loaded some data from an external source – a CSV file, input from a user, from a website, another database, wherever. You need to use some of this data to do a SQL query.

```
result <- dbSendQuery(paste0("SELECT * FROM users ",  
                           "WHERE username = ''", username, "' ",  
                           "AND password = ''", password, "''))
```

Now suppose `username` is the string `''; DROP TABLE users;--`. What does the query look like before we send it to Postgres?

```
SELECT * FROM users  
WHERE username = ''; DROP TABLE users; -- AND password = 'theirpassword'
```

Practicing Safe SQL (cont'd)

We have *injected* a new SQL statement, which drops the table. Because -- represents a comment in SQL, the commands following are not executed.



(source: [xkcd](#))

Practicing Safe SQL (cont'd)

Less maliciously, the username might contain a single quote, confusing Postgres about where the string ends and causing syntax errors. Or any number of other weird characters....

Clever attackers can use SQL injection to do all kinds of things: imagine if the password variable were `foo' OR 1=1` – we'd be able to log in without knowing the right password!

We need a better way of writing queries. Database systems provide *parametrized queries*, where you are explicit about input parameters that are not SQL syntax. For example:

```
username <- "'"; DROP TABLE users;--"  
password <- "walruses"  
query <- sqlInterpolate(con,  
                        "SELECT * FROM users WHERE username = ?user AND password = ?  
                        user = username, pass = password)  
users <- dbGetQuery(con, query)
```

Strings of the form `?var` are replaced with the corresponding `var` in the arguments, but with any special characters escaped so they do not affect the meaning of the query.

In this example, `query` is now

```
SELECT * FROM users WHERE username = '''; DROP TABLE users;--'  
AND password = 'walruses'
```

The single quote at the beginning of `username` is doubled there: that's a standard way of escaping quotation marks, so Postgres recognizes it's a quote inside a string, not the boundary of the string.

Practicing Safe SQL (cont'd)

psycopg2 provides similar facilities:

```
cur.execute("SELECT * FROM users "
            "WHERE username = %(user)s AND password = %(pass)s",
            {"user": username, "pass": password})
```

You should *always* use this approach to insert data into SQL queries. You may think it's safe with your data, but at the least opportune moment, you'll encounter *nasal demons*.

THE END