



Although there are many reasons to shift from using Excel to databases, the most important cornerstones are: (i) Data Integrity, (ii) Redundancy, (iii) Error prone, (iv) User Access and Security and (v) Data Accessibility and Speed.

Data Integrity:

- There are a set of rules that govern the structure of the data, how different information relate to one another and also what input can be put into a certain column of a table (ex. numeric vs string).
- These rules and protocols build a general framework under which everyone must adhere to, and as such, increases reliability of the data.
- It also removes questions like "What do I do when I get new data? Where do I store that data? How must the data look?".
- Excel is limited to 1,000,000 rows of data in a single sheet. If you operating at close to the edge of that amount of data how difficult to you think it is to ensure the quality of the information?



Redundancy:

- I think all of us know version control of files that are called final_analytics_tom_v2.3_client_clean.xlsx. These files eventually just become copies of the same data with small changes that are difficult to track and keep clean.
- Using relational databases also ensures that we separate information out across tables to ensure we do not have multiple versions of the same data in different places. Example would be to separate out customers and their purchases.
 - When updating a customer's information, we only have to update the customers table and not the purchases tables, as the purchases table only links back to the customer via a unique id (also called a key) of some sort.

Error prone:

- Excel sheets (as we saw earlier) is very much susceptible to proliferation of errors, especially when the data gets large. There is no way of know what changed and how when someone accidentally overwrites a cell/row/column.
- Also, because sheets are usually linked together, if the sheet changes unexpectedly (perhaps someone added a column), suddenly the formulas are no longer correct or links to the data are broken.



Access and Security:

- Multiple users can work on a database at the same time. Most of the time end-users only need to collect or aggregate data for their purposes, which means we can create specials user that are only allowed those operations. This ensures that the risk of unknowledgeable users do not accidentally change the underlying data without them knowing.
- With personal information protection laws coming into force, one also wants to restrict access to certain types of information by only granting access to those who have the right clearance¹.

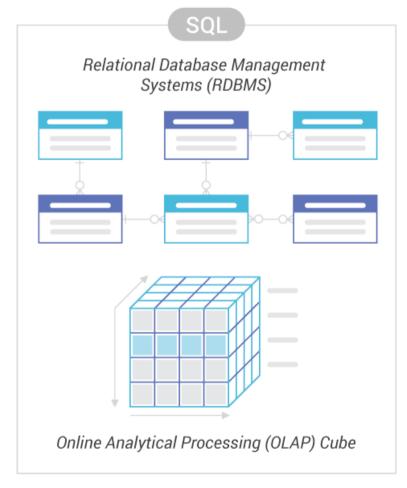
Data Accessibility and Speed:

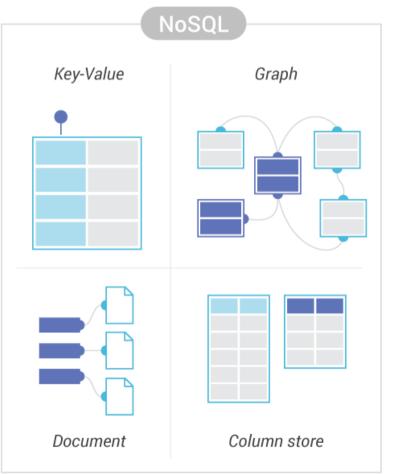
- Ain't no one got time to work on an Excel sheet over 10,000 rows. Having to quickly analyse information using aggregation tools such as *pivot* or *vlookups* becomes a total nightmare. Excel is dynamic, which means every action causes all the information to automatically recalculate.
- Databases allows you to do deep analysis of data over millions (even billions) of rows of data in seconds.
- Because you most likely will be using a relational setup in a database, you only query the data you need, not the whole information set.

¹ Please do NOT ever send personal information in plain text in Excels spreadsheets over email, unless you want to end up in jail. 💀

Different database types



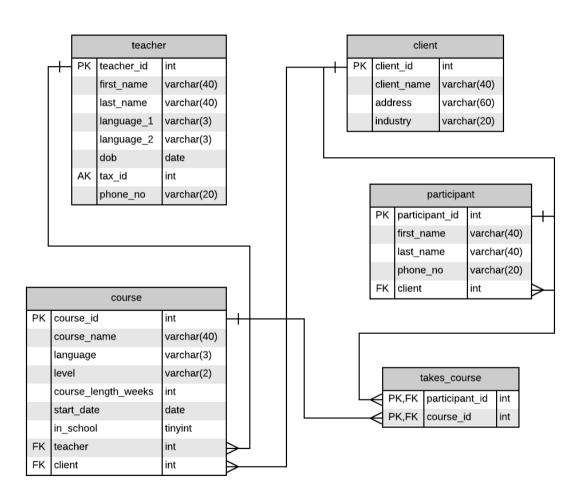




^{*}link to original website

RDMS database example data schema





^{*}link to original website

Different database structures



Do note, that although we do not cover OLAP database in this course, they tend to be a little bit different as they try to avoid complex joins which could slow down analytics.

Database structures, or *schema* design, depends a lot on the application of the database. Although there are different schemas and designs, they do have some common traits:

- Includes the name of the fields in the table.
- The type that the field consists of (numeric, date, varchar etc.).
- Associations and keys linking fields.

Common schemas that you might encounter are:

- Star schema
- Snowflake schema
- Fact constellation

Different database structures



To understand schemas a bit better, we need an understanding of the pieces. The two most important components consist of: fact and dimension tables.

FACT:

- Fact table contains measurements, metrics, and facts about a business process
 - ex. sales or webpage visits.
- Fact tables form the primary table of the design and are usually normalized
 - We assign a numerical number or code to an attribute for better performance
 - An example of this would be where we code GENDER as 1 for male and 2 for female.

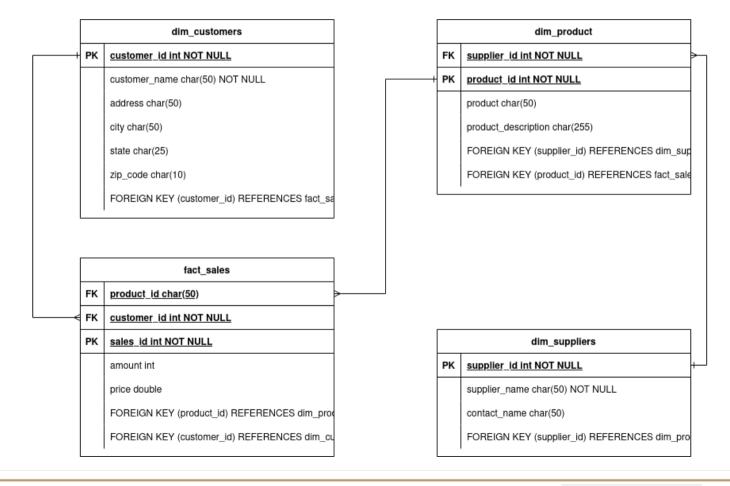
DIMENSION:

- Provides the information about the facts
 - ex. *location* of transaction, *customer*
- These tables are de-normalized and have to be joined to the fact table table before analysis can happen.
- The tables also usually contain descriptions of the field in order to make it easier to understand.

Different database structures: Star Schema



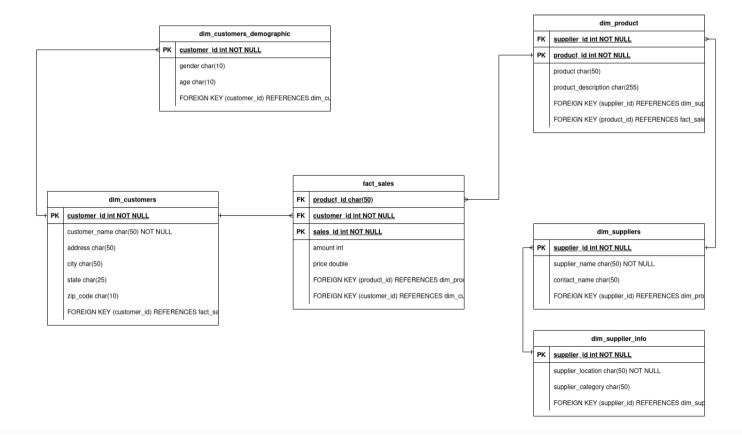
Star schema design has only *one* fact table and multiple dimension tables. This is a very common design as its relational properties are easily understood.



Different database structures: Snowflake Schema



Snowflake schema design extends the Star Schema by only *one* fact table, multiple dimension tables, each with their own dimension tables. This adds another layer of abstration to the dimension tables and could contain additional information about attributes not in use every day.





Basic SQL 🛠

Jumping in with both feet



In an advanced Database course you should cover databases, its setup and optimization in much more depth. But for now, lets focus on getting you writing SQL queries and leave the database setup to the DBA.

To learn SQL, one needs to have a database. I have partitioned a DB for each of you! We are going to need mostly functions for this exercise: db_query and db_write. These functions are going to have the same format:

- 1) Connect to DB using variables from the Environment
- 2) Either read or write and then
- 3) ...make sure that if the function exits, we disconnect



... dont worry if you feel like this ...

Jumping in with both feet





Tip: In copy paste this! No need to reinvent the wheel 💡



Step 0: Specify your project variables in .Renviron, usethis::edit_r_environ()

```
gp data = datascience
gp user = datascience
gp pass = datascience
gp host = XXX
gp port = 600X
```

Step 1: Connect to DB using variables from the Environment using dbbasic

```
remotes::install github("HanjoStudy/dbbasic")
library(dbbasic)
conn ← db connect(db = "psql datascience")
DBI::dbDisconnect(conn)
```

What is dbbasic?



dbbasic is a little passion project from when I worked at 71point4. You can find it HERE

The aim is to assist and facilitate R programmers with interacting with databases without too much cognitive overhead:

- Read
- Write
- Connect

It also has some very nice utilities:

- db_collapes if you need to collapse a vector for an IN style query
- hash creates a SHA256 hash
- tibble_to_sql helps with boilerplate

https://github.com/HanjoStudy/dbbasic/tree/master



Rstudio + dbbasic?

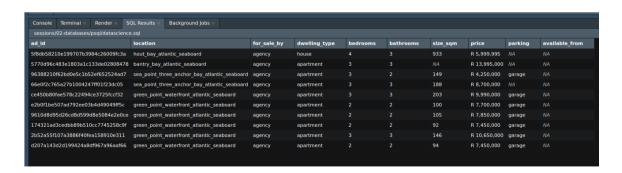


To use dbbasic in R, create a new dev.sql file in your project and add -- !preview conn=db_connect(db = "psql_datascience") to the top of the file.

For this to work your .Renviron file in your project directory must be set up correctly!

```
-- !preview conn=db_connect(db = "psql_datascience")
SELECT * FROM gumtree LIMIT 10
```

And then press CTRL + SHIFT + ENTER for magic 🧙



^{*}You can also do some very clever things with dbbasic + Quarto + SQL chunks ... but thats for another day.

The KING of all statements: SELECT



The way to think about SQL is in terms of english commands. Also, also start from the inside and work your way out (you will see what I mean). It is good practice to always end your statements with LIMIT 10 until you are sure that the correct results is returned. Working on billion row tables and forgetting to limit your results can crash tables.

Lets start with the two statements you will most likely use every day:

• Counting how many rows there are

```
SELECT COUNT(*) FROM gumtree LIMIT 10;
db_query("SELECT COUNT(*) FROM gumtree LIMIT 10;", db = "psql_datascience")
```

• Getting a 10 row sample

```
SELECT * FROM gumtree LIMIT 10;
db_query("SELECT * FROM gumtree LIMIT 10;", db = "psql_datascience")
```

The KING of all statements: SELECT



Previously I decided I wanted to return all the columns (*), but what if I only want to return one or two of the columns?

```
SELECT {column1}, {column2} FROM table LIMIT 10;
Lets only return dwelling_type, size_sqm and price:

SELECT dwelling_type, size_sqm, price FROM gumtree LIMIT 10;
```

de It is good practice to not have long SQL statements in one row.

Code Needs a lot of whitespace.

That is how it breaths

Roger Peng, Jenny Bryan, useR 2018

The KING of all statements: SELECT



Lets build a bigger SELECT statement (I like 3 tab indentation):

```
SELECT
    dwelling_type,
    bedrooms,
    bathrooms,
    parking,
    size_sqm,
    price
FROM gumtree
LIMIT 10;
```

SELECT but with filter criteria



What happens if we only want to return an ad of a certain type?

Well, then we can employ the WHERE statement. We are going to collect the same columns as previously, but now we will specify the WHERE criteria on dwelling_type column where equal to house:

```
SELECT
    dwelling_type,
    bedrooms,
    bathrooms,
    parking,
    size_sqm,
    price
FROM gumtree
WHERE
    dwelling_type = 'house'
LIMIT 10;
```

Quick Practice:, write the code to bring back 100 examples where there is parking and ORDER BY price. TIP: Please google "psql filter NOT NULL on column"

SELECT but with filter criteria and order



In certain circumstances, it is necessary to order your data to get the correct output. For instance if we want to get the top 10 largest size_sqm houses:

```
SELECT
      dwelling type,
      bedrooms,
      bathrooms,
      parking,
      size sqm,
      price
FROM gumtree
WHERE
      dwelling type = 'house'
      AND size sqm IS NOT NULL
ORDER BY
      size sqm DESC
LIMIT 10
```

Lets clean and upload a clean version!



By now you would have seen that price and size_sqm and price is not in the correct format... Lets fix that:

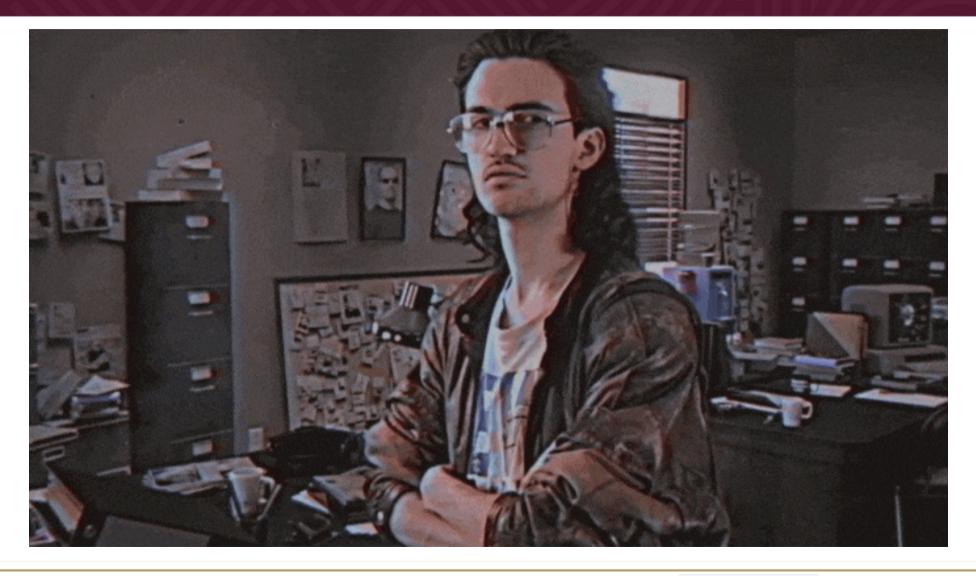
- 1) Pull whole database into R, 2) Fix size_sqm and price, 3) Upload to new table called gumtree_clean
- (1) Pull whole database into R & (2) Fix size_sqm and price

(3) Upload to new table called gumtree_clean

```
db_write(gumtree_clean, "gumtree_clean", db = "psql_datascience")
db_query("DROP TABLE gumtree_clean", db = "psql_datascience")
db_write(gumtree_clean, "gumtree_clean", db = "psql_datascience")
db_query("SELECT * FROM gumtree LIMIT 10;", db = "psql_datascience")
```

WELL DONE!!!





Time for security 🔒



I have set your password for you, but given that you are a super user, you can now make your password whatever you want!

• This will also ensure that someone doesn't mess with your sever

ALTER USER datascience WITH PASSWORD 'XXX';

Aggregations (Pivoting) in SQL



Pivoting forms part of the aggregation function of SQL. This helps us answer questions like:

- What is the average price of houses by dwelling_type?
- Total value and volume per location?

As you can see, aggregations or GROUP BY clauses gets used OFTEN, so learn it well and get comfortable with it.

Aggregations (Pivoting) in SQL



What is the average price of houses by dwelling_type?

```
dwelling_type,
    AVG(price) as mean_price

FROM gumtree_clean

WHERE
    price IS NOT NULL

GROUP BY
    dwelling_type

ORDER BY
    AVG(price) DESC

LIMIT 10
;
```

Once you are comfortable that you have the query correctly specified, drop the LIMIT and scroll through your magnificent piece of work!

Aggregations (Pivoting) in SQL



Total value and volume per location?

```
SELECT
    location,
    COUNT(*) AS volume,
    SUM(price) AS value

FROM
    gumtree

GROUP BY
    location
LIMIT 10
;
```

Notice how I ALIAS my aggregations as {aggregation} then {name}. This will make your life a lot easier and in some case it is mandatory... as in joins.

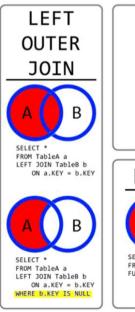
```
db querv("
          SELECT
                location.
                COUNT() AS volume.
                SUM(price) AS value.
                SUM(price)/COUNT() avg price
          FROM
                gumtree clean
          GROUP BY
                location
          ORDER BY
                SUM(price) DESC
          LIMIT 10
          ;", db = "psql datascience")
```

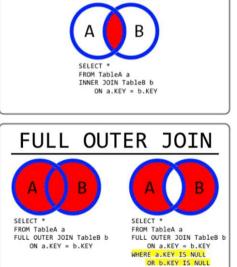


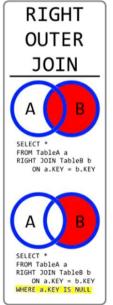
By now you are asking yourself, if we designed our database in the beautiful star schema that we talked about earlier, how do we *join* all the information together again? This is where JOINS come in and there are a multitude of them. Most important one is LEFT JOIN and INNER JOIN:

SQL JOINS

INNER JOIN









Lets attempt a basic join before we combine joins with aggregations. To start off we will JOIN the gumtree_descriptions table onto the gumtree clean table:

• First check what tables are available:

```
db_query("SELECT * FROM pg_catalog.pg_tables WHERE schemaname = 'public'", db = "psql_datascience")
```

• Next get a sample of rows from gumtree_descriptions:

```
db_query("SELECT * FROM gumtree_description LIMIT 10", db = "psql_datascience")
```

From this we see that gumtree_clean.ad_id = gumtree_descriptions.ad_id!



Now that we have the specific keys from the different tables, there are two ways to join: (1) When name is the same in both table, (2) When they are different:

When name is the same in both table

```
SELECT
  *
FROM gumtree_description
LEFT JOIN(
   SELECT * FROM gumtree_clean
) tbl_clean
USING(ad_id);
```

When they are different

```
SELECT
  *
FROM gumtree_descriptions
LEFT JOIN(
   SELECT * FROM gumtree_clean
) tbl_clean
ON gumtree_clean.ad_id = gumtree_descriptions.ad_id
;
```

• Important to note the ALIAS of the inner table called tbl_clean.



There are ways to optimize your joins to be extremely fast. Although we do not cover these in this course, it is interesting and worth knowing none the less.

- One is keys (which is why we use primary and muli keys in tables).
- Another is query optimization through column selection and subqueries.
 - Although we do not cover these in this course, having knowledge of advance backend mechanics can sometimes take your execution time from days to minutes.

Bonus Round: CTEs





Common Table Expressions



This is the %>% equivalent for SQL. The structure is quite simple:

```
WITH cte_one AS(
... some SQL ...
), cte_two AS(
... do something with cte_one ...
)
SELECT * FROM cte_two
```

These are quite 'advanced' topics, but can promise you once you start using them you wont go back to subqueries¹.

Now lets move on to asking practical questions

What are the top words associated with a property advert?

¹If you want to know why I say that, go read up on subqueries.



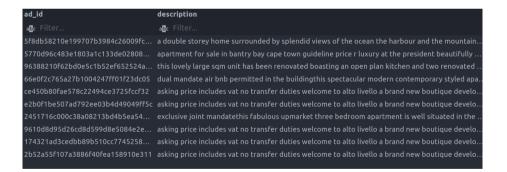
Case Study



Step 1: I know we are going to need to remove useless words (stopwords) such as "a", "the" etc. So lets upload a stopwords table to our DB:

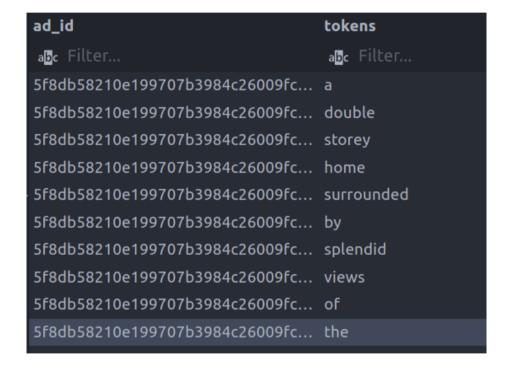
```
db_write(tidytext::stop_words, table_name = "stopwords", db = "psql_datascience")
```

Step 2: Cleaning the text of puncuation:





Step 3: Tokenise the long text into words





Step 4: Remove stopwords and GROUP BY count words per ad_id

```
WITH cte one AS(
    SELECT ad id,
        REGEXP_REPLACE(LOWER(description), '[^A-Za-z ]', '', 'g') AS
    FROM gumtree description
), cte two AS(
    SELECT ad id,
       TRIM(REGEXP SPLIT TO TABLE(LOWER(description), '\s+')) AS to
    FROM cte one
), cte three AS(
    SELECT *
    FROM cte two
    WHERE tokens NOT IN (SELECT DISTINCT word FROM stopwords)
), cte four AS(
    SELECT ad id,
       tokens,
       COUNT(*) as obs,
        ROW_NUMBER() OVER (
            PARTITION BY ad id
            ORDER BY COUNT(*) DESC
        ) as rank
    FROM cte_three
    GROUP BY ad_id,
        tokens
SELECT * FROM cte four LIMIT 10
```

ad_id	tokens	obs	rank
aBc Filter	abc Filter	abc Filter	a b c Filte
000a2cd6656b24763bd1fc416ea01	family	4	1
000a2cd6656b24763bd1fc416ea01	plan		2
000a2cd6656b24763bd1fc416ea01	helena		3
000a2cd6656b24763bd1fc416ea01	bay		4
000a2cd6656b24763bd1fc416ea01	home		5
000a2cd6656b24763bd1fc416ea01	builtin		6
000a2cd6656b24763bd1fc416ea01	st		7
000a2cd6656b24763bd1fc416ea01	perfect	2	8
000a2cd6656b24763bd1fc416ea01	kitchen	2	9
000a2cd6656b24763bd1fc416ea01	coast	2	10
000a2cd6656b24763bd1fc416ea01	bedroom	2	11
000a2cd6656b24763bd1fc416ea01	west	2	12
000a2cd6656b24763bd1fc416ea01	de		13
000a2cd6656b24763bd1fc416ea01	decades		14
000a2cd6656b24763bd1fc416ea01	detailsada		15
000a2cd6656b24763bd1fc416ea01	disappointedproperty		16
000a2cd6656b24763bd1fc416ea01	doors	1	17



Final Step: Concatenate the token and observation and filter where less than 10

```
WITH cte one AS(
    SELECT ad id,
        REGEXP REPLACE(LOWER(description), '[^A-Za-z ]', '', 'g') AS description
    FROM gumtree description
). cte two AS(
    SELECT ad id,
       TRIM(REGEXP SPLIT TO TABLE(LOWER(description), '\s+')) AS tokens
    FROM cte one
). cte three AS(
    SELECT *
    FROM cte two
    WHERE tokens NOT IN (SELECT DISTINCT word FROM stopwords)
), cte four AS(
    SELECT ad id,
        tokens,
        COUNT(*) as obs,
        ROW NUMBER() OVER (
            PARTITION BY ad id
            ORDER BY COUNT(*) DESC
       ) as rank
    FROM cte_three
    GROUP BY ad id,
        tokens
), cte five AS(
    SELECT ad id,
       STRING AGG(tokens || '(' || CAST(obs AS VARCHAR) || ')',', ') AS top words
    FROM cte four
    WHERE rank ≤ 10
    GROUP BY ad id
SELECT * FROM cte five LIMIT 10;
```





VSCode ___

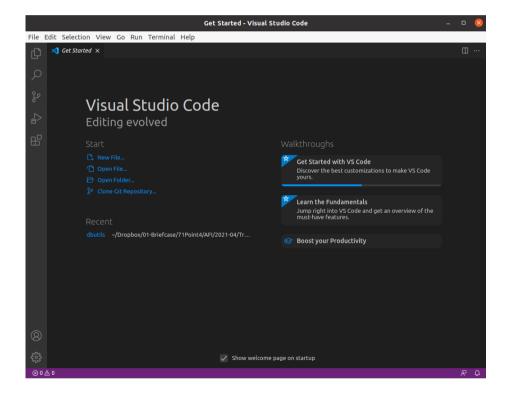
Learning to code in VSCode



Why switch from RStudio to VSCode for SQL development?

The first few things we are gonna do in VSCode is:

- Interact with a remote server
- Connect to database on remote server
- Execute code and download results

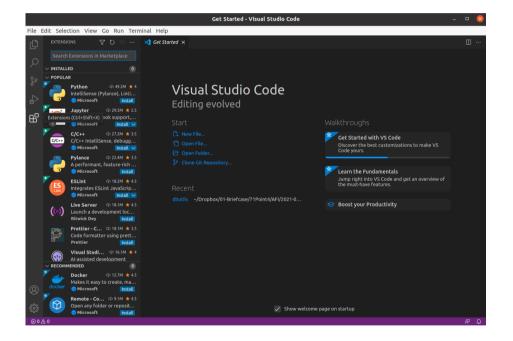


Installing the recommended Extension



Installing *Extensions* in VSCode is pretty straight forward. Just navigation to the search tab using GUI. Then search and install the following:

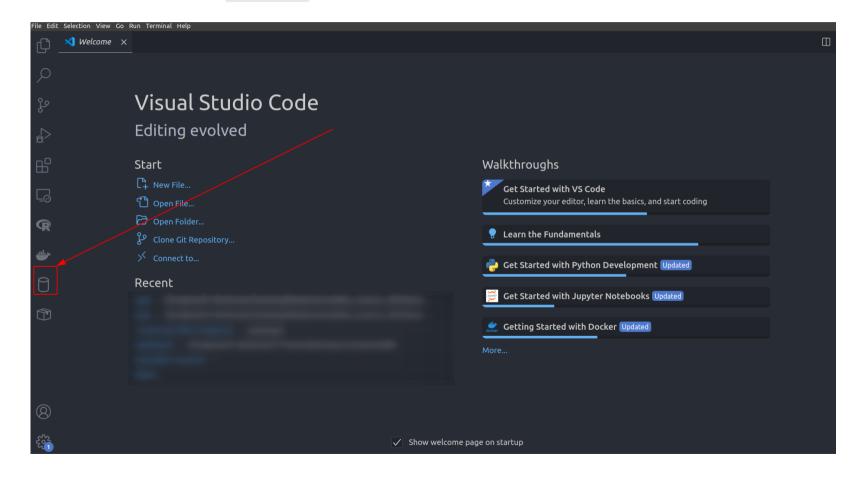
- R Extension for Visual Studio Code
- Spelling Checker for Visual Studio Code
- SQLTools



Connecting to DB



After installing VScode, you should see a SQLTools icon in the left-hand bar:

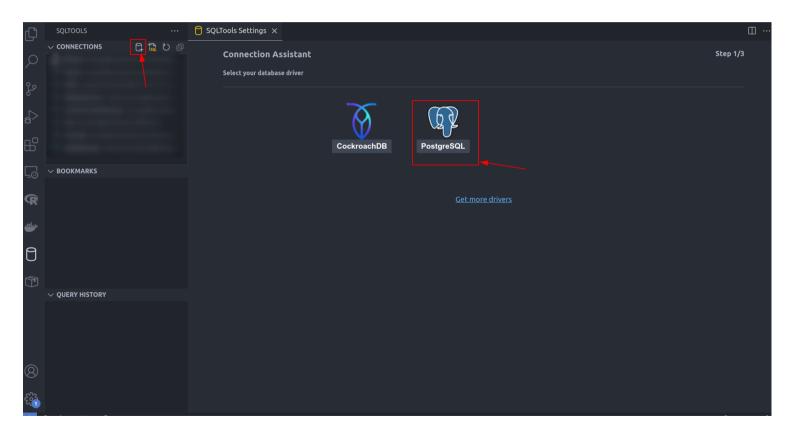


Connecting to DB



We can now add a new connection to a DB.

1 If you do not see PostgreSQL, then just click on 'Get More Drivers' and search for 'SQLTools PostgreSQL/Cockroach Driver'.



Connecting to DB



In the last step we going to fill in the connection strings (On the next slide):

Connection Assistant			< Step 2/3
Connection name*			60
Connection group			Val s
Connect using*	Server and Port		
Server Address*	localhost		
Port*	5432		
Database*			
Username*			
Use password	SQLTools Driver Credentials		
node-pg driver specif	ic options		
SSL	Disabled		
statement_timeout			
query_timeout			
connectionTimeoutMilli			
idleTimeoutMillis			
	Number of milliseconds a client must s and discarded. Default is 10000 (10 sec	it idle in the pool and not be checked out before it is disconnected from the ba conds) - set to 0 to disable auto-disconnection of idle clients	
max			
Connection Timeout			
Show records default limit			
SAVE CONNECTION		TE	ST CONNECTION