

# STAT 19000 Project 11 Examples

---

In this project we will explore another database engine called MySQL. You will find that the learning curve for basic SQL functionality is very low when switching between one of the more popular RDBMS's like MySQL, Postgresql, SQLite, Microsoft SQL Server, etc. This doesn't mean that the engines don't have some *very* different features, but it does mean that learning to navigate and pull some data will be a pretty straightforward process once you've done it on one system.

In addition, one of the primary ways to interact with any database is by using some package or interface within a given programming language. In this project we will learn to connect to both a MySQL and SQLite database from *within* R, and make queries.

**Important note:** This set of examples assumes you are using <https://rstudio.scholar.rcac.purdue.edu/>

## MySQL & RMariaDB

---

One major way that MySQL differs from SQLite is that MySQL has a server running that we must connect to in order to make queries, whereas an SQLite database is represented by a single file.

To connect to a MySQL database:

1. Open a terminal in scholar
2. Execute the following command: `mysql -u <username> -h scholar-db.rcac.purdue.edu -p`

**Note:** read-only usernames and passwords are available for class use in scholar at `/class/datamine/data/spring2020/mysql_credentials.txt`. For now, connect to MySQL using the `imdb` username.

3. Enter the password when prompted, and you will be presented with a prompt that resembles:

```
MariaDB [(none)]>
```

**Note:** MariaDB is a fork of MySQL. You can read about it [here](#).

**Get a list of the databases available to you:**

```
show databases;
```

**Make the imdb database your active database:**

```
use imdb;
```

**Show the tables in the active database:**

```
show tables;
```

**Show important table information:**

```
show columns from <table_name>; or describe <table_name>;
```

**Show any indexes a table may have:**

```
show index from <table_name>;
```

## Create a connection to the MySQL database using RMariaDB:

```
1 # if you haven't already, install RMariaDB
2 install.packages("RMariaDB")
3
4 library(RMariaDB)
5 host <- "scholar-db.rcac.purdue.edu"
6
7 # fill in the following values from the credentials provided on scholar at
8 # /class/datamine/data/spring2020/mysql_credentials.txt
9 # Note that you should not include the following two lines anywhere in your
10 # solutions.
11 user <- "<username>"
12 password <- "<password>"
13 connection <- dbConnect(RMariaDB::MariaDB(), host = host, user = user,
14 password = password)
```

## Connect to the proper database:

```
1 # executes a statement. use when you don't want a return value
2 dbExecute(connection, "USE imdb;")
3
4 # use dbGetQuery if you want the results (and they fit in memory) as a
5 # data.frame
6 dbGetQuery(connection, "SHOW tables;")
```

## Connect directly to the proper database without using the USE command:

```
1 library(RMariaDB)
2 host <- "scholar-db.rcac.purdue.edu"
3 dbname <- "imdb"
4
5 # fill in the following values from the credentials provided on scholar at
6 # /class/datamine/data/spring2020/mysql_credentials.txt
7 # Note that you should not include the following two lines anywhere in your
8 # solutions.
9 user <- "<username>"
10 password <- "<password>"
11 connection <- dbConnect(RMariaDB::MariaDB(), host = host, dbname=dbname,
12 user = user, password = password)
13
14 # this will now work right away
15 dbGetQuery(connection, "SHOW tables;")
```

## Get results directly into a data.frame, all at once:

```
1 # this sends, fetches, and clears for you
2 result <- dbGetQuery(connection, "SELECT * FROM episodes WHERE
3 show_title_id='tt0108778'")
```

## Get a result from a query one piece at a time:

```
1 # this sends, fetches, and clears for you
2 result <- dbGetQuery(connection, "SELECT * FROM episodes WHERE
3 show_title_id='tt0108778'")
```

```

1 res <- dbSendQuery(connection, "SELECT * FROM titles;")
2 some_data <- dbFetch(res, n=10)
3 some_data
4
5 # get the next 20 results (after the first 10 you fetched)
6 more_data <- dbFetch(res, n=20)
7 more_data
8
9 # until you try and dbFetch past the end of the data, this will be FALSE
10 dbHasCompleted(res)
11
12 # For example:
13 dbGetQuery(connection, "SELECT COUNT(*) FROM episodes;")
14
15 # This means we have 4658378 results to fetch before
16 dbHasCompleted(res)
17 # will return TRUE. Let's test it out:
18 # First clear the result
19 dbClearResult(res)
20 res <- dbSendQuery(connection, "SELECT * FROM episodes;")
21 dbHasCompleted(res) # FALSE
22 dat <- dbFetch(res, n=4658377)
23 dbHasCompleted(res) # FALSE
24 last <- dbFetch(res, n=1)
25 dbHasCompleted(res) # FALSE, what? You have to try and fetch *past* the end
26 one_more <- dbFetch(res, n=1)
27 one_more # empty
28 dbHasCompleted(res) # TRUE
29
30 # Last but not least, you must clear your results when using dbSendQuery,
    and dbFetch
31 dbClearResult(res)
32
33 # And if you are done with the database, you should disconnect
34 dbDisconnect(con)

```

### How to make comments in MySQL:

```

1 # this is how you make a comment, the -- notation does not work
2 /* This is how you make a block comment, this still works in MySQL,
3    as you can see.
4 */

```

### How to quit MySQL:

```

1 /* In sqlite we can quit/exit using .exit.
2 In MySQL, in order to quit or exit you simply type exit and then press enter
3 */

```

### How to use the LIMIT clause:

```
1  /* In MySQL, you must not use parenthesis for the limit clause like you can
   in sqlite.
2  */
3  # this is not valid:
4  select * from table limit(5);
5
6  # this IS valid
7  select * from table limit 5;
```

## SQLite & RSQLite

Luckily, there is virtually nothing new to learn in order to use SQLite within R.

**Create a connection to an sqlite database:**

```
1  # if you haven't already, install RMariaDB
2  install.packages("RSQLite")
3
4  library(RSQLite)
5
6  connection <- dbConnect(RSQLite::SQLite(),
   "/class/datamine/data/spring2020/imdb.db")
```

**Get a list of the tables:**

```
1  # this won't work
2  dbGetQuery(connection, ".tables")
3
4  # instead use this
5  dbListTables(connection)
```

And that is the extent of the differences (at least that we will cover).

## Examples

So when are we supposed to use ON vs WHERE?

This is a great question. You want to use ON to *join* data, and WHERE to *filter* data. Here is a great (and easy to read) article that should help: <https://dataschool.com/how-to-teach-people-sql/difference-between-where-and-on-in-sql/>

So let's say we want to join the episodes and title tables from our imdb.db database found on scholar:

```
/class/datamine/data/spring2020/imdb.db
```

In addition, we want to just do this for FRIENDS. In this case we want to join the tables where the episode\_title\_id from the episodes table and the title\_id from the titles table are the same. Then, we'd want to filter out all rows where the show\_title\_id from the episodes table (or alternatively the title\_id from the titles table) is not 'tt0108778'.

So how about this?

```

1 SELECT *
2 FROM episodes AS e
3 WHERE e.show_title_id = 'tt0108778'
4 INNER JOIN titles AS t
5     ON e.episode_title_id = t.title_id;

```

This won't work, the WHERE clause must come *after* the join. How about this?

```

1 SELECT *
2 FROM episodes AS e
3 INNER JOIN titles AS t
4     ON e.episode_title_id = t.title_id
5     WHERE e.show_title_id = 'tt0108778';

```

Yes, this will work. Now some may say, why don't we just move the show\_title\_id='tt0108778' part to the ON clause?

```

1 SELECT *
2 FROM episodes AS e
3 INNER JOIN titles AS t
4     ON e.episode_title_id = t.title_id
5     AND e.show_title_id = 'tt0108778';

```

In fact, you *can* do this and you will get the same result. However, let's look at if you do this and a LEFT JOIN:

```

1 SELECT *
2 FROM episodes AS e
3 LEFT JOIN titles AS t
4     ON e.episode_title_id = t.title_id
5     AND e.show_title_id = 'tt0108778';

```

Here, the result would be every row in the episodes table on the left side. On the right side, we'd have appended title information for only FRIENDS episodes. If instead you had done:

```

1 SELECT *
2 FROM episodes AS e
3 LEFT JOIN titles AS t
4     ON e.episode_title_id = t.title_id
5     WHERE e.show_title_id = 'tt0108778';

```

The interim table (the table that is created *before* the WHERE filtering), would contain every row in the episodes table and every matching piece of title info appended. Then it would be filtered using the WHERE clause and only episodes of FRIENDS would remain.

How about using the BETWEEN clause with the GROUP BY, is this possible? Yes. Here we demonstrate using the chinook database which can be found on scholar:

```
/class/datamine/data/spring2020/chinook.db
```

```
1  --Get the SUM(Total) for invoices in 2009 and invoices in 2013
2  SELECT InvoiceDate, SUM(Total) FROM invoices GROUP BY InvoiceDate BETWEEN
   '2009-01-01' AND '2009-12-31', InvoiceDate BETWEEN '2013-01-01' AND '2013-12-
   31';
```

The result:

InvoiceDate	SUM(Total)
2010-01-08 00:00:00	1428.56
2013-01-02 00:00:00	450.58
2009-01-01 00:00:00	449.46

You can see that in 2009 the total was 449.46 and in 2013 the total was 450 and for the rest of the dates the total was 1428. You can confirm these results with the following queries:

```
1  SELECT SUM(Total) FROM invoices WHERE InvoiceDate BETWEEN '2013-01-01' AND
   '2013-12-31';
2  SELECT SUM(Total) FROM invoices WHERE InvoiceDate BETWEEN '2009-01-01' AND
   '2009-12-31';
3  SELECT SUM(Total) FROM invoices;
```

If at this point you find you would like some more MySQL examples, you can find a good tutorial [here](#).