

STAT 19000 Project 9 Examples

In this series of examples we will learn some SQL basics by exploring the read-only sqlite3 database, chinook.db, [online](#), or on scholar:

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

Note: This database is < 1mb, so don't hesitate to download it onto your own system.

Relational databases and SQL

Relational databases are digital databases that store data into tables of columns and rows where a unique key identifies each row.

Typically, users use SQL (Structured Query Language) for managing (adding, deleting, updating, etc.) data inside the database.

SQLite

SQLite is a database engine that allows us to connect to the database and use SQL to manage the data. There are a variety of database engines, SQLite is one brand that happens to have a low learning curve when it comes to database-engine-specific tooling. It's simple to set up, connect to a database, and use.

From <https://sqlite.org>:

SQLite is a C-language library that implements a [small](#), [fast](#), [self-contained](#), [high-reliability](#), [full-featured](#), SQL database engine. SQLite is the [most used](#) database engine in the world. SQLite is built into all mobile phones and most computers and comes bundled inside countless other applications that people use every day.

The "self-contained" feature is particularly useful! From <https://sqlite.org>:

SQLite is "stand-alone" or "self-contained" in the sense that it has very few dependencies. It runs on any operating system, even stripped-down bare-bones embedded operating systems. SQLite uses no external libraries or interfaces (other than a few standard C-library calls described below). The entire SQLite library is encapsulated in a [single source code file](#) that requires no special facilities or tools to build.

Normally, a database engine has its own set of special dependencies that need to be properly installed and managed, making the time to get "up-and-running" higher.

In addition, the database itself (for sqlite), is typically stored in a single `.db` file, making it very convenient to transport, locate, and open.

SQLite does have its disadvantages, mainly its inability to have multiple simultaneous writers to the database. With that being said, we will only be reading from the database, for now, so it suits our needs well. Plus, other than database-engine unique syntax, [and a few omitted SQL features](#), the syntax is the same as any other relational database.

Enough is enough, let's get started.

The first step is to log in to scholar, and either connect to the database directly:

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

Or, copy and paste the database to a local directory:

```
mkdir -p ~/project09 && cp /class/datamine/data/spring2020/chinook.db "$_"
```

To connect directly, open a terminal and type the following:

```
sqlite3 /class/datamine/data/spring2020/chinook.db
```

To connect locally, open a terminal and type the following:

```
cd ~/project09; sqlite3 chinook.db
```

You should be presented with a simple prompt and message with the SQLite version. Great!

w3schools.com has a [good introduction](#) to SQL, including some exercises. For now, we are going to show some examples of the following SQL statements: `SELECT`, `FROM`, `WHERE`, `AS`, `SELECT DISTINCT`, `AND`, `OR`, `NOT`, `ORDER BY`, `LIMIT`, `MIN`, `MAX`, `COUNT`, `AVG`, `SUM`, `OFFSET`, `DESC`, `ASC` and `BETWEEN`. Additionally, we are going to show you how to make comments like we do in R using `#`.

```
1  -- Let's start by talking about comments. This line is a single-line
   comment,
2  -- and so is this one, because they start with two dashes "--".
3
4  /* A more appropriate solution when writing multiple lines is to use the
5  multi-line comments. Multi-line comments start with "/*" and end with
6  the two characters reversed, so "*" and then "/". Anything in between these
   sets of symbols is recognized by sqlite as a comment, great! */
7
8  -- All of the commands given here are to be run in the sqlite prompt.
9  -- Once you've started the sqlite prompt you should be able to see
10 -- something like:
11 --  sqlite>
12
13 /* We've loaded the chinook database. Databases can be pretty complicated,
   but, in general, when dealing with a relational database, you really need
   to think about tables and indexes.
14
15 Think of tables like an excel spreadsheet. The spreadsheet has a header row
   with the names of the columns, and each remaining row contains one datum
   per column that conforms to a certain type (i.e. in a spreadsheet, one
   column may have floating point integers, another may have integers, another
   may have strings, etc.). In the same way, a table has columns (or fields),
   with headers. Tables have rows of data that conform to a certain, pre-
   defined data type.
```

```
16
```

```

17 Great. Now what about indexes. Think of an index just like an index you
    would find in a text book. Those last 10-20 pages of a 800 page book that
    allows you to find some term quickly (because it's in alphabetical order),
    and immediately flip to the page in the book where the topic is discussed.
    If you didn't have the index you'd have to flip through every page to find
    what you are looking for. The same concept applies to databases. Sometimes,
    if you have a table with large amounts of data, creating an index on a
    certain column, or combination of columns will greatly increase query
    performance. Great, so why don't we create indexes for everything? Indexes
    take up storage space, if the table is small an index won't necessarily
    speed things up, if the table gets updated frequently, the indexes will
    need updated too, which takes time.
18
19 Great, now let's actually do something. */
20
21 -- What are the tables within the database
22 .tables
23 /*albums          employees          invoices          playlists
24 artists           genres             media_types       tracks
25 customers          invoice_items      playlist_track*/
26
27 -- Are there any indexes?
28 .indices
29
30 /*IFK_AlbumArtistId          IFK_PlaylistTrackTrackId
31 IFK_CustomerSupportRepId     IFK_TrackAlbumId
32 IFK_EmployeeReportsTo       IFK_TrackGenreId
33 IFK_InvoiceCustomerId       IFK_TrackMediaTypeId
34 IFK_InvoiceLineInvoiceId     sqlite_autoindex_playlist_track_1
35 IFK_InvoiceLineTrackId*/
36
37 -- Let's turn on headers so it is easy to tell the data from our queries
38 .headers on
39
40 /* Let's test it out by selecting all of the data, every column (*), from
    the first employee only. */
41 SELECT * FROM employees LIMIT(1);
42 /*EmployeeId|LastName|FirstName|Title|ReportsTo|BirthDate|HireDate|Address
    |City|State|Country|PostalCode|Phone|Fax|Email
43 1|Adams|Andrew|General Manager||1962-02-18 00:00:00|2002-08-14
    00:00:00|11120 Jasper Ave NW|Edmonton|AB|Canada|T5K 2N1|+1 (780) 428-
    9482|+1 (780) 428-3457|andrew@chinookcorp.com */
44
45 -- Great, we can easily see the headers: EmployeeId, LastName, Title, etc.
46
47 -- Let's break this query down into parts:
48
49 -- The SELECT command tells sqlite that we want to get (or "select") data
    from the database.
50 -- The * immediately following the SELECT tells sqlite to select all of the
    columns.
51
52 -- Instead of * you could do:
53 SELECT FirstName, LastName FROM employees LIMIT(1);
54 /*FirstName|LastName
55 Andrew|Adams*/
56

```

```

57  -- Note, that this will also work as column names and sql commands are not
    case sensitive.
58  select firstname, lastname from employees limit(1);
59
60  -- With that being said, capitalizing the commands makes things easier to
    read.
61  -- The FROM command tells sqlite from what table do we want to get every
    column of
62  -- information?
63  -- The "employees" part following FROM is the table.
64  -- LIMIT(1) tells sqlite to limit the output to a single row.
65  -- LIMIT(2) would limit the output to two rows, etc.
66  -- Lastly, the semi-colon tells sqlite that this is the end of the
    statement.
67
68  -- Let's try that LIMIT(2) thing out...
69  SELECT FirstName, LastName FROM employees LIMIT(2);
70  /*FirstName|LastName
71  Andrew|Adams
72  Nancy|Edwards*/
73
74  -- Great, and if we look they are indeed the first two:
75  SELECT FirstName, LastName FROM employees;
76  /*FirstName|LastName
77  Andrew|Adams
78  Nancy|Edwards
79  Jane|Peacock
80  Margaret|Park
81  Steve|Johnson
82  Michael|Mitchell
83  Robert|King
84  Laura|Callahan*/
85
86  -- What if we wanted the LastName and then the FirstName?
87  SELECT LastName, FirstName FROM employees;
88  /*LastName|FirstName
89  Adams|Andrew
90  Edwards|Nancy
91  Peacock|Jane
92  Park|Margaret
93  Johnson|Steve
94  Mitchell|Michael
95  King|Robert
96  Callahan|Laura*/
97
98  -- Cool, so you just put them in whatever order you want them.
99  -- What if we, for some reason wanted something repeated like James, James
    Bond
100 SELECT FirstName, FirstName, LastName FROM employees;
101 /*FirstName|FirstName|LastName
102 Andrew|Andrew|Adams
103 Nancy|Nancy|Edwards
104 Jane|Jane|Peacock
105 Margaret|Margaret|Park
106 Steve|Steve|Johnson
107 Michael|Michael|Mitchell
108 Robert|Robert|King
109 Laura|Laura|Callahan*/

```

```

110
111 -- Nice.
112
113 -- What if instead of the first two names, we wanted to get the first two
names
114 -- that occur after the 5th name?
115 SELECT FirstName, LastName FROM employees LIMIT(2) OFFSET(5);
116 /*FirstName|LastName
117 Michael|Mitchell
118 Robert|King*/
119
120 -- What if we wanted the last two names? Well, for this database we could
do:
121 SELECT FirstName, LastName FROM employees LIMIT(2) OFFSET(6);
122
123 -- BUT, as soon as the database adds a new row, this won't work.
124 -- So, first of all, how is the order currently determined? It's a
125 -- small table, so let's look:
126 SELECT * FROM employees;
127 /*EmployeeId|LastName|FirstName|Title|ReportsTo|BirthDate|HireDate|Address
|City|State|Country|PostalCode|Phone|Fax|Email
128 1|Adams|Andrew|General Manager||1962-02-18 00:00:00|2002-08-14
00:00:00|11120 Jasper Ave NW|Edmonton|AB|Canada|T5K 2N1|+1 (780) 428-
9482|+1 (780) 428-3457|andrew@chinookcorp.com
129 2|Edwards|Nancy|Sales Manager|1|1958-12-08 00:00:00|2002-05-01 00:00:00|825
8 Ave SW|Calgary|AB|Canada|T2P 2T3|+1 (403) 262-3443|+1 (403) 262-
3322|nancy@chinookcorp.com
130 3|Peacock|Jane|Sales Support Agent|2|1973-08-29 00:00:00|2002-04-01
00:00:00|1111 6 Ave SW|Calgary|AB|Canada|T2P 5M5|+1 (403) 262-3443|+1 (403)
262-6712|jane@chinookcorp.com
131 4|Park|Margaret|Sales Support Agent|2|1947-09-19 00:00:00|2003-05-03
00:00:00|683 10 Street SW|Calgary|AB|Canada|T2P 5G3|+1 (403) 263-4423|+1
(403) 263-4289|margaret@chinookcorp.com
132 5|Johnson|Steve|Sales Support Agent|2|1965-03-03 00:00:00|2003-10-17
00:00:00|7727B 41 Ave|Calgary|AB|Canada|T3B 1Y7|1 (780) 836-9987|1 (780)
836-9543|steve@chinookcorp.com
133 6|Mitchell|Michael|IT Manager|1|1973-07-01 00:00:00|2003-10-17
00:00:00|5827 Bowness Road NW|Calgary|AB|Canada|T3B 0C5|+1 (403) 246-
9887|+1 (403) 246-9899|michael@chinookcorp.com
134 7|King|Robert|IT Staff|6|1970-05-29 00:00:00|2004-01-02 00:00:00|590
Columbia Boulevard West|Lethbridge|AB|Canada|T1K 5N8|+1 (403) 456-9986|+1
(403) 456-8485|robert@chinookcorp.com
135 8|Callahan|Laura|IT Staff|6|1968-01-09 00:00:00|2004-03-04 00:00:00|923 7
ST NW|Lethbridge|AB|Canada|T1H 1Y8|+1 (403) 467-3351|+1 (403) 467-
8772|laura@chinookcorp.com*/
136
137 -- Clearly, the table is ordered by EmployeeId. We will dig into that a
little bit later.
138 -- So, one way to do this would be to use DESC and ORDER BY
139 SELECT EmployeeId, FirstName, LastName FROM employees ORDER BY EmployeeId
DESC;
140 /*FirstName|LastName
141 Laura|Callahan
142 Robert|King
143 Michael|Mitchell
144 Steve|Johnson
145 Margaret|Park
146 Jane|Peacock

```

```

147 Nancy|Edwards
148 Andrew|Adams*/
149
150 -- Great, just add the LIMIT back now
151 SELECT FirstName, LastName FROM employees ORDER BY EmployeeId DESC
LIMIT(2);
152 /*FirstName|LastName
153 Laura|Callahan
154 Robert|King*/
155
156 -- By default, ORDER BY goes in ascending order, but you can specify using
ASC.
157 -- You can also sort by other columns:
158 SELECT FirstName FROM employees ORDER BY FirstName ASC;
159 /*FirstName
160 Andrew
161 Jane
162 Laura
163 Margaret
164 Michael
165 Nancy
166 Robert
167 Steve*/
168
169 -- Let's look at the Title column
170 SELECT Title FROM employees;
171 /*Title
172 General Manager
173 Sales Manager
174 Sales Support Agent
175 Sales Support Agent
176 Sales Support Agent
177 IT Manager
178 IT Staff
179 IT Staff*/
180
181 /* This result is a little bit redundant. Imagine if we just wanted a list
of titles but there were 10000 employees! This is where SELECT DISTINCT can
be useful. SELECT DISTINCT returns only unique results or unique
combinations depending on what column(s) you SELECT DISTINCT.*/
182 SELECT DISTINCT Title FROM employees;
183 /*Title
184 General Manager
185 Sales Manager
186 Sales Support Agent
187 IT Manager
188 IT Staff*/
189
190 -- But if you wanted FirstName/Title combos (of which all are unique),
you'd get:
191 SELECT DISTINCT FirstName, Title FROM employees;
192 /*FirstName|Title
193 Andrew|General Manager
194 Nancy|Sales Manager
195 Jane|Sales Support Agent
196 Margaret|Sales Support Agent
197 Steve|Sales Support Agent
198 Michael|IT Manager

```

```

199 Robert|IT Staff
200 Laura|IT Staff*/
201
202 -- SQL also has some built in functions including but not limited to: MIN,
    MAX, COUNT, SUM, AVG
203 SELECT MIN(Total) FROM invoices;
204 /*MIN(Total)
205 0.99*/
206
207 -- As you can see the invoice table has what appears to be a variety of
    different types.
208 -- Let's check them out.
209 .schema invoices
210 /*
211 CREATE TABLE IF NOT EXISTS "invoices"
212 (
213     [InvoiceId] INTEGER PRIMARY KEY AUTOINCREMENT NOT NULL,
214     [CustomerId] INTEGER NOT NULL,
215     [InvoiceDate] DATETIME NOT NULL,
216     [BillingAddress] NVARCHAR(70),
217     [BillingCity] NVARCHAR(40),
218     [BillingState] NVARCHAR(40),
219     [BillingCountry] NVARCHAR(40),
220     [BillingPostalCode] NVARCHAR(10),
221     [Total] NUMERIC(10,2) NOT NULL,
222     FOREIGN KEY ([CustomerId]) REFERENCES "customers" ([CustomerId])
223         ON DELETE NO ACTION ON UPDATE NO ACTION
224 );
225 CREATE INDEX [IFK_InvoiceCustomerId] ON "invoices" ([CustomerId]);
226 */
227
228 /* Interesting! This is the SQL code that is used to build the invoices
    table. This is a lot to take in, but for now, the main takeaway should be
    that each column has a type explicitly declared.
229
230 For example, InvoiceId is an INTEGER, InvoiceDate is a DATETIME, many of
    these are NVARCHAR(X) where X is the length of the NVARCHAR, and Total is
    NUMERIC. If interested, you can read about SQLite types here:
    https://www.sqlite.org/datatype3.htmls
231
232 For now, ignore the rest, but note that using .schema can give you a pretty
    good idea what the data will look like.
233 */
234
235 SELECT MAX(Total) FROM invoices;
236 /*MAX(Total)
237 25.86*/
238
239 -- You can still use LIMIT...
240 SELECT MAX(Total) FROM invoices LIMIT(2);
241 /*MAX(Total)
242 25.86*/
243
244 -- BUT it still finds the maximum of the entire column, "Total"
245 -- AVG and SUM work as expected
246
247 SELECT AVG(Total) FROM invoices;
248 /*AVG(Total)

```

```

249 5.65194174757282*/
250
251 -- But again, AVGS the whole column
252 SELECT AVG(Total) FROM invoices LIMIT(2);
253 /*AVG(Total)
254 5.65194174757282*/
255
256 -- SUM, same behavior
257 SELECT SUM(Total) FROM invoices;
258 /*SUM(Total)
259 2328.6*/
260
261 SELECT SUM(Total) FROM invoices LIMIT(2);
262 /*SUM(Total)
263 2328.6*/
264
265 -- COUNT is useful too, and returns a count of the number of rows when
266 -- COUNT(*), otherwise if it is COUNT(column_header), it returns the count
267 -- of non null values.
268
269 SELECT COUNT(*) FROM invoices;
270 /*COUNT(*)
271 412*/
272
273 -- What if we want to add a condition to our queries. For example we want
274 -- only invoices
275 -- for Brazil? Well, the WHERE clause let's us do this!
276 SELECT * FROM invoices WHERE BillingCountry="Brazil";
277
278 -- And it can be combined with function too, finding the total for Brazil
279 -- invoices
280 -- is simple then:
281 SELECT SUM(Total) FROM invoices WHERE BillingCountry="Brazil";
282 /*SUM(Total)
283 190.1*/
284
285 -- And of course, you can find the average and sum at the same time:
286 SELECT AVG(Total), SUM(Total) FROM invoices WHERE BillingCountry="Brazil";
287 /*AVG(Total)|SUM(Total)
288 5.43142857142857|190.1*/
289
290 -- BETWEEN can be useful as well, and selects the rows that are between two
291 -- values,
292 -- inclusively (which means if we want numbers between 1-10 we would
293 -- include 1 and 10 in the list).
294
295 -- Let's get invoices that total between 13-14.
296 SELECT * FROM invoices WHERE Total BETWEEN 13 AND 14;
297
298 -- BETWEEN works on numbers, dates, and strings.
299
300 -- Date example: get invoices between 1/1/2013 and 2/1/2013
301 SELECT * FROM invoices WHERE InvoiceDate BETWEEN '2013-1-1' AND '2013-2-1';
302
303 -- Uh oh! This is wrong. Why? Because you need to zero pad your month and
304 -- day values:

```



```

301 SELECT * FROM invoices WHERE InvoiceDate BETWEEN '2013-01-01' AND '2013-
    02-01';
302
303 -- Great! How many of those were in the UK?
304 SELECT * FROM invoices WHERE InvoiceDate BETWEEN '2013-01-01' AND '2013-
    02-01' AND BillingCountry="United Kingdom";
305 /*335|53|2013-01-15 00:00:00|113 Lupus St|London||United Kingdom|SW1V
    3EN|0.99
306 336|54|2013-01-28 00:00:00|110 Raeburn Pl|Edinburgh ||United Kingdom|EH4
    1HH|1.98*/
307
308 -- Excellent! Although not explicitly mentioned earlier, AND, OR, and NOT
    function just like
309 -- normal logical operators in any language. You can use parenthesis to
    force order of
310 -- operations as well.
311
312 SELECT * FROM invoices WHERE Total > 20 OR Total < 5 AND
    BillingCountry="Brazil";
313
314 -- Uh oh! We wanted invoices where the Total was either greater than 20 or
    less
315 -- than 5 and from Brazil. Due to the fact that AND has precedence over OR
    (in the
316 -- same way multiplication has precedence over addition and subtraction),
317 -- what this statement did was equivalent to:
318 SELECT * FROM invoices WHERE Total > 20 OR (Total < 5 AND
    BillingCountry="Brazil");
319
320 -- To fix, add parenthesis:
321 SELECT * FROM invoices WHERE (Total > 20 OR Total < 5) AND
    BillingCountry="Brazil";
322
323 -- Now the last topic we will touch on is called aliasing (using AS). I'd
    recommend reading:
324 -- https://www.w3schools.com/sql/sql\_alias.asp
325
326 -- You can alias a column...
327 SELECT BillingCountry AS Country FROM invoices LIMIT(1);
328 /*Country
329 Germany*/
330
331 -- or multiple columns
332 SELECT BillingCountry AS Country, Total AS DollarValue FROM invoices
    LIMIT(1);
333 /*Country|DollarValue
334 Germany|1.98*/
335
336 -- You can alias a table, and access that tables individual values
337 SELECT inv.BillingCountry FROM invoices AS inv LIMIT(1);
338 /*BillingCountry
339 Germany*/
340
341 -- Which allows you to query multiple tables at once
342 SELECT i.BillingCountry, c.FirstName, c.LastName FROM invoices AS i,
    customers AS c WHERE i.CustomerId = c.CustomerId;
343

```

```

344 -- That is pretty cool. We printed the invoices billing country, and
      customer first and last
345 -- name every place that the customer id matches in the tables.
346
347 -- You can also "combine" columns into a single string. Note that the
      syntax varies
348 -- depending on the database engine. In sqlite, we use || to concatenate.
349 SELECT InvoiceId, BillingCity || ', ' || BillingCountry AS Place FROM
      invoices;
350
351 -- If the alias name contains space you need to use ""
352 SELECT InvoiceId AS "Invoice ID" FROM invoices LIMIT(1);
353 /*Invoice ID
354 98*/
355
356 /* A cool feature of sqlite is it is really easy to export a query result
357 into some file using the .mode and .once dot commands.*/
358
359 -- Check what mode we are currently in:
360 .mode
361
362 -- Export entire invoices table to a csv
363 .mode csv
364 .output invoices.csv
365 SELECT * FROM invoices;
366
367 -- To go back to using | separator
368 .separator |
369 .mode list
370
371 -- Very cool! By default the .mode is "list". Some other useful modes
      include "tabs"
372 -- and "html". This is a very reasonable way to go from sqlite database ->
      analyzing a
373 -- table in R.
374
375 -- Feel free to explore SQL and sqlite3 features in more depth!

```

You can see a full list of dot commands by typing `.help`:

```

1  .backup ?DB? FILE      Backup DB (default "main") to FILE
2  .bail ON|OFF           Stop after hitting an error.  Default OFF
3  .databases             List names and files of attached databases
4  .dump ?TABLE? ...      Dump the database in an SQL text format
5                          If TABLE specified, only dump tables matching
6                          LIKE pattern TABLE.
7  .echo ON|OFF           Turn command echo on or off
8  .exit                  Exit this program
9  .explain ?ON|OFF?      Turn output mode suitable for EXPLAIN on or off.
10                          With no args, it turns EXPLAIN on.
11 .header(s) ON|OFF      Turn display of headers on or off
12 .help                  Show this message
13 .import FILE TABLE    Import data from FILE into TABLE
14 .indices ?TABLE?       Show names of all indices
15                          If TABLE specified, only show indices for tables

```

16		matching LIKE pattern TABLE.
17	.load FILE ?ENTRY?	Load an extension library
18	.log FILE off	Turn logging on or off. FILE can be stderr/stdout
19	.mode MODE ?TABLE?	Set output mode where MODE is one of:
20		csv Comma-separated values
21		column Left-aligned columns. (See .width)
22		html HTML <table> code
23		insert SQL insert statements for TABLE
24		line One value per line
25		list Values delimited by .separator string
26		tabs Tab-separated values
27		tcl TCL list elements
28	.nullvalue STRING	Use STRING in place of NULL values
29	.output FILENAME	Send output to FILENAME
30	.output stdout	Send output to the screen
31	.print STRING...	Print literal STRING
32	.prompt MAIN CONTINUE	Replace the standard prompts
33	.quit	Exit this program
34	.read FILENAME	Execute SQL in FILENAME
35	.restore ?DB? FILE	Restore content of DB (default "main") from FILE
36	.schema ?TABLE?	Show the CREATE statements
37		If TABLE specified, only show tables matching
38		LIKE pattern TABLE.
39	.separator STRING	Change separator used by output mode and .import
40	.show	Show the current values for various settings
41	.stats ON OFF	Turn stats on or off
42	.tables ?TABLE?	List names of tables
43		If TABLE specified, only list tables matching
44		LIKE pattern TABLE.
45	.timeout MS	Try opening locked tables for MS milliseconds
46	.trace FILE off	Output each SQL statement as it is run
47	.vfsname ?AUX?	Print the name of the VFS stack
48	.width NUM1 NUM2 ...	Set column widths for "column" mode
49	.timer ON OFF	Turn the CPU timer measurement on or off