

Loading Taxi Data into Google Cloud SQL

 partner.cloudskillsboost.google/course_sessions/221410/labs/61489

Overview

In this lab, you will learn how to import data from CSV text files into Cloud SQL and then carry out some basic data analysis using simple queries.

The dataset used in this lab is collected by the NYC Taxi and Limousine Commission and includes trip records from all trips completed in Yellow and Green taxis in NYC from 2009 to present, and all trips in for-hire vehicles (FHV) from 2015 to present. Records include fields capturing pick-up and drop-off dates/times, pick-up and drop-off locations, trip distances, itemized fares, rate types, payment types, and driver-reported passenger counts.

This dataset can be used to demonstrate a wide range of data science concepts and techniques and will be used in several of the labs in the Data Engineering curriculum.

Objectives

- Create Cloud SQL instance
- Create a Cloud SQL database
- Import text data into Cloud SQL
- Check the data for integrity

Setup and Requirements

For each lab, you get a new Google Cloud project and set of resources for a fixed time at no cost.

1. Sign in to Qwiklabs using an **incognito window**.
2. Note the lab's access time (for example, 02:00:00), and make sure you can finish within that time. There is no pause feature. You can restart if needed, but you have to start at the beginning.
3. When ready, click **Start lab**.
4. Note your lab credentials (**Username** and **Password**). You will use them to sign in to the Google Cloud Console.
5. Click **Open Google Console**.

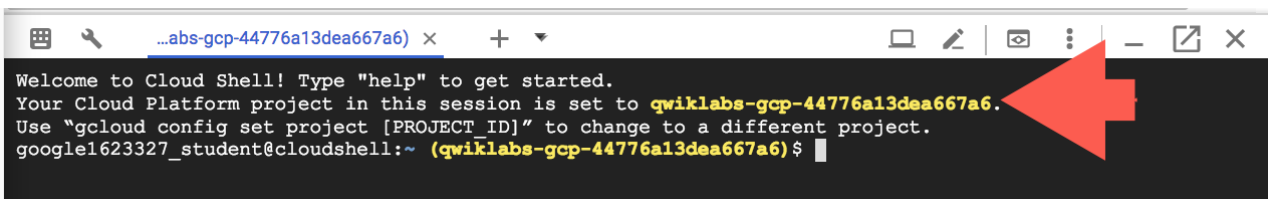
6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts. If you use other credentials, you'll receive errors or **incur charges**.
7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you have finished the lab or want to restart it. This clears your work and removes the project.

Activate Cloud Shell

Cloud Shell is a virtual machine that contains development tools. It offers a persistent 5-GB home directory and runs on Google Cloud. Cloud Shell provides command-line access to your Google Cloud resources. `gcloud` is the command-line tool for Google Cloud. It comes pre-installed on Cloud Shell and supports tab completion.

1. In the Google Cloud Console, in the navigation pane, click **Activate Cloud Shell** (image).
2. Click **Continue**. It takes a few moments to provision and connect to the environment. When you are connected, you are also authenticated, and the project is set to your *PROJECT_ID*. For example:



```
...abs-gcp-44776a13dea667a6) x + v
Welcome to Cloud Shell! Type "help" to get started.
Your Cloud Platform project in this session is set to qwiklabs-gcp-44776a13dea667a6.
Use "gcloud config set project [PROJECT_ID]" to change to a different project.
google1623327_student@cloudshell:~ (quwiklabs-gcp-44776a13dea667a6) $
```

Sample commands

List the active account name:

```
gcloud auth list
```

(Output)

Credentialed accounts: - <myaccount>@<mydomain>.com (active)
(Example output)

Credentialed accounts: - google1623327_student@qwiklabs.net
List the project ID:

```
gcloud config list project
```

(Output)

[core] project = <project_ID>
(Example output)

[core] project = qwiklabs-gcp-44776a13dea667a6

For full documentation of `gcloud`, see the [gcloud command-line tool overview](#).

Preparing your Environment

Create **environment** variables that will be used later in the lab for your project ID and the storage bucket that will contain your data:

```
export PROJECT_ID=$(gcloud info --format='value(config.project)') export  
BUCKET=${PROJECT_ID}-ml
```

Create a Cloud SQL instance

Enter the following **commands** to **create a Cloud SQL instance**:

```
gcloud sql instances create taxi \ --tier=db-n1-standard-1 --activation-policy=ALWAYS
```

This will take a **few minutes to complete**.

Test Completed Task

Click **Check my progress** to verify your performed task. If you have completed the task successfully you will be granted with an assessment score.

Create a Cloud SQL instance.

Set a root password for the Cloud SQL instance:

```
gcloud sql users set-password root --host % --instance taxi \ --password Password
```

When prompted for the **password** type **Password** and **press enter** this will **update root password**.

Now create an **environment variable** with the **IP address of the Cloud Shell**:

```
export ADDRESS=$(wget -qO - http://ipecho.net/plain)/32
```

Whitelist the Cloud Shell instance for management access to your SQL instance.

```
gcloud sql instances patch taxi --authorized-networks $ADDRESS
```

When prompted **press Y** to **accept the change**.

Test Completed Task

Click **Check my progress** to verify your performed task. If you have completed the task successfully you will be granted with an assessment score.

Whitelist the Cloud Shell instance to access your SQL instance.

Get the **IP address of your Cloud SQL instance** by **running**:

```
MYSQLIP=$(gcloud sql instances describe \ taxi --  
format="value(ipAddresses.ipAddress)")
```

Check the variable **MYSQLIP**:

```
echo $MYSQLIP
```

you should get an IP address as an output.

Create the `taxi trips` table by logging into the `mysql` command line interface.

```
mysql --host=$MYSQLIP --user=root \ --password --verbose
```

When prompted for a password enter `Password`. Paste the following content into the command line to create the schema for the `trips` table:

```
create database if not exists bts; use bts; drop table if exists trips; create table trips (
vendor_id VARCHAR(16), pickup_datetime DATETIME, dropoff_datetime DATETIME,
passenger_count INT, trip_distance FLOAT, rate_code VARCHAR(16),
store_and_fwd_flag VARCHAR(16), payment_type VARCHAR(16), fare_amount
FLOAT, extra FLOAT, mta_tax FLOAT, tip_amount FLOAT, tolls_amount FLOAT,
imp_surcharge FLOAT, total_amount FLOAT, pickup_location_id VARCHAR(16),
dropoff_location_id VARCHAR(16) );
```

Test Completed Task

Click **Check my progress** to verify your performed task. If you have completed the task successfully you will be granted with an assessment score.

Create a `bts` database and `trips` table.

In the `mysql` command line interface check the import by entering the following commands:

```
describe trips;
```

Query the `trips` table:

```
select distinct(pickup_location_id) from trips;
```

This will return an empty set as there is no data in the database yet.

Exit the `mysql` interactive console:

```
exit
```

Add data to Cloud SQL instance

Now you'll copy the New York City taxi trips CSV files stored on Cloud Storage locally. To keep resource usage low, you'll only be working with a subset of the data (~20,000 rows).

Run the following in the command line:

```
gsutil cp gs://cloud-training/OCBL013/nyc_tlc_yellow_trips_2018_subset_1.csv
trips.csv-1 gsutil cp gs://cloud-
training/OCBL013/nyc_tlc_yellow_trips_2018_subset_2.csv trips.csv-2
```

Import the CSV file data into Cloud SQL using `mysql` :

```
mysqlimport --local --host=$MYSQLIP --user=root --password \ --ignore-lines=1 --fields-
terminated-by=',' bts trips.csv-*
```

When prompted for a password enter `Password`.

Connect to the `mysql` interactive console:

```
mysql --host=$MYSQLIP --user=root --password
```

When prompted for a password enter `Password`.

Checking for data integrity

Whenever data is imported from a source it's always important to check for data integrity. Roughly, this means making sure the data meets your expectations.

In the `mysql` interactive console select the database:

```
use bts;
```

Query the `trips` table for unique pickup location regions:

```
select distinct(pickup_location_id) from trips;
```

This should return 159 unique ids. Let's start by digging into the `trip_distance` column. Enter the following query into the console:

```
select max(trip_distance), min(trip_distance) from trips;
```

One would expect the trip distance to be greater than 0 and less than, say 1000 miles. The maximum trip distance returned of 85 miles seems reasonable but the minimum trip distance of 0 seems buggy. How many trips in the dataset have a trip distance of 0?

```
select count(*) from trips where trip_distance = 0;
```

There are 155 such trips in the database. These trips warrant further exploration. You'll find that these trips have non-zero payment amounts associated with them. Perhaps these are fraudulent transactions? Let's see if we can find more data that doesn't meet our expectations. We expect the `fare_amount` column to be positive. Enter the following query to see if this is true in the database:

```
select count(*) from trips where fare_amount < 0;
```

There should be 14 such trips returned. Again, these trips warrant further exploration. There may be a reasonable explanation for why the fares take on negative numbers. However, it's up to the data engineer to ensure there are no bugs in the data pipeline that would cause such a result.

Finally, let's investigate the `payment_type` column.

```
select payment_type, count(*) from trips group by payment_type;
```

The results of the query indicate that there are four different payment types, with:

- payment type = 1 has 13863 rows
- payment type = 2 has 6016 rows
- payment type = 3 has 113 rows
- payment type = 4 has 32 rows

Digging into the documentation, a payment type of 1 refers to credit card use, payment type of 2 is cash, and a payment type of 4 refers to a dispute. The figures make sense.

Exit the 'mysql' interactive console:

exit

End your lab

When you have completed your lab, click **End Lab**. Qwiklabs removes the resources you've used and cleans the account for you.

You will be given an opportunity to rate the lab experience. Select the applicable number of stars, type a comment, and then click **Submit**.

The number of stars indicates the following:

- 1 star = Very dissatisfied
- 2 stars = Dissatisfied
- 3 stars = Neutral
- 4 stars = Satisfied
- 5 stars = Very satisfied

You can close the dialog box if you don't want to provide feedback.

For feedback, suggestions, or corrections, please use the **Support** tab.

©2021 Google LLC All rights reserved. Google and the Google logo are trademarks of Google LLC. All other company and product names may be trademarks of the respective companies with which they are associated.