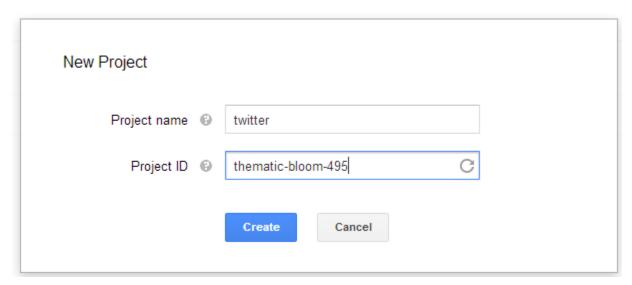# Google BigQuery

## I.    Load the data into a Google BigQuery endpoint for analysis

Download the sample files by copying and pasting these URLs into browser to download these two files and then upload to Google Cloud Storage
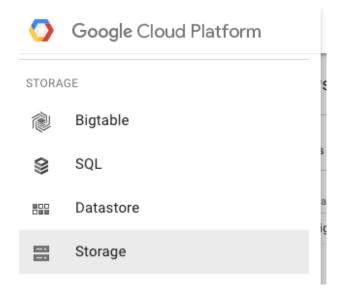
https://storage.googleapis.com/uwclddata230publichw2017/week07hw/seq0fix.csv

https://storage.googleapis.com/uwclddata230publichw2017/week07hw/seq1fix.csv

To load the data into Google BigQuery the first thing you need to do is upload the files into Google Cloud Storage
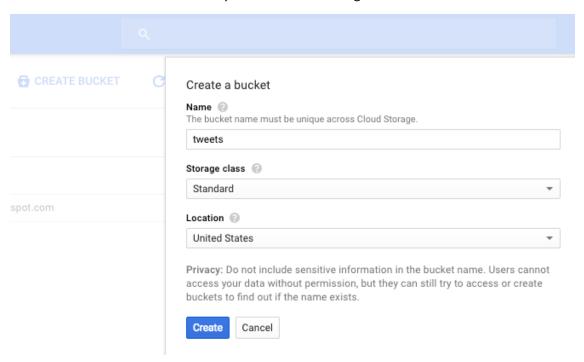
If you have not already done so you might have to create a New Project since the Cloud Services and billing are Project based.  Here is an example to creating a project called twitter
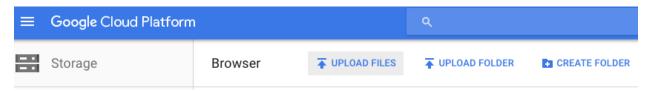


Create a Storage Bucket

Click on **CREATE BUCKET** and name your bucket something like tweets and click **Create**



As I mentioned earlier Google has a nice browse-based upload for Files or the contents of a Folder
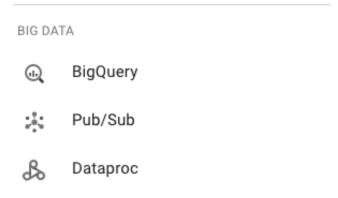


Upload the 2 files or the folder with the 2 files

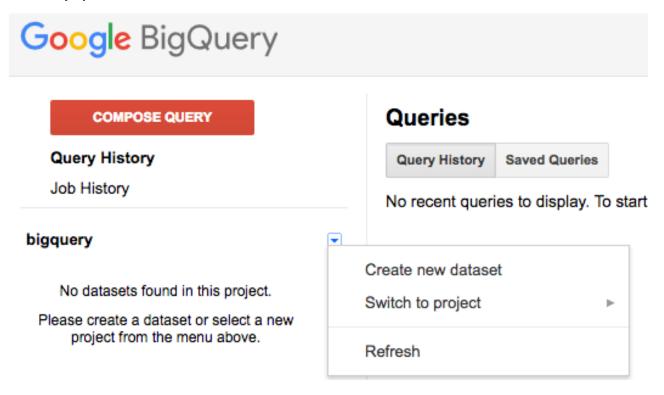When you select the folder the files start to upload and provide you with a status

Upload 0 of 2 complete                                        Cancel    —

seq0fix.csv                                        ━━━━━         52%        ×

seq1fix.csv                                        ▪━━━          2%        ×

You only have seq0fix and seq1fix

Buckets  /  tweetsuwclddata230  /  input  /  seq0fix

| | Name | Size |
|---|---|---|
| ☐ | 📄 seq0fix.csv | 9.18 MB |
| ☐ | 📄 seq1fix.csv | 197.75 MB |

Click on BigQuery in the Big Data section in the left panel of the Google Developer Console

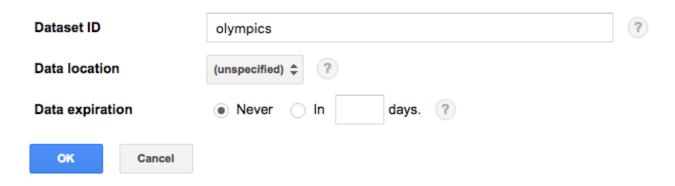BIG DATA

◉  BigQuery

⁘  Pub/Sub

⅋  Dataproc

Once the files are uploaded to Cloud Storage you can create a BigQuery Dataset (mine is called **olympics**) and upload the data into a BigQuery Table (mine is called **summer2012**)

Create **olympics** Dataset
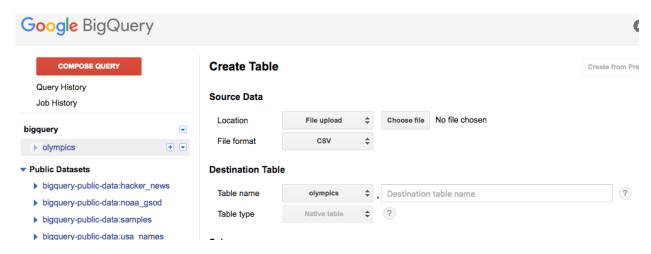


Enter Dataset Name and click the OK button



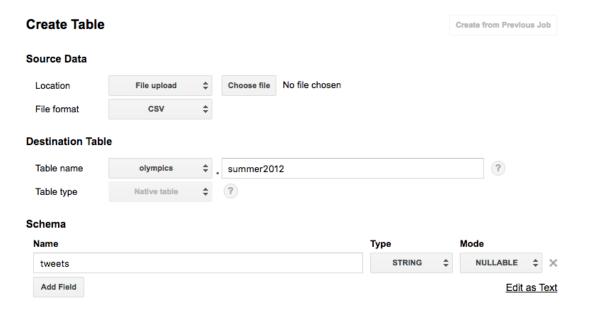Create the **summer2012** Table

Pick Google Cloud Storage as the location and provide a path the file.  Mine is gs://tweetsuwclddata230/input/seq0fix/seq0fix.csv

**gs://<bucket>/<folderIfYouHaveOne>/seq0fix.csv**)



Enter Table name as **summer2012** and Schema Name as **tweets** and then click Edit as Text

Next it asks for the Schema which is as follows:

id:INTEGER,created_at:STRING,created_at_date:STRING,created_at_year:STRING,created_at_month:STRING,created_at_day:STRING,created_at_time:STRING,in_reply_to_user_id_str:STRING,contributors:STRING,retweeted:STRING,truncated:STRING,coordinates:STRING,source:STRING,retweet_count:INTEGER,url:STRING,first_hashtag:STRING,first_user_mention:STRING,screen_name:STRING,name:STRING,followers_count:INTEGER,listed_count:INTEGER,friends_count:INTEGER,lang:STRING,user_location:STRING,time_zone:STRING,profile_image_url:STRING

**Schema** ?

id:INTEGER,created_at:STRING,created_at_date:STRING,created_at_year:STRING,created_at_month:STRING,created_at_day:STRING,created_at_time:STRING,in_reply_to_user_id_str:STRING,contributors:STRING,retweeted:STRING,truncated:STRING,coordinates:STRING,source:STRING,retweet_count:INTEGER,url:STRING,first_hashtag:STRING,first_user_mention:STRING,screen_name:STRING,name:STRING,followers_count:INTEGER,listed_count:INTEGER,friends_count:INTEGER,lang:STRING,user_location:STRING,time_zone:STRING,profile_image_url:STRING
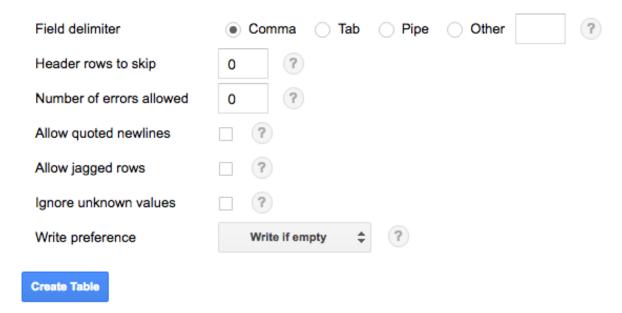
Edit as Fields

Click on Edit as fields to see what that looks like

**Schema**

| Name | Type | Mode | |
|---|---|---|---|
| id | INTEGER | NULLABLE | × |
| created_at | STRING | NULLABLE | × |
| created_at_date | STRING | NULLABLE | × |
| created_at_year | STRING | NULLABLE | × |
| created_at_month | STRING | NULLABLE | × |
| created_at_day | STRING | NULLABLE | × |
| created_at_time | STRING | NULLABLE | × |
| in_reply_to_user_id_str | STRING | NULLABLE | × |
| contributors | STRING | NULLABLE | × |
| retweeted | STRING | NULLABLE | × |
| truncated | STRING | NULLABLE | × |

Check your Options and click Create Table

**Options**

| | |
|---|---|
| Field delimiter | ◉ Comma  ○ Tab  ○ Pipe  ○ Other  [    ]  (?) |
| Header rows to skip | [ 0 ]  (?) |
| Number of errors allowed | [ 0 ]  (?) |
| Allow quoted newlines | ☐  (?) |
| Allow jagged rows | ☐  (?) |
| Ignore unknown values | ☐  (?) |
| Write preference | [ Write if empty ⇕ ]  (?) |

**Create Table**

This is a CSV file so choose Comma as the Field delimiter and enter the Number of errors allowed.  I choose 100 usually.

You will see that the job is running the and table is loading

**Recent Jobs**

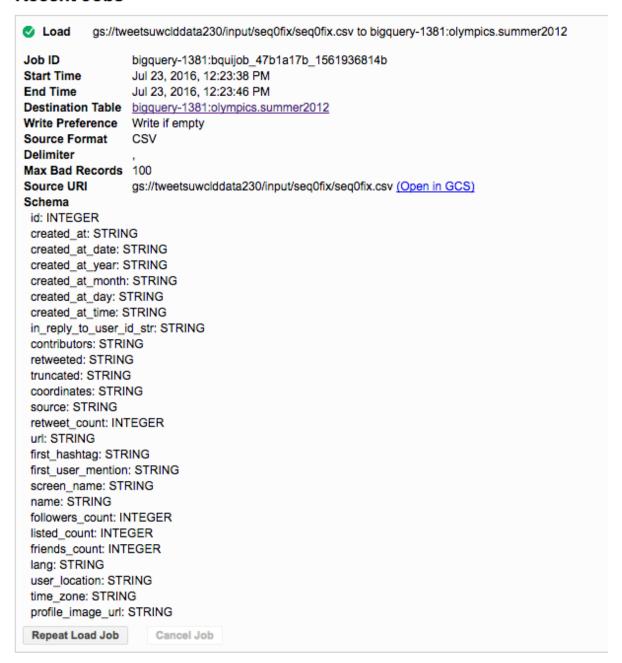✅ Load     gs://tweetsuwclddata230/input/seq0fix/seq0fix.csv to bigquery-1381:olympics.summer2012

Once the load completes you will be able to review the Job History.  If you have errors or exceed the error threshold it will provide guidance on how to solve the issue.

Click on the Repeat load job and load the second file in the same way as you did above.

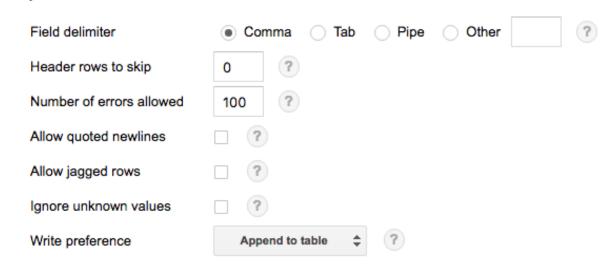Remember to change the file to seq1fix.csv

## Recent Jobs

✅ **Load**    gs://tweetsuwclddata230/input/seq0fix/seq0fix.csv to bigquery-1381:olympics.summer2012

**Job ID**              bigquery-1381:bquijob_47b1a17b_1561936814b
**Start Time**          Jul 23, 2016, 12:23:38 PM
**End Time**            Jul 23, 2016, 12:23:46 PM
**Destination Table**   bigquery-1381:olympics.summer2012
**Write Preference**    Write if empty
**Source Format**       CSV
**Delimiter**           ,
**Max Bad Records**     100
**Source URI**          gs://tweetsuwclddata230/input/seq0fix/seq0fix.csv (Open in GCS)
**Schema**
  id: INTEGER
  created_at: STRING
  created_at_date: STRING
  created_at_year: STRING
  created_at_month: STRING
  created_at_day: STRING
  created_at_time: STRING
  in_reply_to_user_id_str: STRING
  contributors: STRING
  retweeted: STRING
  truncated: STRING
  coordinates: STRING
  source: STRING
  retweet_count: INTEGER
  url: STRING
  first_hashtag: STRING
  first_user_mention: STRING
  screen_name: STRING
  name: STRING
  followers_count: INTEGER
  listed_count: INTEGER
  friends_count: INTEGER
  lang: STRING
  user_location: STRING
  time_zone: STRING
  profile_image_url: STRING

[ Repeat Load Job ]    [ Cancel Job ]

## Create Table

### Source Data

| Location | Google Cloud Storage ⇕ | gs://tweetsuwclddata230/input/seq0fix/seq1fix.csv | ? |
| File format | CSV ⇕ | | View Files |

Also make sure you change the Write preference in Options to **Append to table**



Click **Create Table**



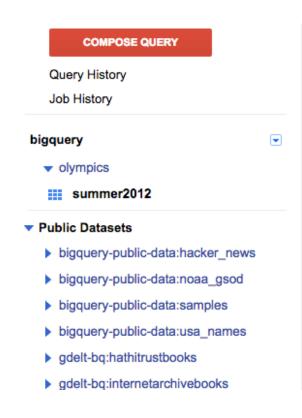If you click on the table summer2012 you can see the Schema



As well as the Details

and a Preview of the table



To run a query, click on the **COMPOSE QUERY** button

Enter a query like:

Select time_zone, count(id) as numTweets, sum(retweet_count) as sumRetweets,
sum(retweet_count)/count(id) as averageRetweetPerTweet
from olympics.summer2012
where time_zone <> '\\N'
group by time_zone order by 2 desc;


Then click the Run Query button

New Query  (?)

```
1  Select time_zone, count(id) as numTweets
2  , sum(retweet_count) as sumRetweets
3  |, sum(retweet_count)/count(id) as averageRetweetPerTweet
4  from olympics.summer2012
5  where time_zone <> '\\N'
6  group by time_zone order by 2 desc;
```

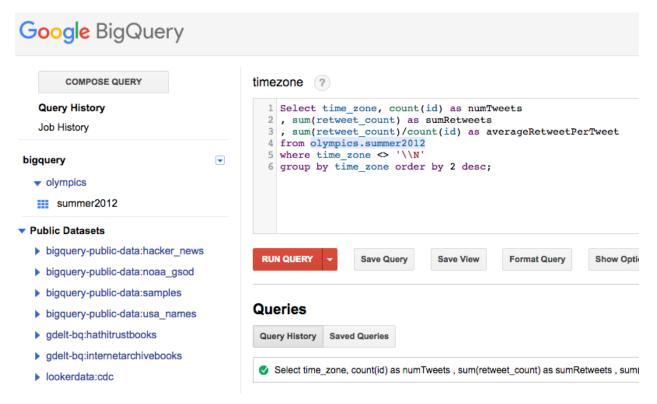**RUN QUERY** ▼     Save Query     Save View     Format Query     Show Options

Which will return records.

Results    Explanation

| Row | time_zone | numTweets | sumRetweets | averageRetweetPerTweet |
|---|---|---|---|---|
| 1 | London | 78913 | 18007430 | 228.19345355011214 |
| 2 | Eastern Time (US & Canada) | 41474 | 8415849 | 202.91867193904616 |
| 3 | Central Time (US & Canada) | 38355 | 10887792 | 283.8689088775909 |
| 4 | Amsterdam | 23320 | 8008156 | 343.40291595197255 |
| 5 | Pacific Time (US & Canada) | 22627 | 5245110 | 231.80757502099263 |
| 6 | Quito | 22268 | 5907937 | 265.3106251122687 |
| 7 | Hawaii | 13469 | 5260572 | 390.5688618308709 |

If you like the Query you can save it by clicking the Save Query button, name it and click OK

You can retrieve the saved query and other recent queries by clicking on Query History in the Left navigation panel



Have fun writing some queries and look at homework exercises over the weekend.

Also feel free to try on the Public Datasets