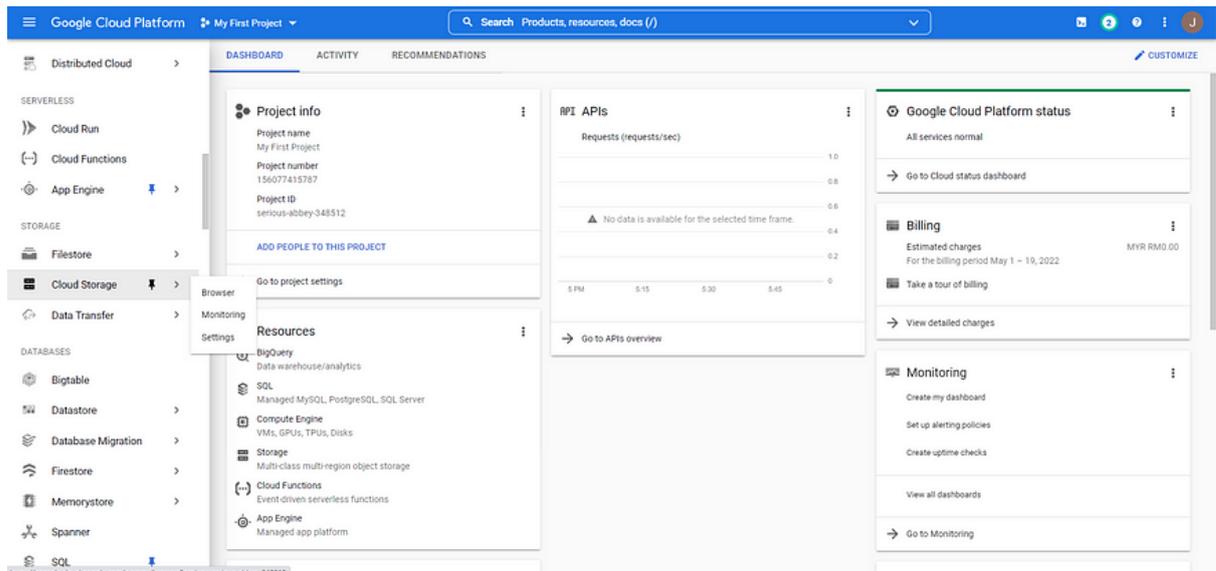# Data Analytics Case Study

# (using BigQuery SQL)

This separate article would document how I used BigQuery SQL instead of RStudio for the analysis process of Prepare & Process.

As a side note, I will not be re-explaining the context of the scenario again, so if you haven't read my original post already, I highly recommend doing so over here!
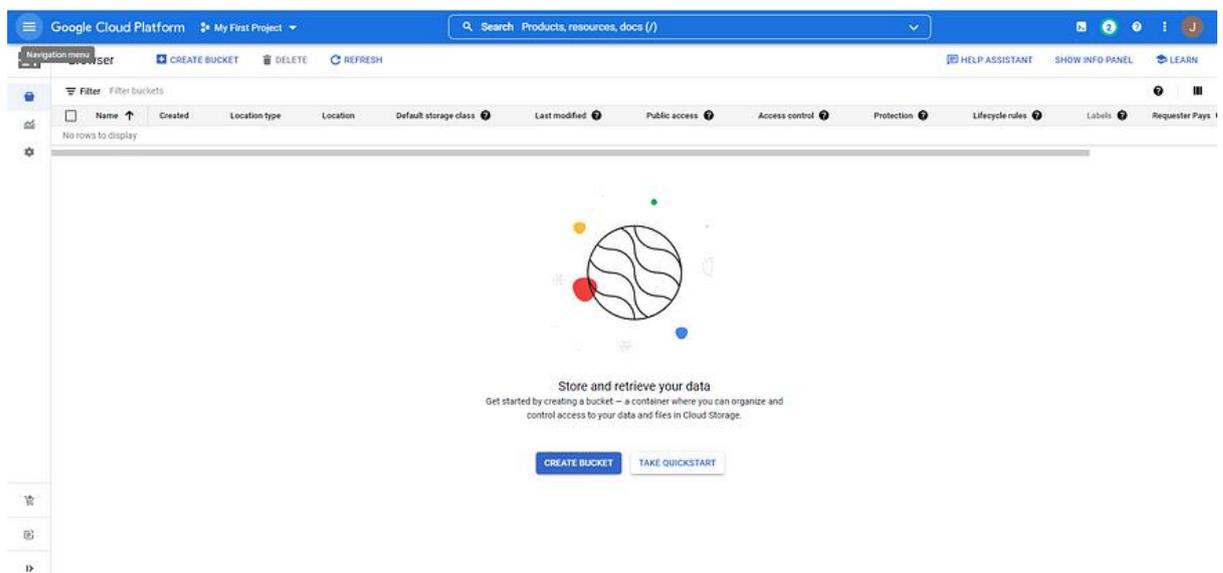
The first section of the article would be dedicated to how I uploaded the datasets for use in BigQuery's system, as so far, the majority of the other students that enrolled in this course have had trouble doing so. Hopefully, this will serve as a guide as well!

The rest of the article would be documenting how I used BigQuery SQL to process the data.
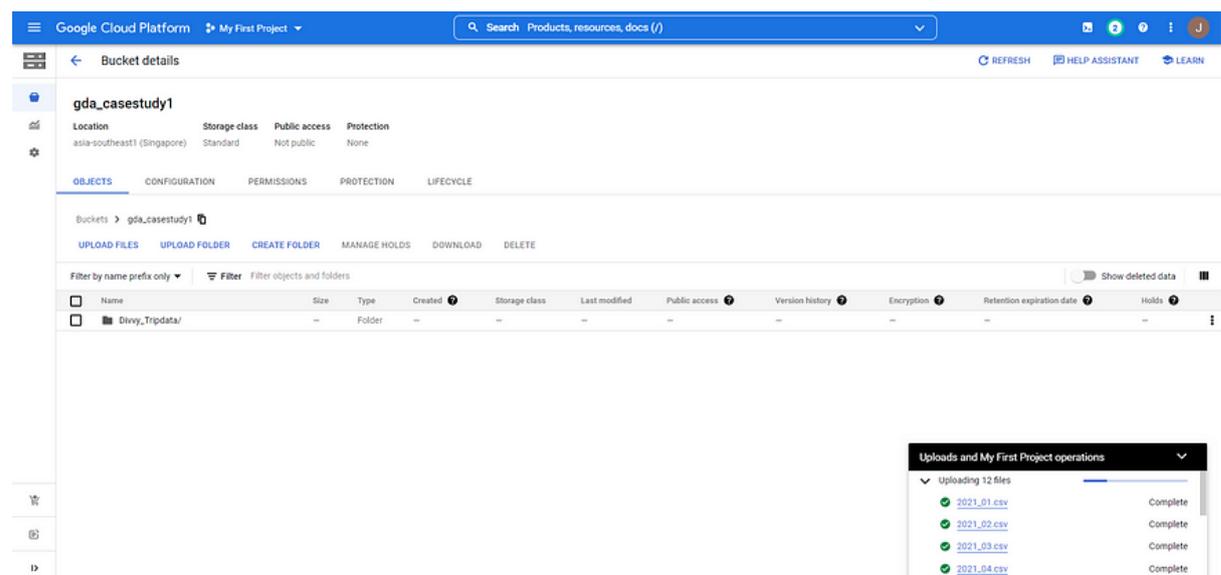
Prepare - Uploading datasets

First, log in to Google Cloud Platform and click the navigation menu on the top left, scroll down till you see the 'Storage' section, and press on 'Cloud Storage'.

Next, create your bucket and give it a name as well.

Now we need to choose where to store our data.

Select 'Region — lowest latency within a single region' and select the region you're the closest to. Proceed with the basic settings for the rest and create your bucket.



Now, upload the files/datasets that you've obtained from the pdf . I would recommend uploading/creating a folder first so that you can organize your data.

Accessing datasets

Now, head over to your BigQuery dashboard/homepage and create a table. Select the Google Cloud Storage option when you press the dropdown menu from 'Create table from'

gda_casestudy1 is my folder name and the file are contained inside

Import the files one by one into your table (remember to name them as well) and make sure to enable 'Auto Detect Schema'

This whole process will take about 10 minutes.
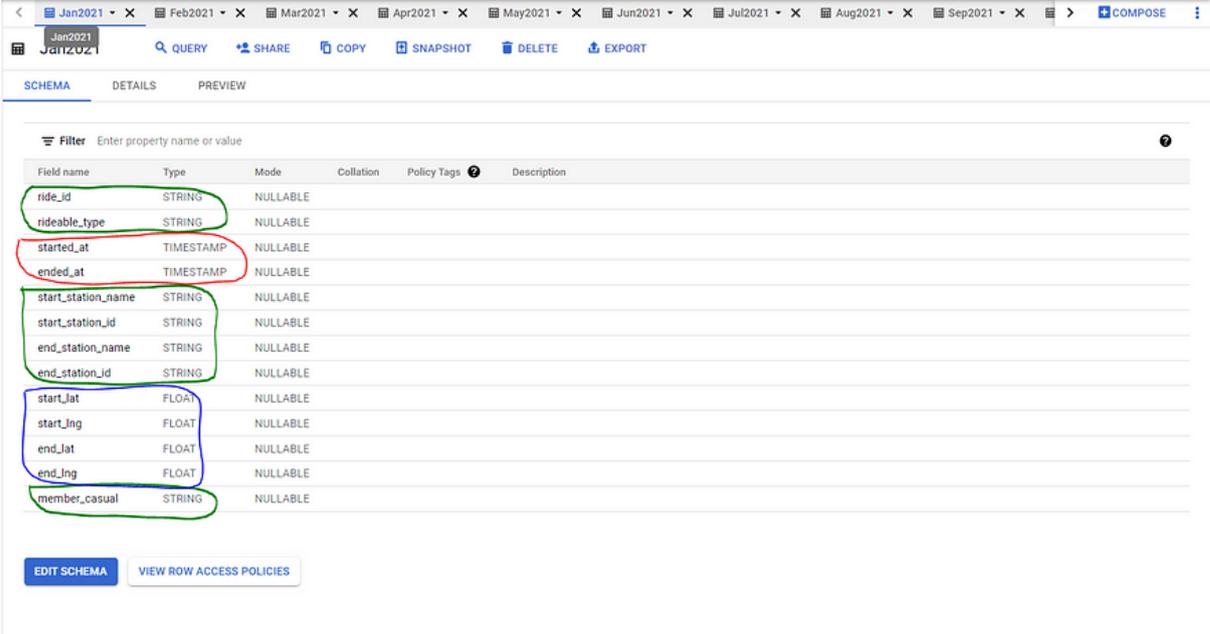
Now that we've imported our datasets, it's time to move on to the Prepare phase!

Process

First, let's check the schema of all the tables that we've imported. We should check if:

- The format of each field is identical
- The naming of each field is identical



After confirming its identical, we will merge all the tables together into one dataset(which we will be calling dataframe from now on) by using UNION ALL. As to why we're using this instead of JOIN, joins will combine data into new columns, which means in our final dataset, we would have ride_id,ride_id2,ride_id3, and so on.

Unions, on the other hand, will combine new data into new rows while staying in the same column (given that the column names are identical).

more information regarding unions can be found here.

The syntax for UNION is shown as the following:

SELECT column_name(s) FROM table1

UNION ALL

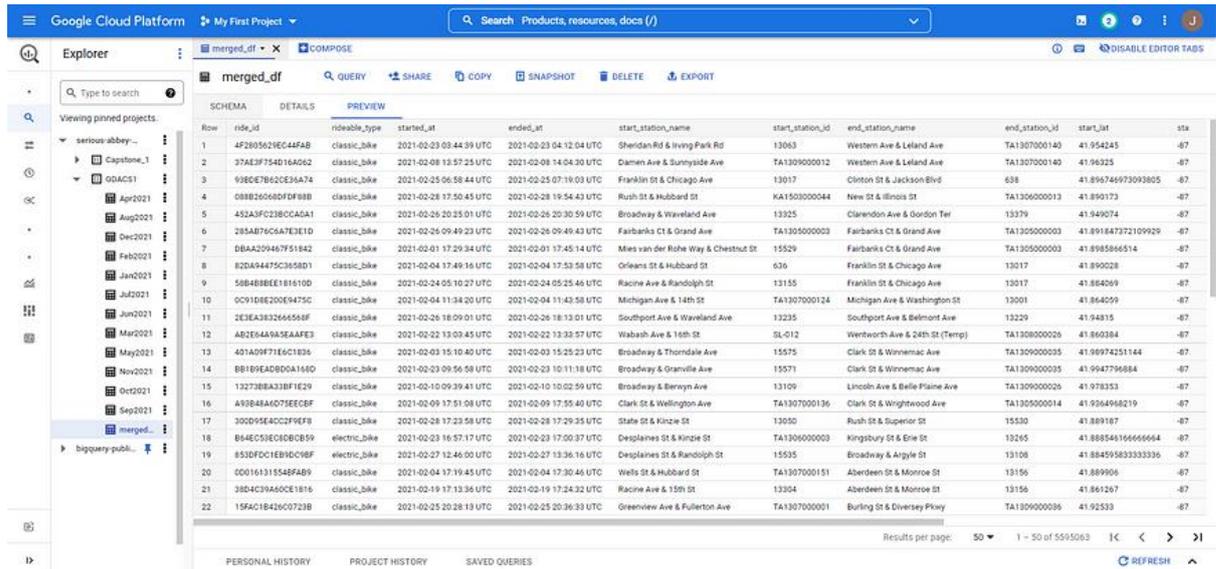SELECT column_name(s) FROM table2;

This is how my query looked like.

Now let's save our query as a new table called 'merged_df', which we do by pressing more > query settings: set a destination table and name the table.

The new dataset, 'merged_df'

Based on the current information available, it is just simply not enough to perform more intricate analysis, therefore we need to create more columns with the following:

- Day of week — By using EXTRACT() & CASE (explained shortly)

for more information, click on EXTRACT & DAYOFWEEK

- Starting hour & Month — By using EXTRACT()
- Trip duration — By using TIMESTAMP_DIFF

for more information, click on TIMESTAMP_FUNCTIONS

I will be saving our new query by overwriting the old one. We also need to check if our new columns have incorrect formatting as well.

| Field name | Type | Mode | Collation | Policy Tags | Description |
|---|---|---|---|---|---|
| ride_id | STRING | NULLABLE | | | |
| rideable_type | STRING | NULLABLE | | | |
| started_at | TIMESTAMP | NULLABLE | | | |
| ended_at | TIMESTAMP | NULLABLE | | | |
| start_station_name | STRING | NULLABLE | | | |
| start_station_id | STRING | NULLABLE | | | |
| end_station_name | STRING | NULLABLE | | | |
| end_station_id | STRING | NULLABLE | | | |
| start_lat | FLOAT | NULLABLE | | | |
| start_lng | FLOAT | NULLABLE | | | |
| end_lat | FLOAT | NULLABLE | | | |
| end_lng | FLOAT | NULLABLE | | | |
| member_casual | STRING | NULLABLE | | | |
| day_of_week | STRING | NULLABLE | | | |
| starting_hour | INTEGER | NULLABLE | | | |
| month | INTEGER | NULLABLE | | | |
| trip_duration | INTEGER | NULLABLE | | | |

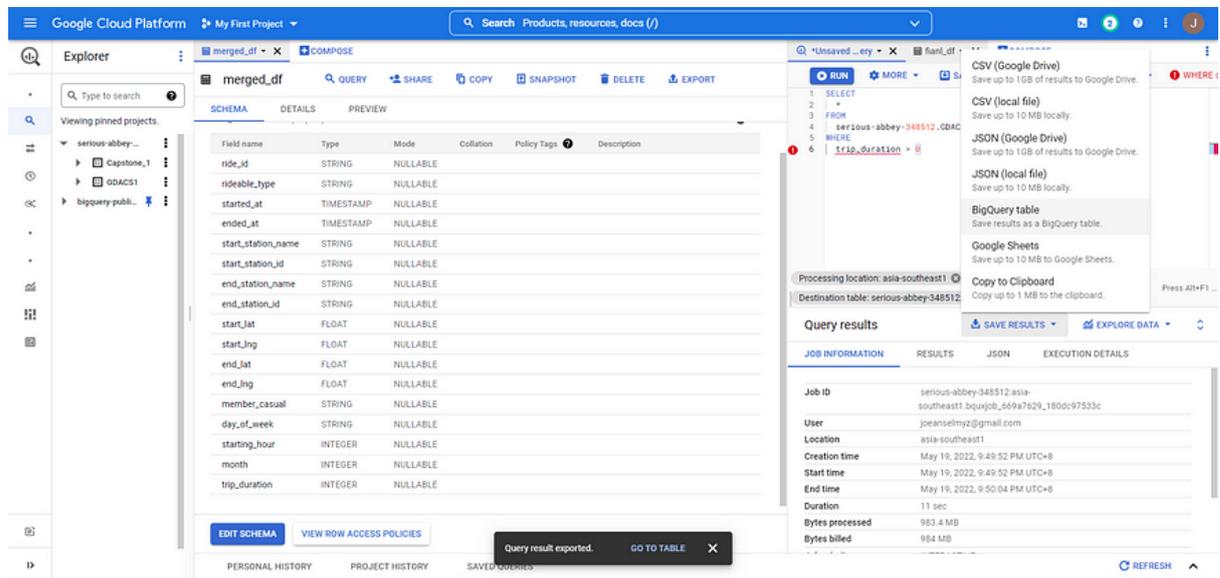All seems good. We will be filtering all trip_durations which are 0 seconds and less by using the following query:

SELECT * FROM [table_name] WHERE trip_duration > 0

Share

Now that we're done, it's time to export the file for visualization using Tableau.

First, click on your latest query, and save the results to a BigQuery table. There may be compatibility issues, where you might need to

create a new dataset (i named my exports) to save the new dataframe into.



Next, open your newly created table(the latest one), and press export to GCS as a CSV with or without GZIP compression (compress to save bandwidth but exports will take a lot longer about 20–30 minutes).

Select your output destination, and again, give your file a name (ENDING IN .CSV). The exporting process can take up to 30minutes. have a quick bite or get a cup of coffee while you're at it.

After it's done, press the navigation menu, open your Google Cloud Storage, navigate to your file destination, and you can download it to share, or to use for visualization.

To follow up with the visualization, do head over to my main article where i use Tableau to create striking visuals from this dataset over here!