# CS4830 - BIG DATA LABORATORY
# LAB - 5
# VISHAL RISHI MK - CH18B013

## 1. Download the dataset and upload it into your bucket.

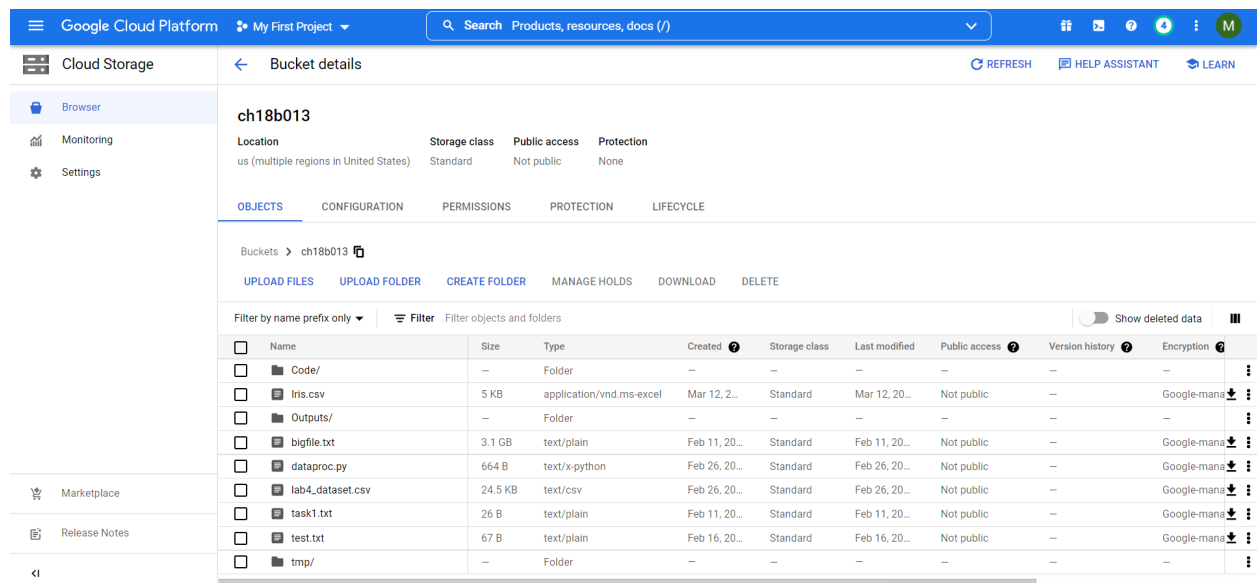The screenshots after uploading the iris dataset in the bucket and as a BigQuery table are shown below.



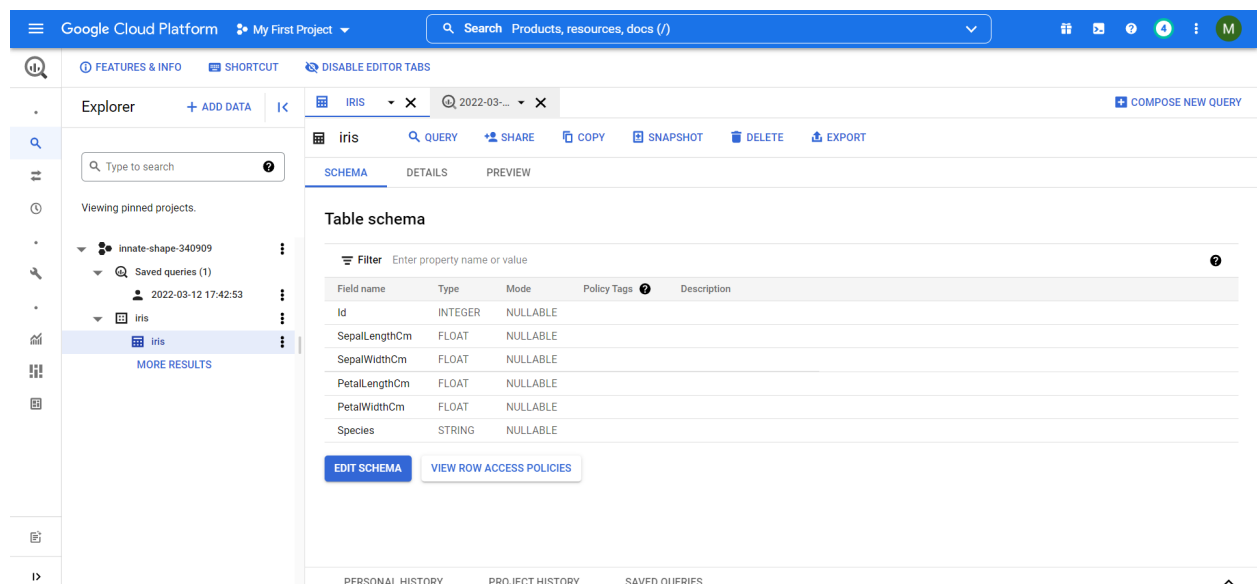*Figure 1: Uploading the iris dataset in the bucket*



*Figure 2: Uploading the iris dataset as a BigQuery table*

**2. Count the number of Iris Virginica flowers that have sepal width greater than 3 cm and petal length smaller than 2 cm. (use BigQuery)**
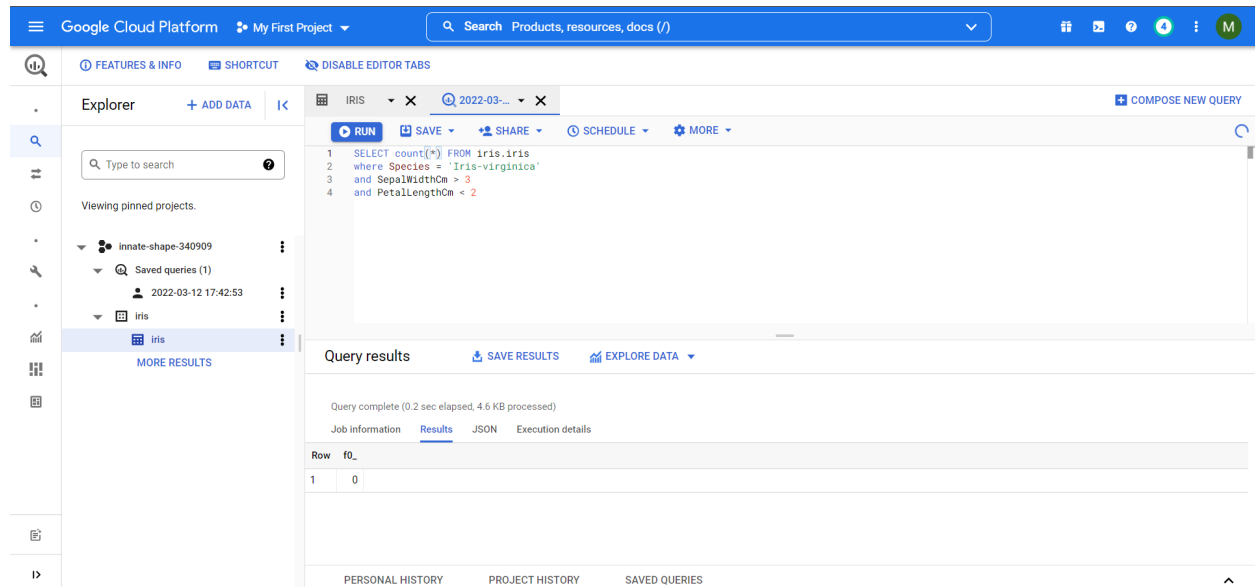The screenshot below shows the SQL query and the query result.



*Figure 3: SQL query and the result (count: 0)*

Hence, using BigQuery on the Iris dataset, we get that there are no Iris Virginica flowers that have sepal width greater than 3 cm and petal length smaller than 2 cm. **The answer is zero.**

**3. Train an ML model on the dataset to predict the class of the iris plant and report the accuracy for different preprocessing techniques and models. Provide the details of data exploration and feature engineering steps.**
Here, we train the dataset to find the species of the flowers using the following steps.

Step 1: We load the data from the BigQuery Table.
Step 2: We separate the data into features (all columns except 'Species') and labels (the column 'Species').
Step 3: We create different pipelines for feature engineering and learning. The details are given below.
Step 4: We split the data into train and test sets.
Step 5: We train different pipelines on the training dataset.
Step 6: We evaluate the pipelines on the test set.

Pipeline 1:
We normalize the features using Standard Scaler. The normalized features are transformed into three-dimensional vectors using PCA. The transformed features are then fed into a random forest classifier to learn the patterns in the data.

Pipeline 2:
We normalize the features using Standard Scaler. The features are then fed into a random forest classifier to learn the patterns in the data.

Pipeline 3:
We normalize the features using Standard Scaler. The normalized features are transformed into two-dimensional vectors using PCA. The transformed features are then fed into a random forest classifier to learn the patterns in the data.

Pipeline 4:
We normalize the features using Standard Scaler. The normalized features are transformed into one-dimensional vectors using PCA. The transformed feature is then fed into a random forest classifier to learn the patterns in the data.

Given below is the output of the job in data proc:



*Figure 4: Output of the job in data proc*

**Accuracy on the test dataset - Pipeline 1: 1.0**
**Accuracy on the test dataset - Pipeline 2: 1.0**
**Accuracy on the test dataset - Pipeline 3: 1.0**
**Accuracy on the test dataset - Pipeline 4: 1.0**