# UNIVERSITY OF NORTH TEXAS, DENTON

## Assignment 7:
### Spark for the Machine Learning

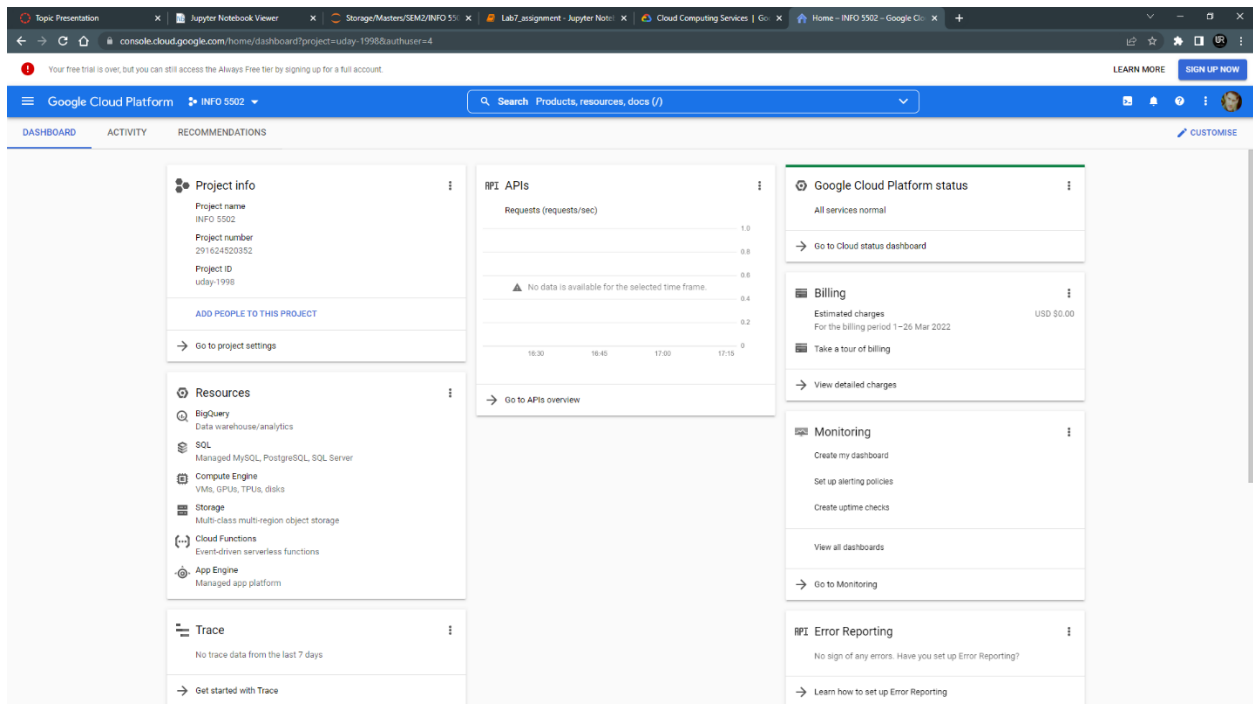Team Name:
Uday Raj Suthapalli
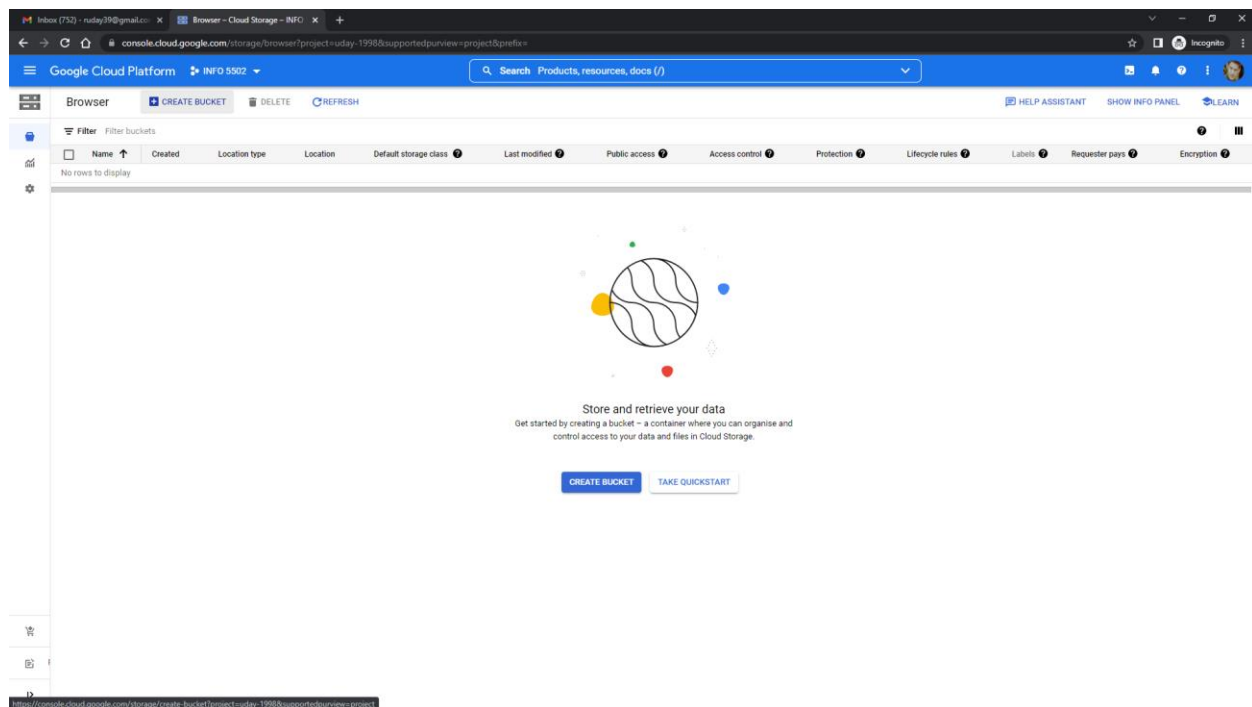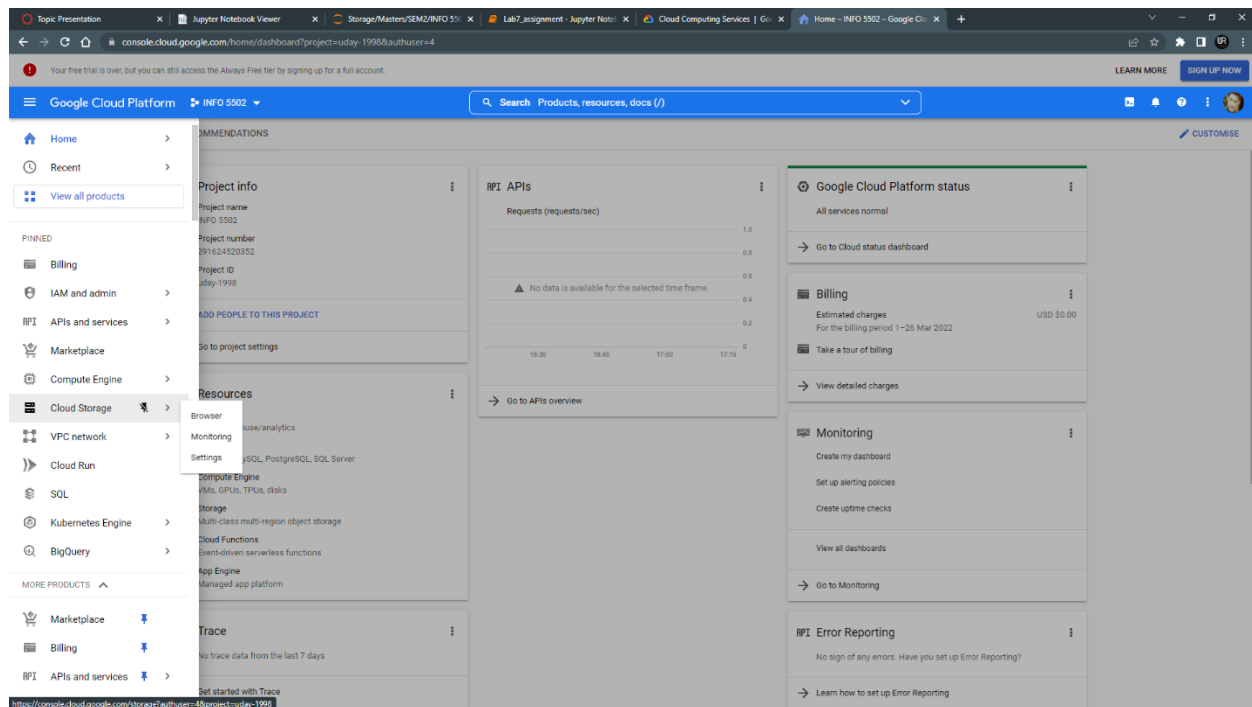Venkat Akshay Reddy Jaggavarapu
Vashisth Hareshkumar Sukhadiya

# Create Storage Bucket in GCP

**Step1: Create a New Storage Bucket.**

- ➢ Open Google Cloud Platform.
- ➢ Go to GCP console.
- ➢ Click on the Navigation Pane on left side.
- ➢ Select Cloud Storage Open Browser tab.
- ➢ Click on the Create Bucket.

## Step2: Create a GCP Bucket.

### Step1:-

> ➢ Name your Bucket (Note: - Do not right 'Uppercase', 'no space' and not like that 'bah-vas-hkasj')

➢ Click on 'Continue' button.

**Step2:-**

➢ Next step 'Choose Where to Store Your Data' as the defualt data (You can not cahange any information) Like, 'Multi-Region' and Location : US.

**Step3:-**

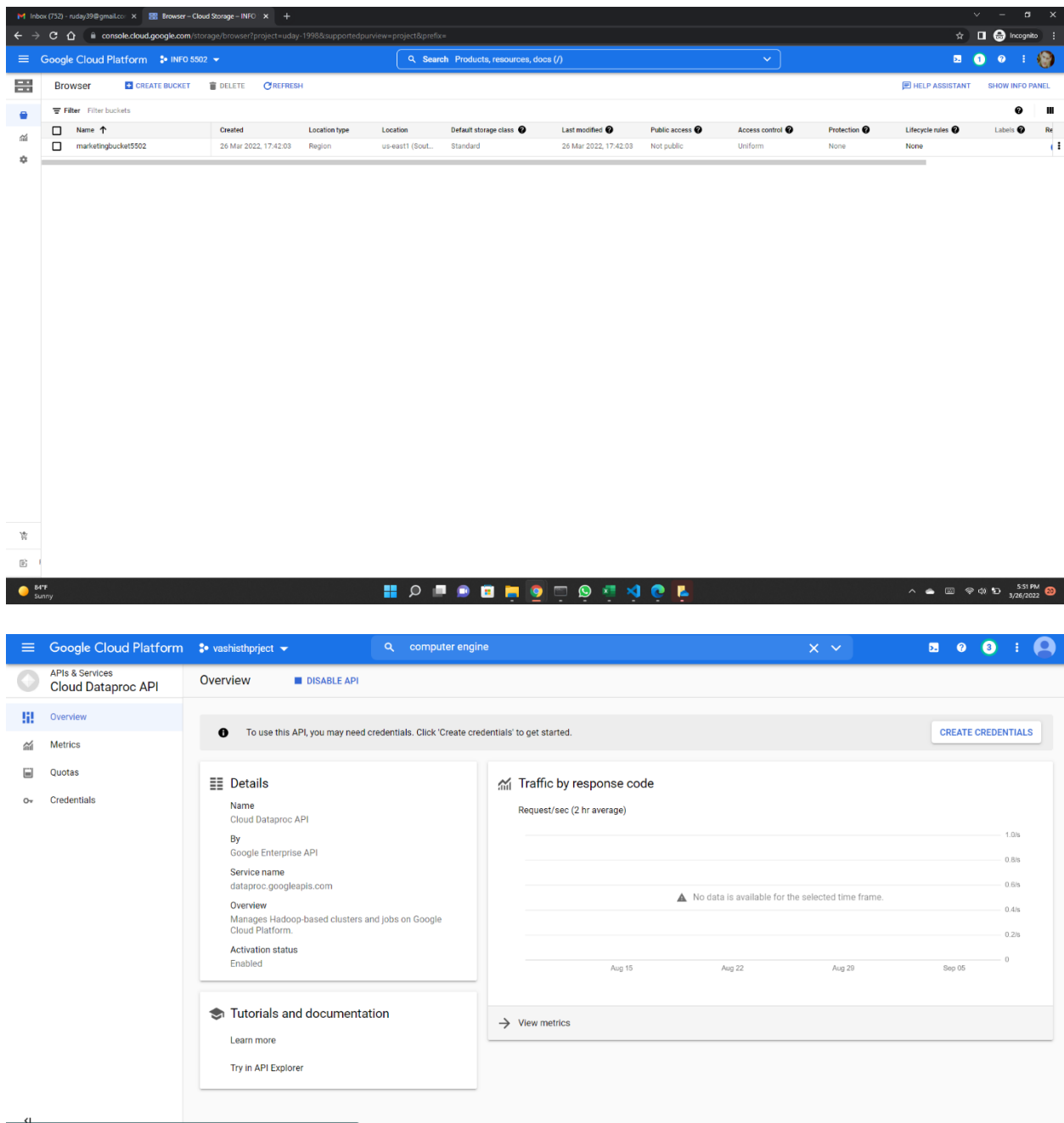➢ 'Choose your defualt storage class for your data'. Select Standard option.

**Step4:-**

➢ **'**Choose how to control access to objects'. Select 'Fine-grained' means Set Object-level and Bucket-level: Permission can be granted at either level – bucket or folder.
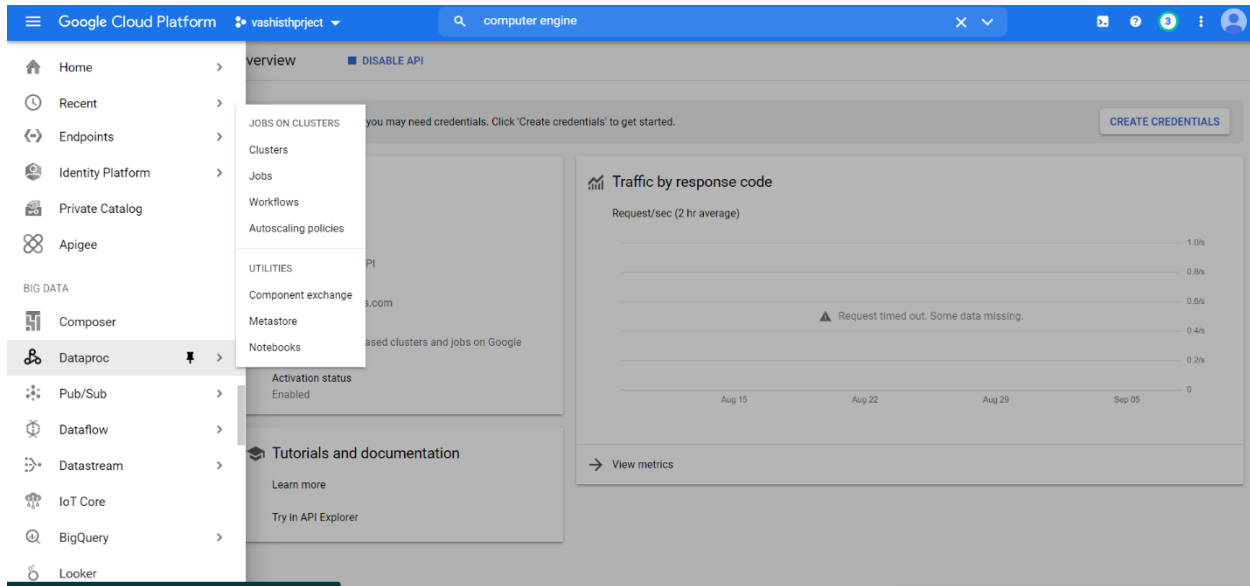


**Step5:-**

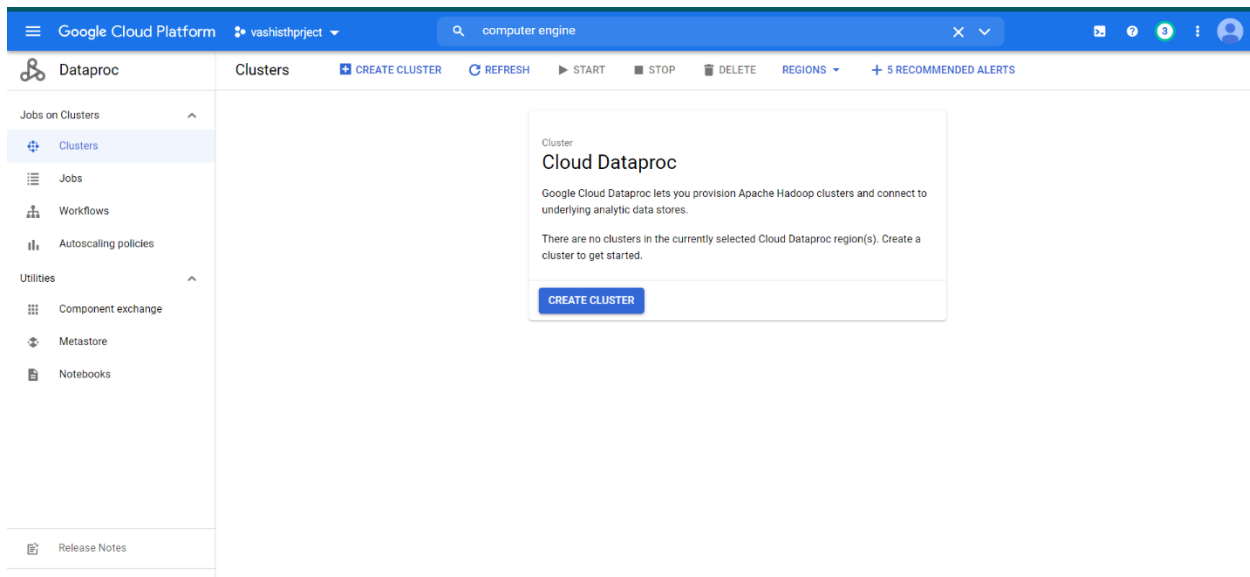➢ Advance Setting (Optional), Click on 'Create' Button.

**Step3: -** Create a GCP DataProc Cluster.

- ➢ Navigate to DataProc
  - • Click the navigation menu to open the menu.
  - • Scroll over and go to the Bigdata, select the "DataProc".
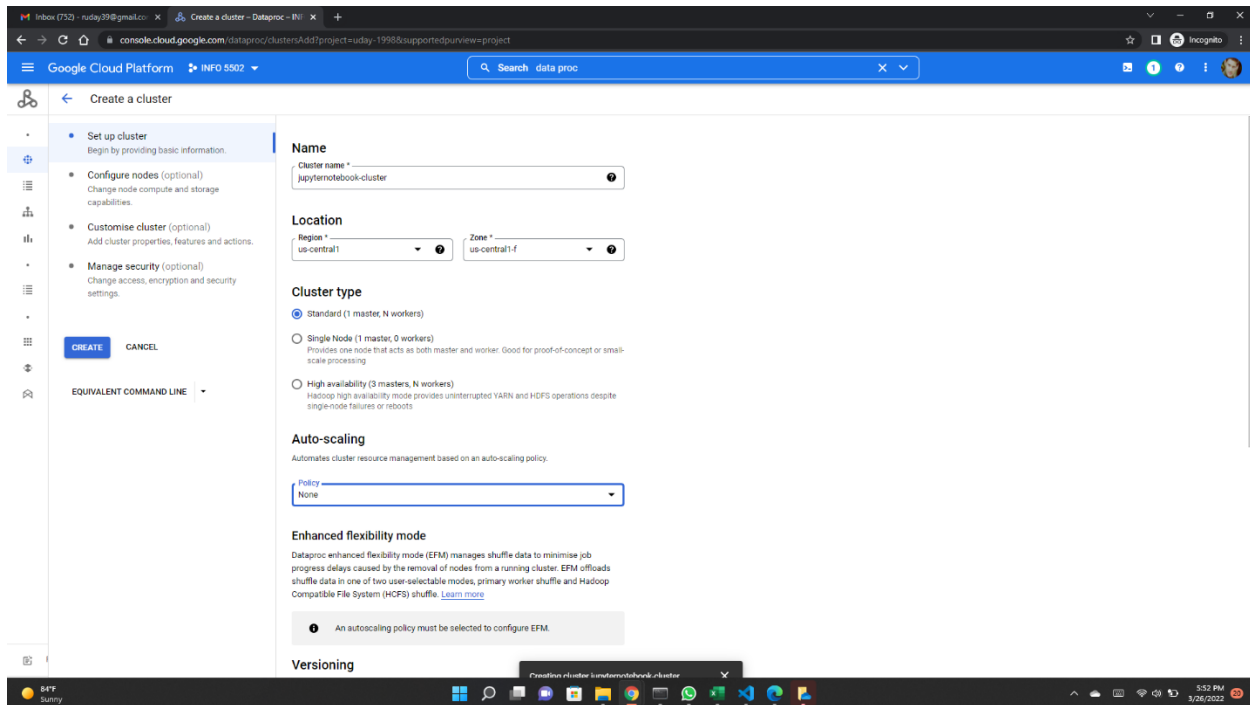  - • Click on "Cluster"
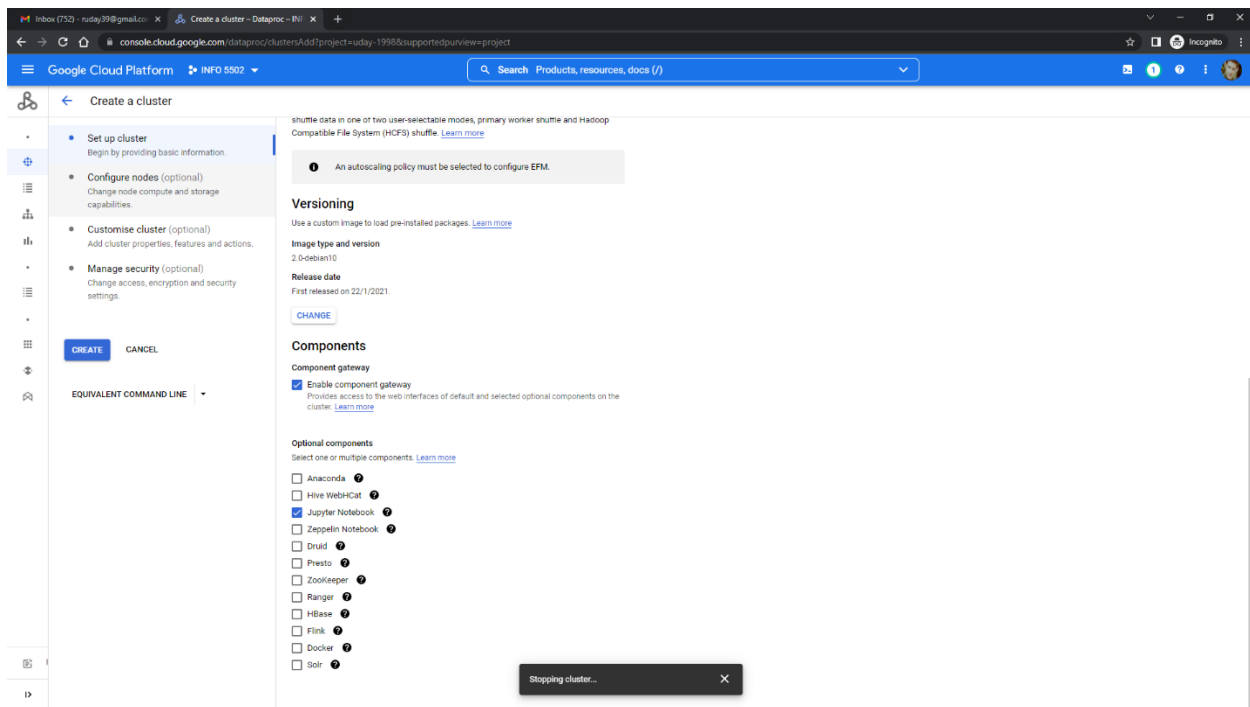
➢ Click "Create Cluster" button.



➢ Set up Cluster
  • Enter "Name" as "jupyternotebook-cluster".
  • Enter "Location" : Region as "us-central1" and Zone as "us-central1-f"
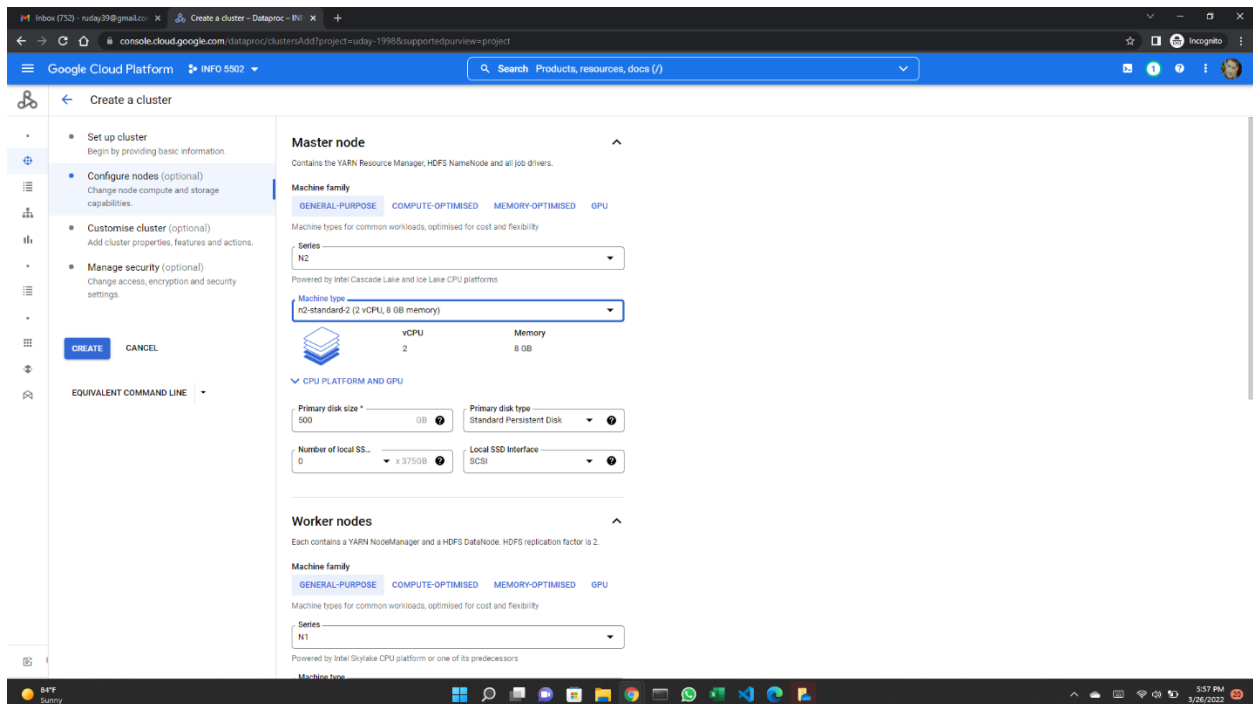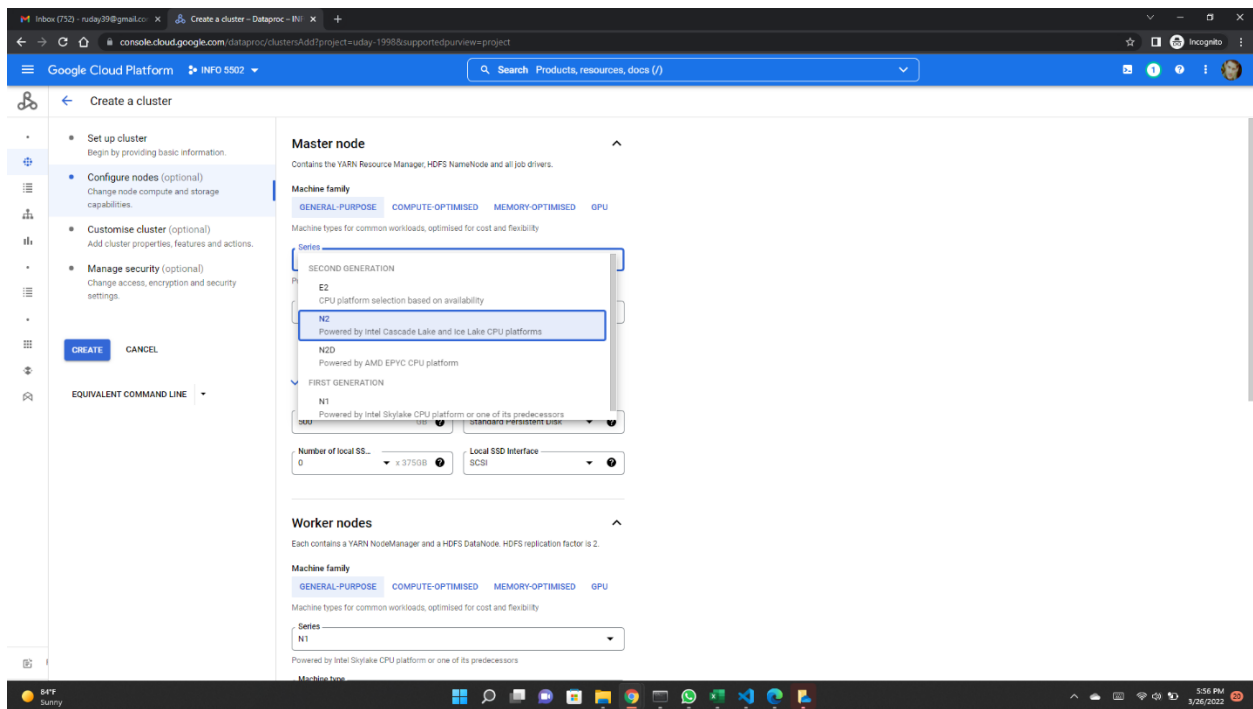  • Enter "Cluster Type" as "Standard (1 master, N workers)".

- Optional component select Jupyter notebook.
- Change system version.
  - Scroll down to "Versioning"
  - Image type and version 2.0-debian10
  - Components: - enable component gateway
  - Optional component we can select Jupyter Notebook
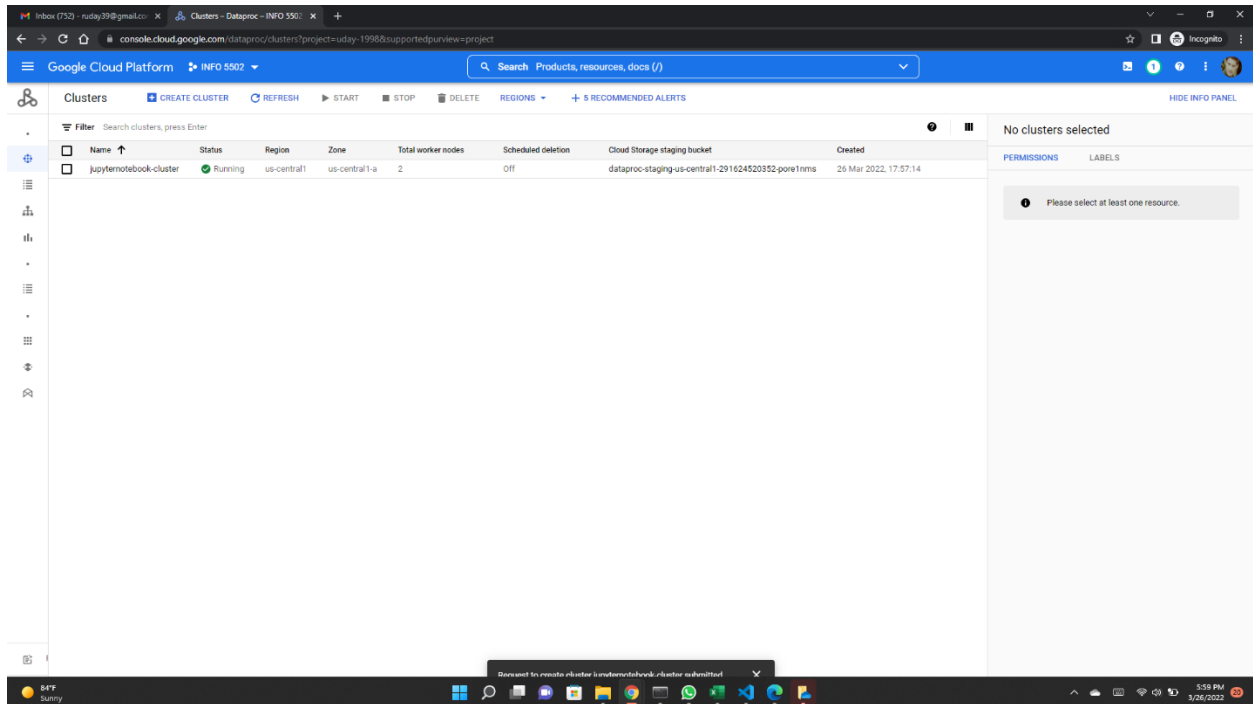  - Here is the screen shot versioning change.

➢ Configure master and work nodes.

- Click "Configure nodes"
- Select "GENERAL-PORPUSE" for the "Machine Family"
- Select "N2" for "Series", Select "n2-standard-2 (2 vCPU, 8GB memory)" for the Machine Type"
- Select "500 GB" for the "Primary disk size (min 10GB)"
- Select "Standard Persistent Disk" for the "Primary disk type"
- Scroll down to "Worker nodes".
- Select "GENERAL-PORPUSE" for the "Machine Family"
- Keep "2" for the "Number of Worker Node"
- Select "N1" for "Series", Select "n2-standard-2 (2 vCPU, 8GB memory)" for the Machine Type"
- Select "500 GB" for the "Primary disk size (min 10GB)"
- Select "Standard Persistent Disk" for the "Primary disk type"

➢ In Cloud Storage Staging Bucket, select browse and click on storage bucket Which is name as marketingbucket5502 .

➢ Click on select

➢ Now, Click on Create



**Step4: - Upload Data to GCP Storage Buckets.**

➢ Open the data folder by clicking on the folder name 'data'.
➢ Click on "Upload files".
➢ Browse for the file to be uploaded, highlight the files, and click "open"
➢ Upload files

➢ We have opened clusters.

➢ We have opened Jupyter Notebook from data proc cluster.
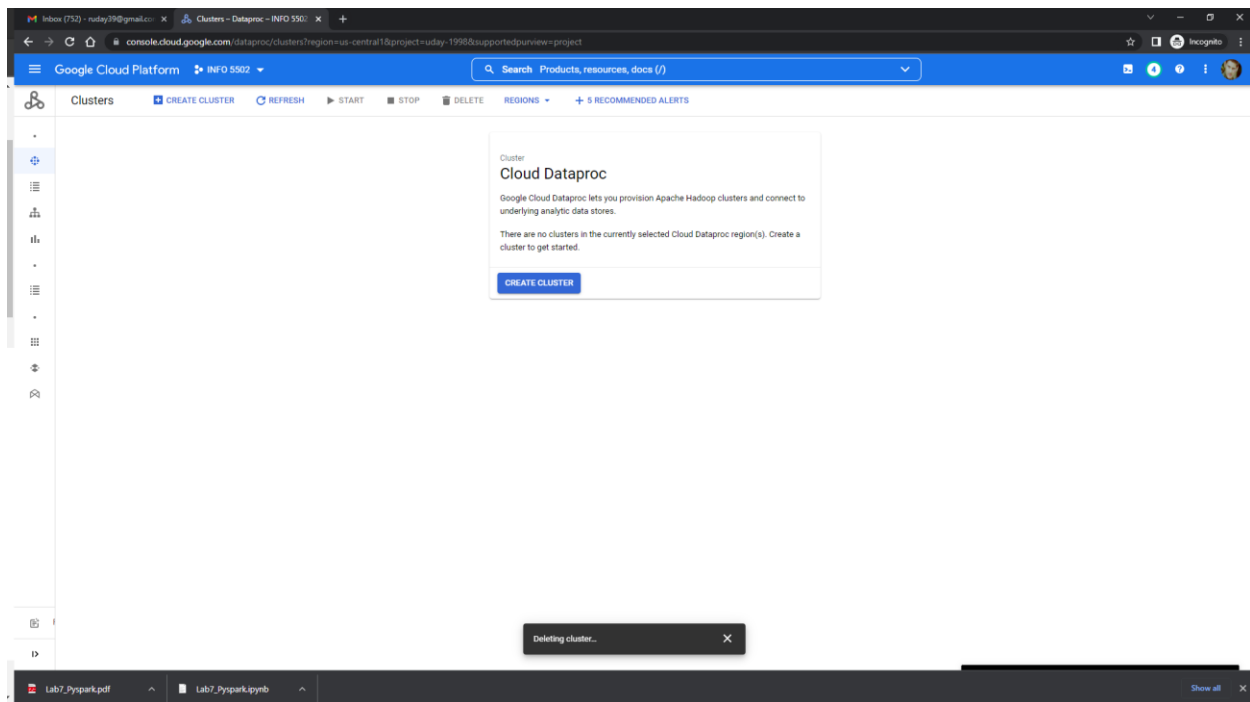
➢ Also, created a logistic ML model using pyspark.

# How to Stop and delete Cluster Nodes and bucket in GCP

**Step1: -** Stop a Cluster Node in GCP

- ➢ Now, select three vertical dots near Connect and the click on the Stop.
- ➢ Then select Stop.
- ➢ As well as Delete cluster.

**Step2: -** Delete Bucket in GCP

- ➢ Delete Cluster which I was created.