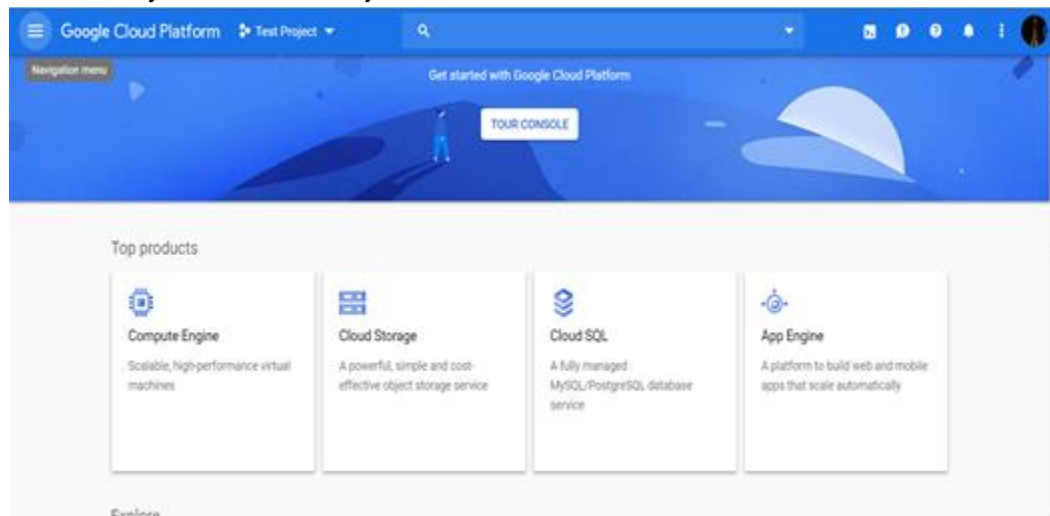# How to Create a VM on Google Dataproc and how to submit a spark job

Google Cloud Dataproc gives the user a provision to implement Apache Hadoop Cluster and also lets user to connect to the various analytic data stored. Using the VM a user can directly submit a spark script to start a job as the Vm environment would be installed with spark, python2 & python3.

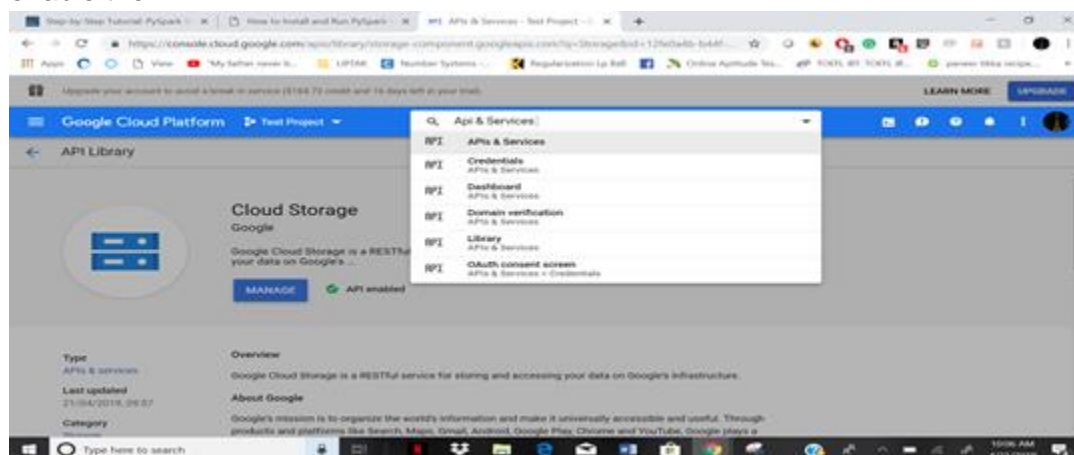Step1: Create a Google Cloud Data platform free account using your Gmail id.
· When you do so you would be given a $300 credit for the next 12 months for usage
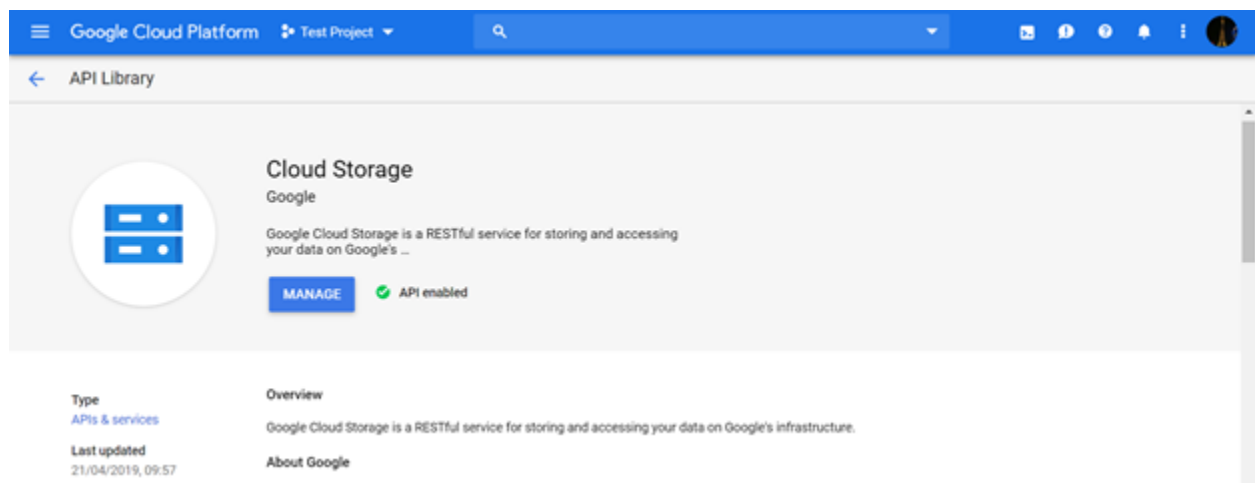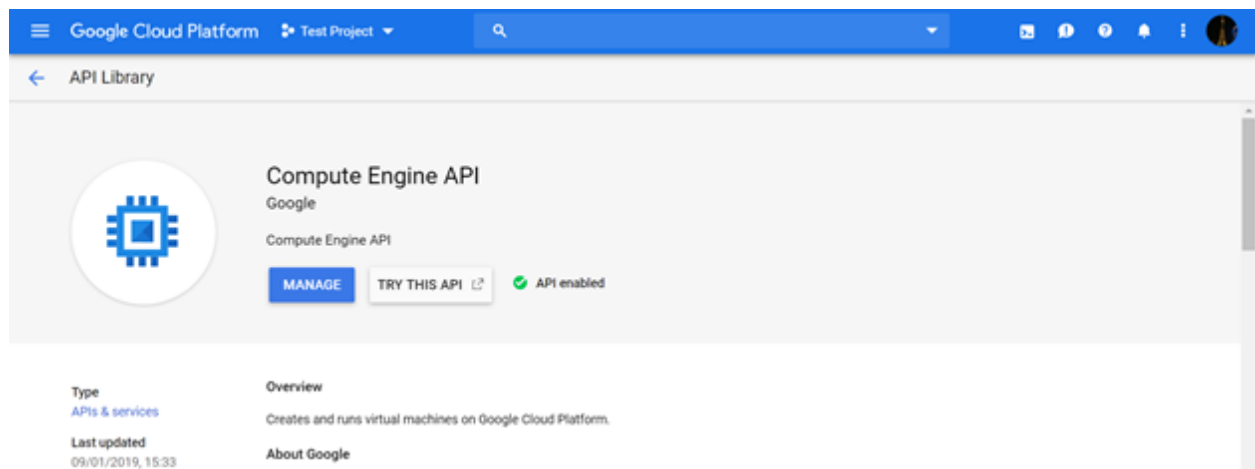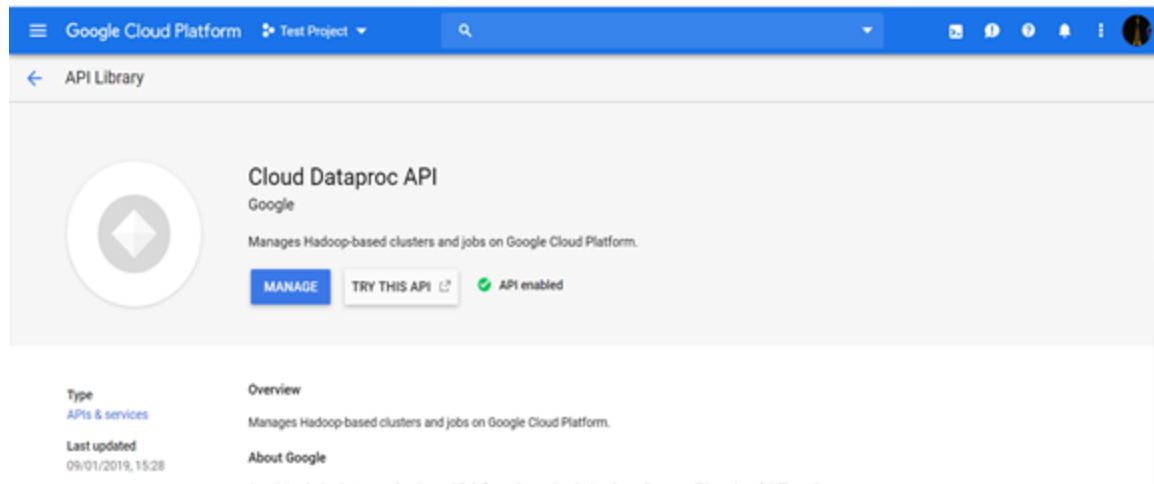· Create you first new Project



On the top you can see 'Test Project' which was created
· **Enabling API's** : We need to enable API to access various like Cloud Dataproc, Compute Engine & Data Storage
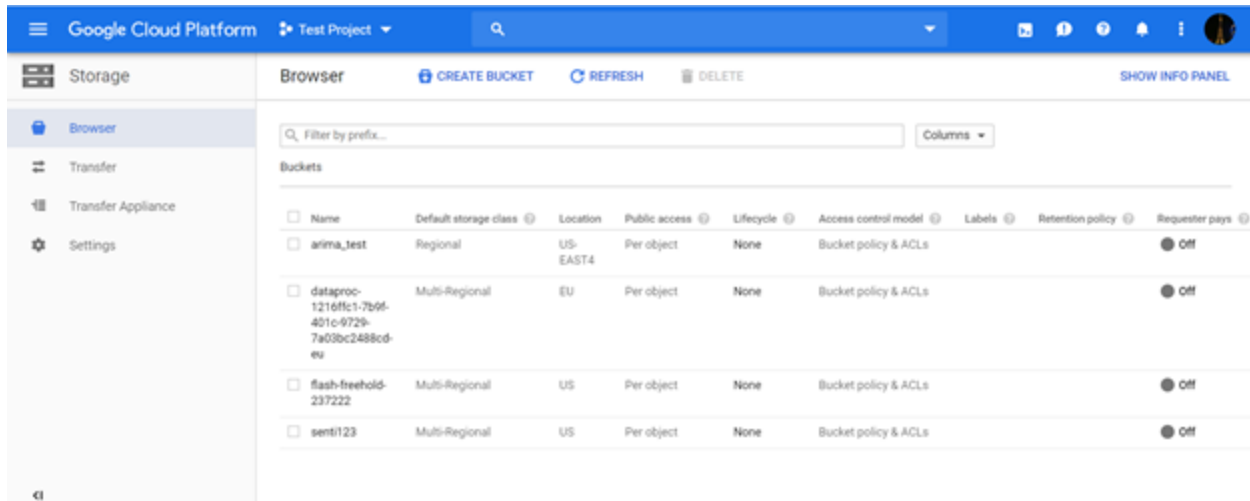To Access we need to API & Services go the Library and search for the above services and enable them
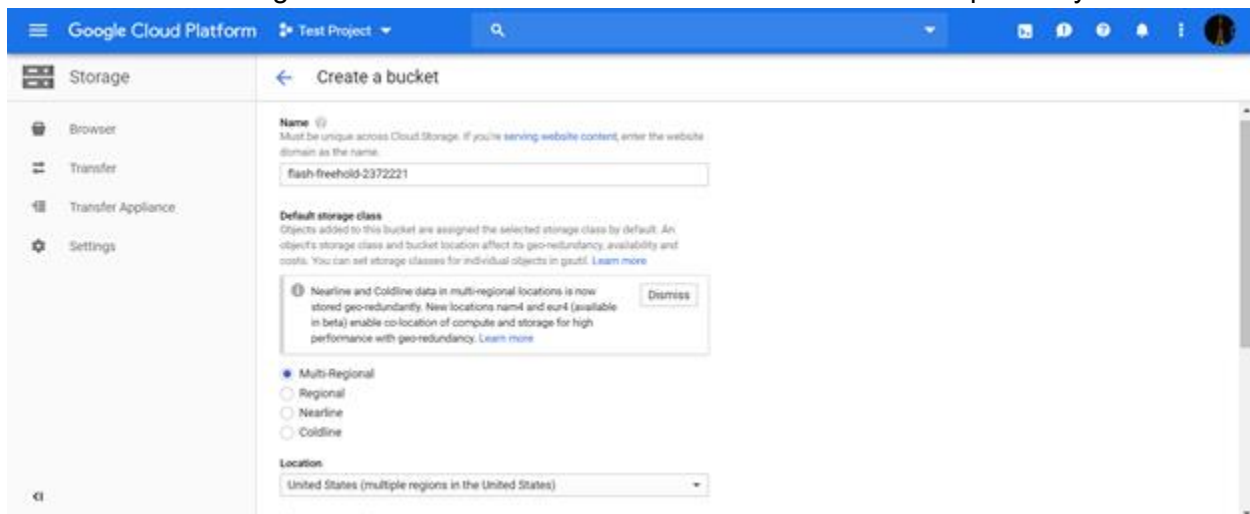
We also may install Google Cloud SDK to access the services using the command line interface

**Step2: Creating A Storage Bucket in Google Cloud Platform**
Go to 'Storage Bucket' in Google Cloud Platform

· On the Top we can see a '**CREATE BUCKET**' option, things to remember when filling in the create bucket page: for easy remember and access try filling the project-Id as the bucket name so won't be confusing or difficult when a user wants to access it. For example in my case



· For easy accessibility for future parts of the code, you can run this simple command on Google Cloud Sdk or your VM instance

> export Project_ID = 'enter_Project_ID'

This Command is for creating a bucket using the command line of Google cloud Platform

> gsutil mb gs://${Project_ID}

This will create a bucket with default settings or you can go to the console and create the bucket as shown earlier.

This is the commands to select the cluster region

```
#ex1) multi-region europe
gsutil mb -l eu gs://${PROJECT_ID}
#ex)2 region europe-west1
gsutil mb -l europe-west1 gs://${PROJECT_ID}
```

**Step 3: Preparing Script to load data into the bucket**

We need to create a script to run our code. This script would be responsible for uploading of data into the cloud storage bucket.
The command to run the script, for example, can be

> ./data_prep.sh or any file name you decide to create

After the script is successfully run, we need to check in the Web console by going into the bucket -> if those files have been uploaded or not.



**How To create a Cloud Dataproc Cluster for running Apache Spark or even Apache Hadoop Clusters**

We can create a new Cluster and give assign the number of worker nodes:



Assign cluster name and number of worker nodes along with memory requirements



By selecting the advance options you can set the Image which basically is the Ubuntu version you want to work within the cluster along with the spark & hadoop version

After doing the above procedure you would have successfully set your Cloud Dataproc cluster

**Step 4: Submit Job**

Click on the Jobs of the Dataprocs page

Click on Submit a Job



You need to give the cloud location of the by specifying the bucket location with the following command in the main class or jar file

 gs://bucket_name/file_location/python_file.py

Also, we can add additional files in the argument section

Click on '**Submit**'

After success full creation of the Job, we can see the result, with the time taken to complete the Job

Important gcloud commands that can be used when you creating a Job using VM instance:
Some important parameter's to remember:

1. Bucket name :
2. Cluster name:
3. Loading  data from a source, for example, git: - git clone 'githib link'
4. After that, we need to check if a file of the GitHub has been created or not : cd filename/
5. Now copying a command into the Google storage bucket: gsutil cp filename.py gs://bucket name/foldername/filename.py
6. Check in the console if the file is uploaded or not
7. Command: gs submit jobs pyspark cluster -- name gs://bucket name/folder name/ filename.py

We have tried implementing it on Google Dataproc but have fallen short of it as we weren't able to solve this error



```
❗ d5bf2f69ac2742cdbf93fdb953101184
Start time: 26 Apr 2019, 22:53:03   Elapsed time: 39 sec   Status:

❗ Google Cloud Dataproc Agent reports job failure. If logs are available, they can be found in 'gs://flash-freehold-237222/google-cloud-dataproc-metainfo/2817a774-de90-49b5-9aef-
  2cab99eff665/jobs/d5bf2f69ac2742cdbf93fdb953101184/driveroutput'.

Output    Configuration

☐ Line wrapping                                                                                                                              Equivalent con

        at com.sun.proxy.$Proxy13.mkdirs(Unknown Source)
        at org.apache.hadoop.hdfs.DFSClient.primitiveMkdir(DFSClient.java:2333)
        ... 26 more
19/04/27 02:53:38 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@91264be{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
19/04/27 02:53:38 WARN org.apache.spark.scheduler.cluster.YarnSchedulerBackend$YarnSchedulerEndpoint: Attempted to request executors before the AM has registered!
19/04/27 02:53:39 WARN org.apache.spark.metrics.MetricsSystem: Stopping a MetricsSystem that is not running
Traceback (most recent call last):
  File "/tmp/d5bf2f69ac2742cdbf93fdb953101184/ARIMA_PySpark.py", line 38, in <module>
    sc =SparkContext.getOrCreate()
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/context.py", line 349, in getOrCreate
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/context.py", line 118, in __init__
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/context.py", line 180, in _do_init
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/context.py", line 288, in _initialize_context
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 1525, in __call__
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/protocol.py", line 328, in get_return_value
py4j.protocol.Py4JJavaError: An error occurred while calling None.org.apache.spark.api.java.JavaSparkContext.
: org.apache.hadoop.hdfs.server.namenode.SafeModeException: Cannot create directory /user/root/.sparkStaging/application_1556243957219_0020. Name node is in safe mode.
The reported blocks 0 needs additional 16 blocks to reach the threshold 0.9990 of total blocks 17.
The number of live datanodes 0 has reached the minimum number 0. Safe mode will be turned off automatically once the thresholds have been reached. NamenodeHostName:freehold-237222-m
        at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.newSafemodeException(FSNamesystem.java:1413)
        at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.checkNameNodeSafeMode(FSNamesystem.java:1400)
        at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.mkdirs(FSNamesystem.java:2989)
        at org.apache.hadoop.hdfs.server.namenode.NameNodeRpcServer.mkdirs(NameNodeRpcServer.java:1096)
```