# Quick set up guide for Jupyter Notebook

# and Pyspark on Google Cloud

**Intro**: Similar to Amazon EMR, Google Dataproc is the platform designed for distributed clusters computation, which is built for big data analysis including Hadoop, Spark, HBase and others.

Compared with common line tools, Jupyter Notebook is an easier way to interact with Spark, as the notebook can easily visualize all the codes and the results.

But by default, Jupyter is not installed on both platforms, and users need to run custom script during cluster initialization if they wish to.

**Advantages of Dataproc compared with EMR:**

|  | Google Dataproc | AWS EMR |
| --- | --- | --- |
| Default Provisioning | 1 minute | 8 minutes |
| Custom Provisioning | 30+ minutes | Less than 5 minutes |
| Accessibility | Both Web shell and SSH | Only SSH |
| Charging cycle | By Minute | By Hour |
| Pricing[1] | $0.24 per hour | $0.336 per hour |

[1] Default Provisioning with 4 vCPUs and 15 GB of RAM

We set up a trial to compare the performance and cost of a typical Spark workload. The trial used clusters with one master and five core instances of AWS's m3.xlarge and GCP's n1-standard-4.
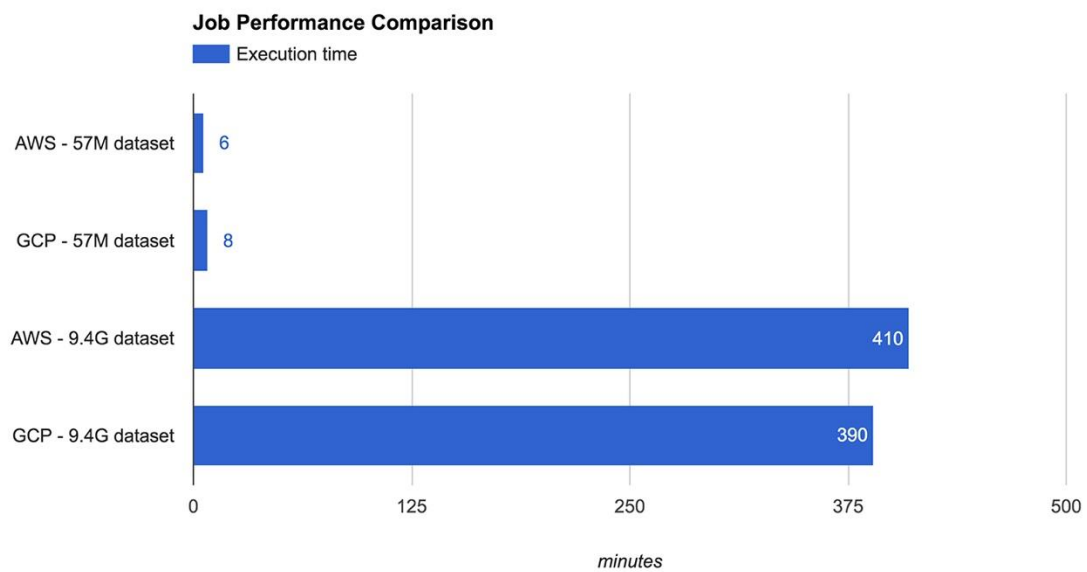
**Job Performance Comparison**

Execution time

| | minutes |
|---|---|
| AWS - 57M dataset | 6 |
| GCP - 57M dataset | 8 |
| AWS - 9.4G dataset | 410 |
| GCP - 9.4G dataset | 390 |

Figure 1. Computation Efficiently Credit: Michael Li and Ariel M'ndange-Pfupfu.

**Job Cost Comparison**

Instance cost    EMR/dataproc

| | dollars |
|---|---|
| AWS - 57M dataset | |
| GCP - 57M dataset | |
| AWS - 9.4G dataset | |
| GCP - 9.4G dataset | |

Figure 2. Cost Comparison Credit: Michael Li and Ariel M'ndange-Pfupfu.

**Set up**

1. You can sign up for free trial on Google Cloud, which would offer $300 credit
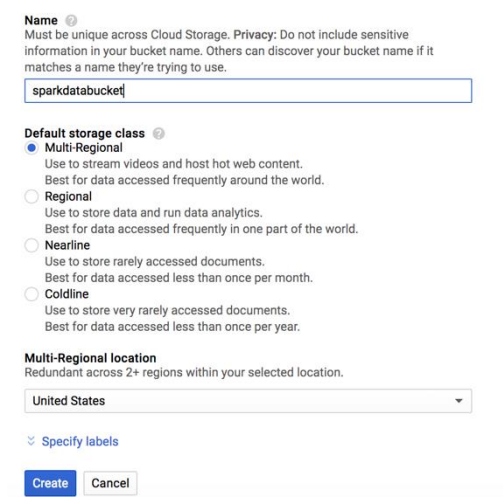and free tier service.



2. Select or create a Cloud Platform project. All the service in Google cloud is
based on 'Project'. And you should also enable billing for the project and the
Cloud Dataproc APIs.



3.Create the Cloud Storage bucket. Multi-regional is for better access.

4. You can use either Cloud SDK(Google Cloud Commend line tool) or Web console to create the cluster. For easier visualization, we would use the web console.



For the 'initialization actions' part, we would use a bash shell initialization script provided by Dataproc team on Cloud Storage. This script and other initialization

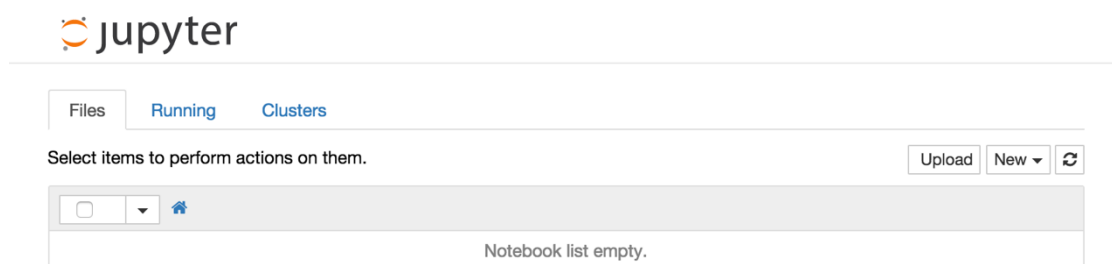scripts are co-located at the GitHub [GoogleCloudPlatform/dataproc-initialization-actions repository](#).

5. Establish SSH tunnel. You can use Cloud SDK to establish the SSH tunnel. It would configure keys and authentication automatically.

*$ gcloud compute ssh "**cluster-name-m**" \\*

  *--project **project-id** \\*

  *--zone=**cluster-zone** \\*

  *-- -D 10000 –N*

6. SSH tunnel supports traffic proxying using the SOCKS protocol, and we need to set up browser to access the tunnel. You should use Terminals or cmd.

*/Applications/Google\ Chrome.app/Contents/MacOS/Google\ Chrome \\*

  *"http://**<cluster-name>**-m:8123" \\*

  *--proxy-server="socks5://localhost:10000" \\*

  *--host-resolver-rules="MAP * 0.0.0.0 , EXCLUDE localhost" \\*

  *--user-data-dir=/tmp*


The opening page of the Jupyter notebook displays in your browser. You can use Jupyter for Pyspark now.

**Reference**

Michael Li &Ariel M'ndange‑Pfupfu. Spark comparison: AWS vs. GCP. O'Reilly, August 30,2016

Install and run a Jupyter notebook in a Cloud Dataproc cluster, Google Cloud Documentation, available at: https://cloud.google.com/dataproc/docs/tutorials/jupyter‑notebook

**From UMN MSBA6330 Trends Marketplace Project by:**

Abhijit Patil: patil074@umn.edu,

Abhilasha Pandey: pande130@umn.edu,

Ankit Agarwal: agarw145@umn.edu,

Ao Liu: liux4399@umn.edu,

Suzanne Kaminski: kamin089@umn.edu