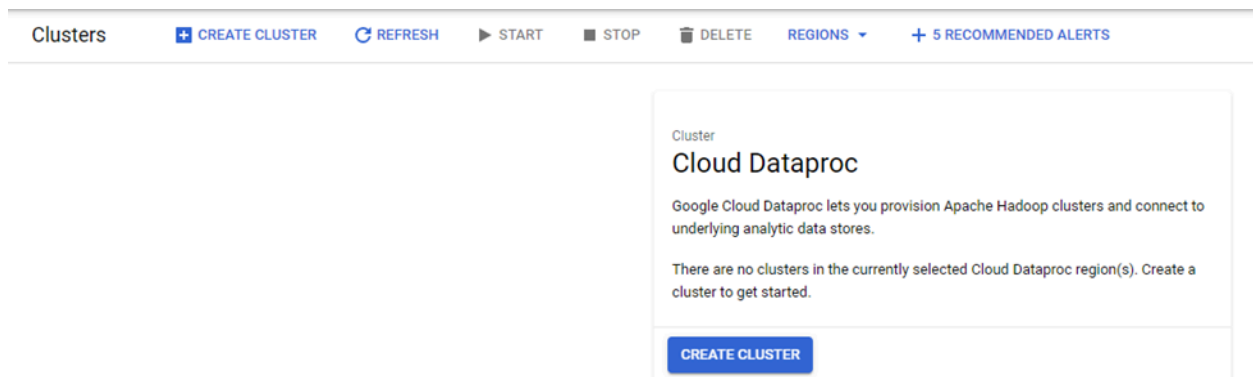# Learning Spark with Google Dataproc and Jupyter Notebook

In this instruction, we will show you how to install and run the Jupyter Notebook with Google Dataproc. Google Dataproc is the managed Hadoop and Spark framework for creating either single-node or autoscaling clusters by the Google Cloud Platform.

## Create a Dataproc Cluster with Jupyter Notebook

Start at the Google Dataproc console to create a new cluster: [console](console)



Start Dataproc cluster creation

Click "Create Cluster" and select the cluster options like Cluster Type, Name of Cluster, Location, Auto-Scaling Options, and more.

Parameters required for Cluster

Select the **Single Node Cluster** option, this means that auto-scaling is disabled as the cluster consists of only 1 master node.

Next, in the "Components" section, be sure to select "Enable Component Gateway" and "Jupyter Notebook". This will enable the web interface of Jupyter Notebook to be available.

# Components

**Component Gateway**

☑ Enable component gateway
Provides access to the web interfaces of default and selected optional components on the cluster. Learn more ⬈

**Optional components**
Select one or multiple components. Learn more ⬈

☐ Anaconda ❓
☐ Hive WebHCat ❓
☑ Jupyter Notebook ❓
☐ Zeppelin Notebook ❓
☐ Druid ❓
☐ Presto ❓
☐ ZooKeeper ❓
☐ Ranger ❓
☐ HBase ❓
☐ Flink ❓
☐ Docker ❓
☐ Solr ❓
☐ Hudi ❓

When you click "Create", it'll start creating the cluster.

After a few minutes the cluster with 1 master node will be ready for use.



| Clusters | | ➕ CREATE CLUSTER | ⟳ REFRESH | ▶ START | ■ STOP | 🗑 DELETE | REGIONS ▾ | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ☰ Filter | Search clusters, press Enter | | | | | | | |
| ☐ | Name ↑ | Status | Region | Zone | Total worker nodes | Scheduled deletion | | |
| ☐ | first-data-proc-cluster | ✅ Running | us-central1 | us-central1-f | 0 | On | | |

Cluster Up and Running

# Open the JupyterLab on Master Node

Since we have created a single-node Dataproc cluster, we can simply launch a Jupyter Notebook on the Master Node to access and control the Spark cluster. We will be using JupyterLab, which is more advanced web interface of Jupyter Notebook.

To access the pre-installed JupyterLab environment, click the name of the Dataproc cluster you just created and switch to the tab of "Web interfaces":



Click the link of "JupyterLab" to access the web interface.

## Check Out the Spark Exercises and Access in Jupyter Notebook

Next, click the Terminal button in the main window to start a Terminal in the browser:
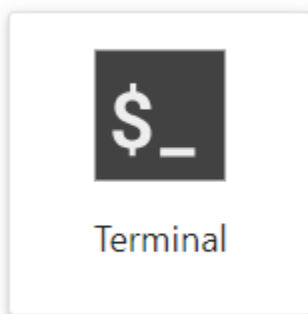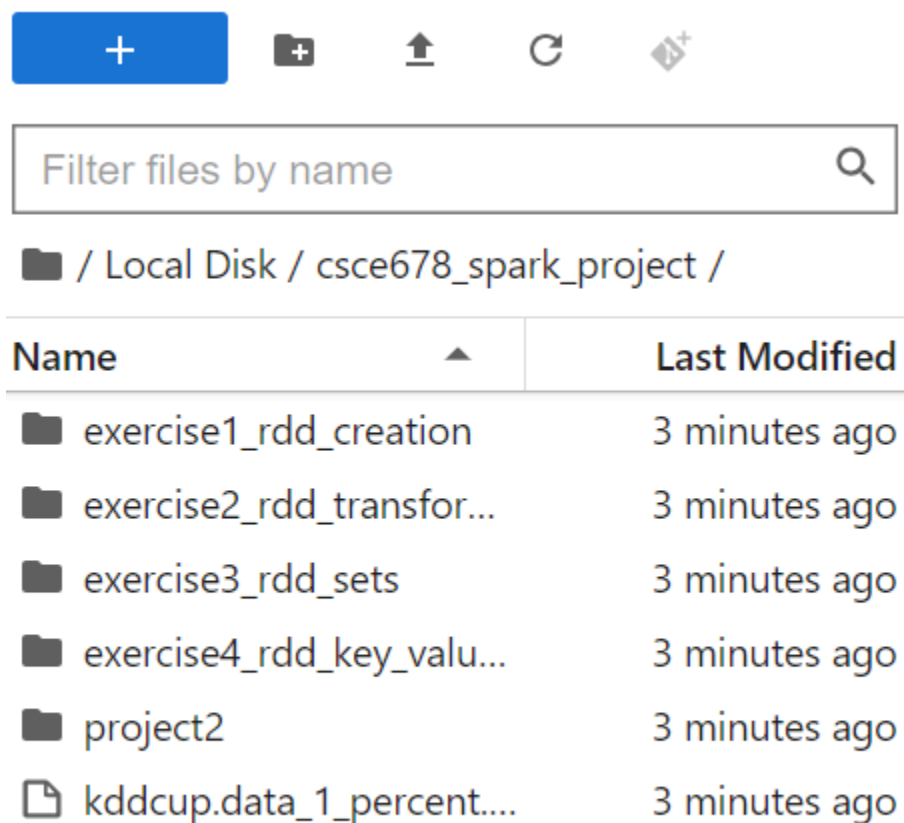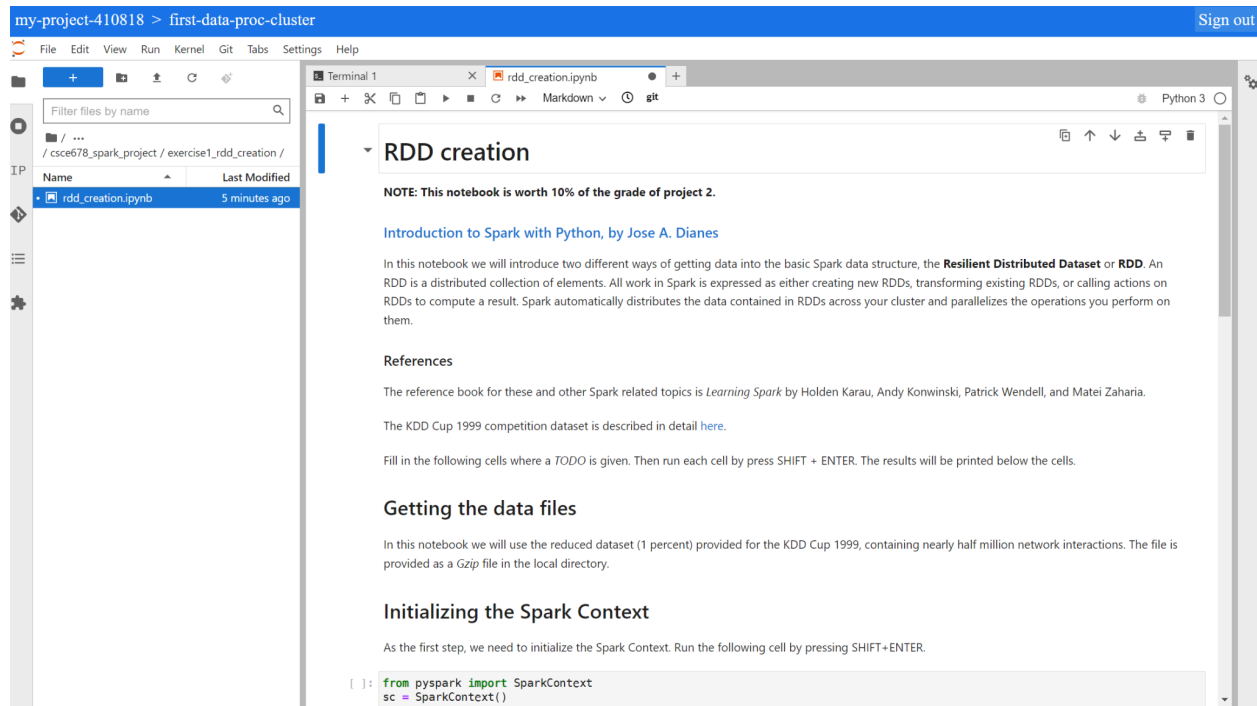


In the Terminal window, enter the following command:

```
git clone https://github.com/chiache/csce678_spark_project
```

```
s_ Terminal 1                    ×   +

root@first-data-proc-cluster-m:/# git clone https://github.com/chiache/csce678_spark_project
Cloning into 'csce678_spark_project'...
remote: Enumerating objects: 91, done.
remote: Counting objects: 100% (91/91), done.
remote: Compressing objects: 100% (51/51), done.
remote: Total 91 (delta 35), reused 67 (delta 22), pack-reused 0
Receiving objects: 100% (91/91), 439.25 KiB | 4.31 MiB/s, done.
Resolving deltas: 100% (35/35), done.
root@first-data-proc-cluster-m:/# ▊
```

Use the file browser in the left panel and naviate to / Local Disk / csce678_spark_project:



| Name | Last Modified |
|---|---|
| 📁 exercise1_rdd_creation | 3 minutes ago |
| 📁 exercise2_rdd_transfor... | 3 minutes ago |
| 📁 exercise3_rdd_sets | 3 minutes ago |
| 📁 exercise4_rdd_key_valu... | 3 minutes ago |
| 📁 project2 | 3 minutes ago |
| 📄 kddcup.data_1_percent.... | 3 minutes ago |

Enter the first folder (exercise1_rdd_creation) and open the Jupyter Notebook file named **exercise1_rdd_creation.ipynb**.

Follow the instructions in the Jupyter Notebook to learn how to program with PySpark. You can use "SHIFT+ENTER" to execute the cell and advance to the next cell. Be sure to finish the code in the cells that contain "# TODO".

In this exercise, you only need to finish the jupyter notebooks in the first three folders in the repository (**exercise1_rdd_creation**, **exercise2_rdd_transformation**, and **exercise3_rdd_sets**). The rest of the exercises are optional.

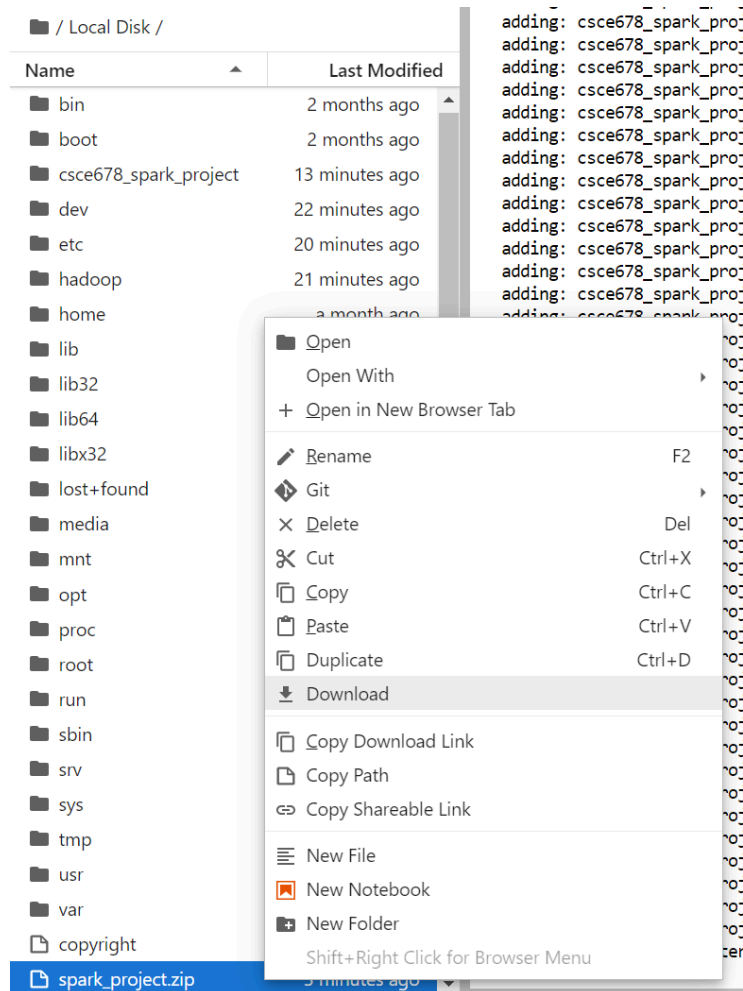## Shut Down the Dataproc Cluster (if necessary)

When you are not working on the Jupyter Notebooks, it is recommended to shut down (NOT deleting) the cluster to prevent unnecessary charges. Open the Dataproc Console and select the cluster you created for Spark. And then click ■ STOP to stop the cluster.

# Download and Submit the Jupyter Notebooks

Once you have completed the Jupyter Notebook, switch back to the terminal and run the following command:

```
zip -r spark_project.zip csce678_spark_project/
```

In the file browser, switch to / Local Disk / and find "spark_project.zip". Right click the file to choose "Download".



Once You have downloaded spark_project.zip, head over to Canvas and upload the file to the assignment "Exercise Spark Jupyter Notebooks (Submission)".

# Delete the Dataproc Cluster and Google Cloud Storage

Finally, you can delete the dataproc cluster and the bucket you created in the Google Cloud Storage.