

Deploy Spark Cluster on Google Cloud Dataproc:

Cloud Dataproc is a fast, easy-to-use, fully-managed cloud service for running Apache Spark and Apache Hadoop clusters in a straightforward and cost-efficient way.

Create a bucket (if you haven't already done so) – The following steps need to be taken **only once**:

1. Login to your Google Cloud account (<https://console.cloud.google.com>) and go to Cloud Storage.
2. Create a new bucket. Call it qst843-xx (xx you initial)
3. Select **multi-region**, and for location, select **us**, leave the rest as default, and click **Create**.
4. Run the following query to test if everything is working fine from Cloud Shell.
 - a. `gsutil ls`
 - b. This should list all the buckets in your project, including the one you just made.

Create a cluster – Login to your Google Cloud (<https://console.cloud.google.com>) and open a Google Cloud Shell. Make sure you are in the intended project.

1. For added security, clusters are set up without an external IP address, which results in a lack of internet access. We can use Cloud NAT service to allow instances that don't have external IP addresses to connect to the internet securely. Create Cloud Router instances in *us-central1* region. Enable Private Google Access on the default Subnetwork and internet access to the VM – **Add Cloud NAT** for your VPC/subnet/region (default/default/us-central1):

```
gcloud compute networks subnets update default \
  --region us-central1 \
  --enable-private-ip-google-access

gcloud compute routers create nat-router-us-central1 \
  --network="default" --region="us-central1"

gcloud compute routers nats create nat-config \
  --router=nat-router-us-central1 \
  --router-region="us-central1" \
  --nat-custom-subnet-ip-ranges="default" \
  --auto-allocate-nat-external-ips
```

We should now have internet access on the master node, which will enable us to perform tasks such as software updates or git clones. This step needs to be done only once per project/region.

2. The following step must be performed **to create a new cluster**. Run the following script in Cloud Shell to launch a new cluster named *single-spark353*. Copy this command to an editor of your choice and make sure to replace <BUCKET-NAME> with the bucket you created in the previous section:

```
gcloud dataproc clusters create single-spark353 \
  --bucket <BUCKET-NAME> \
  --region us-central1 \
  --master-machine-type n2-highmem-4 \
```

```
--master-boot-disk-type pd-standard \
--master-boot-disk-size 30 \
--single-node \
--no-address \
--image-version 2.2-debian12 \
--optional-components JUPYTER \
--initialization-actions gs://dataproc-initialization-actions/python/pip-install.sh \
--metadata 'PIP_PACKAGES=google-cloud-storage' \
--enable-component-gateway
```

(This script is also accessible in 05-A-Tour-of-Spark/Deploy-Dataproc-Cluster-Single-Node.txt)

- This will give you a single-node cluster with 50 GB HDD, 4 CPU cores, and 32GB Memory (n2-highmem-4). This cluster will cost about \$0.26/hour.
3. Confirm that your cluster is up and running from the Dataproc page (located under Big Data under the Google Cloud menu).

Connecting to the Cluster

- From the Dataproc tab, select the cluster, and under web interfaces, select JupyterLab.
- OR - From the Cloud SDK shell, execute the following command to establish a secure SSH tunnel:

```
gcloud compute ssh --zone us-central1-X single-spark353-m -- -L 2222:localhost:8123 -L 8088:localhost:8088
```

Replace X with the region dedicated to the node.

Note: You may occasionally get disconnected. If that happens, repeat the step above and refresh your browser. Your work will be saved in your cluster and bucket under the *notebook* folder.

Stopping your Cluster

To ensure you don't run out of money, stop your instance once you are not using it. You can do this by going to the Dataproc page and stopping the cluster.

Clean up

To make sure we won't be charged for any of the resources we are not using, delete the cluster after each use:

- Select the cluster From the Dataproc page and click the DELETE button.
 - Alternatively, you can also use the following command from a **new** Cloud SDK terminal (please check from the UI to make sure the cluster is being deleted):

```
gcloud dataproc clusters delete single-spark353 --region us-central1
```

* Notice that even after deleting the cluster, your notebooks will persist in the bucket, and when you create a new cluster that points to the same bucket, you can reuse those notebooks.

Deploying a Large Cluster

To create a large cluster (**CAUTION**) with multiple nodes, use the following command:

```
gcloud dataproc clusters create multi-spark353 \  
  --bucket <BUCKET-NAME> \  
  --region us-central1 \  
  --master-machine-type n2-highmem-2 \  
  --master-boot-disk-type pd-standard --master-boot-disk-size 30 \  
  --num-workers 2 \  
  --worker-machine-type n2-highmem-2 \  
  --worker-boot-disk-type pd-standard --worker-boot-disk-size 30 \  
  --secondary-worker-type spot \  
  --num-secondary-workers 0 \  
  --secondary-worker-boot-disk-size 30 --num-secondary-worker-local-ssds 0 \  
  --image-version 2.2-debian12 \  
  --no-address \  
  --optional-components JUPYTER \  
  --initialization-actions gs://dataproc-initialization-actions/python/pip-install.sh \  
  --metadata 'PIP_PACKAGES=google-cloud-storage' \  
  --enable-component-gateway
```

(This script is also accessible in 05-A-Tour-of-Spark/Deploy-Dataproc-Cluster-Multi-Node.txt)

Important note: Please be advised of its cost. Cost can be calculated using Dataproc cost calculator (make sure to include GCE is selected): <https://cloud.google.com/products/calculator/>

The cluster above has 1 master node, 2 primary workers, and 4 secondary workers (12 CPUs & 96 GB memory). Two of these worker nodes are spot/preemptible instances that can lower our cost significantly. A cluster with preemptible instances cannot be stopped, so it's not ideal for our use-case. Preemptible instances are offered with a big discount (~70%) but will not last for more than 24 hours. The following line includes 2 preemptible workers:

```
--num-secondary-workers 2
```

This cluster will cost about \$0.90/hour.

Note about preemptible worker pool: One can easily add/remove preemptible workers from the UI (or in the command line). Go to the cluster of interest in Dataproc (make sure the cluster is running) > Configuration > Edit > add/remove secondary workers > save. After ~30 seconds your cluster will reflect the change.

Connecting to the Cluster

Same as the single node cluster mentioned above.



Clean up

```
gcloud dataproc clusters delete single-spark353 --region us-central1
```