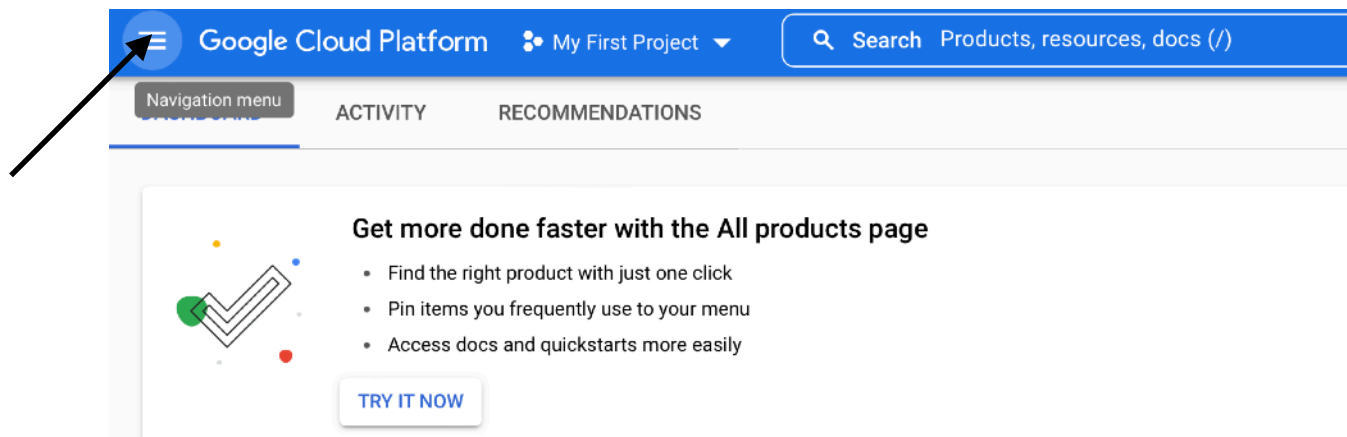


Big data processing on Google Cloud

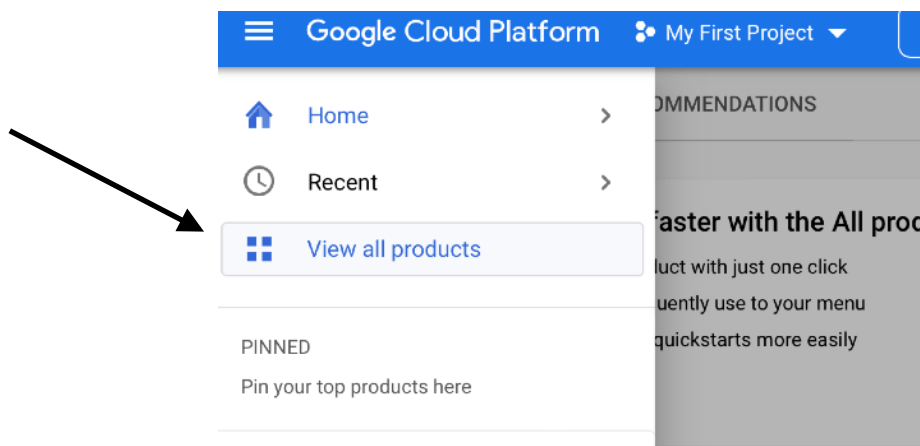
Creating and launching an Hadoop cluster, 1 Master and 2 Slaves

Please follow the steps illustrated by following images

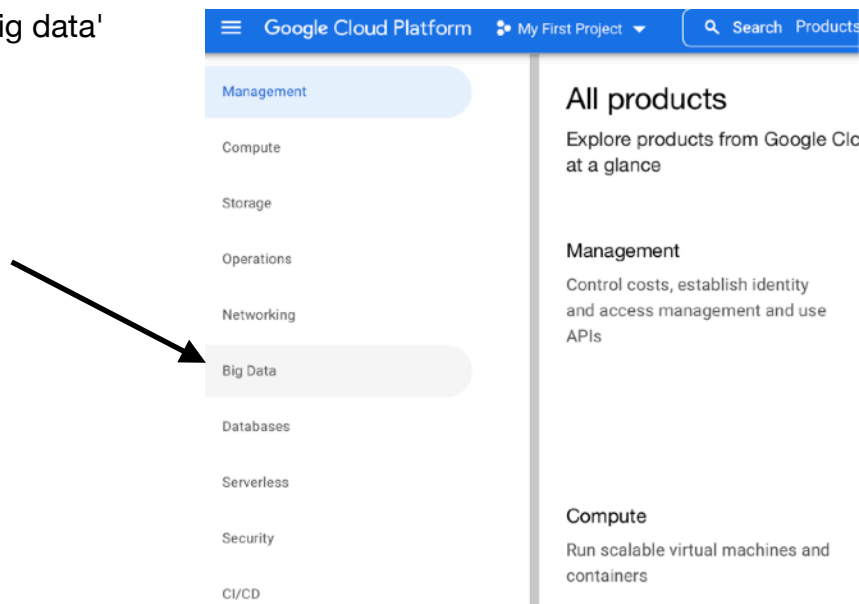
Starting from the initial home page (after having logged in) click on the Navigation menu



Then click on 'View all products'

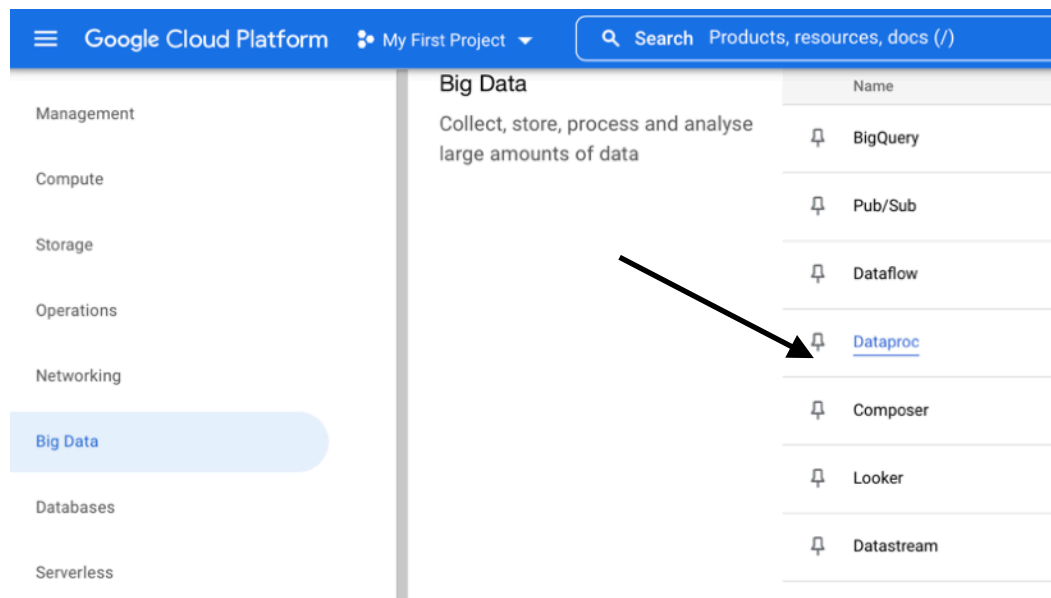


Now click on 'Big data'

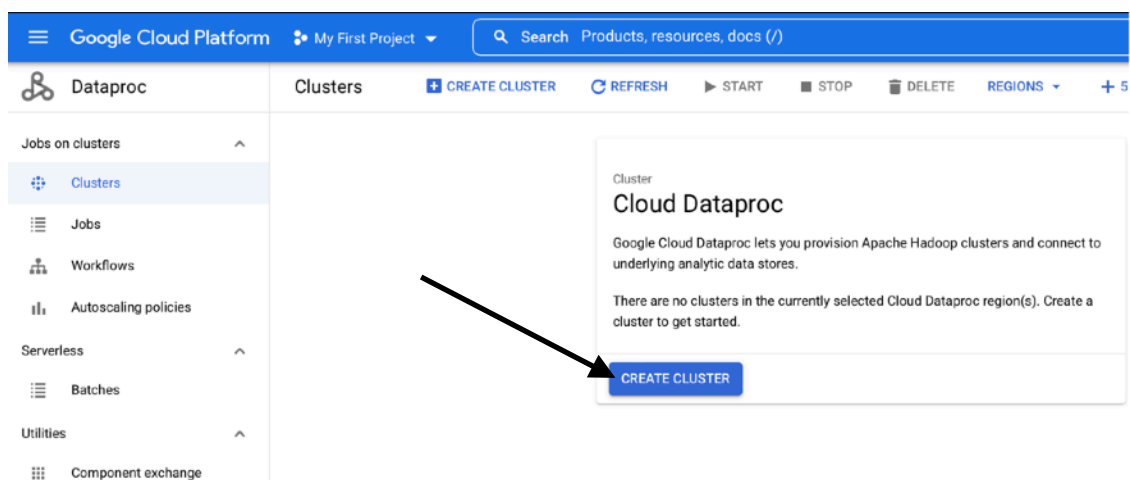


Big data processing on Google Cloud

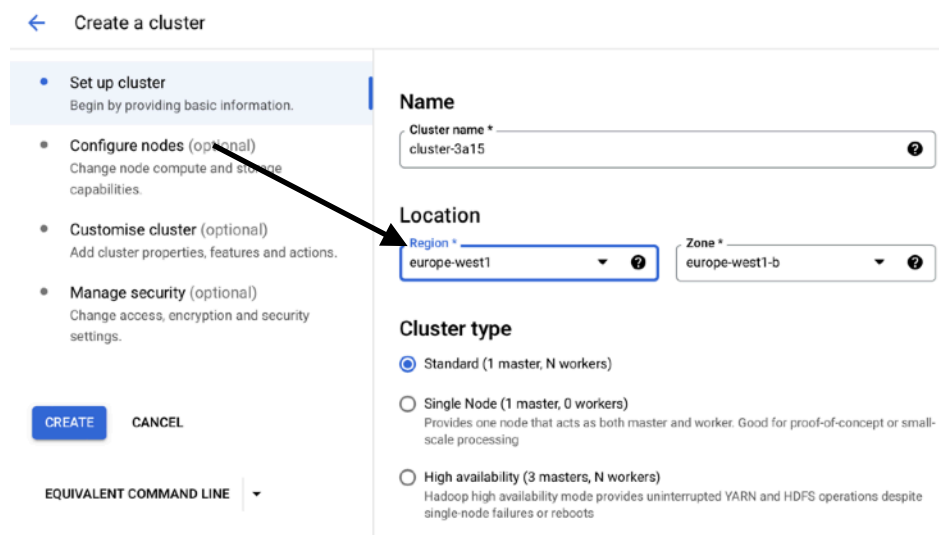
Then opt for Dataproc (which will allow for the creation of an Hadoop cluster).



For cluster creation click on CREATE CLUSTER (opt for compute engine)



Now fill the regions like indicated below and leave the rest as pre-filled



Big data processing on Google Cloud

Concerning 'Configure node' you can observe that as a default setting you have 2 slaves (worker nodes). You can keep all these default settings.

Set up cluster
Begin by providing basic information.

Configure nodes (optional)
Change node compute and storage capabilities.

Customise cluster (optional)
Add cluster properties, features and actions.

Manage security (optional)
Change access, encryption and security settings.

CREATE CANCEL

EQUIVALENT COMMAND LINE

Worker nodes

Each contains a YARN NodeManager and a HDFS DataNode. HDFS replication factor is 2.

Machine family

GENERAL-PURPOSE COMPUTE-OPTIMISED MEMORY-OPTIMISED

Machine types for common workloads, optimised for cost and flexibility

Series
N1

Powered by Intel Skylake CPU platform or one of its predecessors

Machine type
n1-standard-4 (4 vCPU, 15 GB memory)

	vCPU	Memory
	4	15 GB

✓ CPU PLATFORM AND GPU

Number of worker nodes *
2

You can now CREATE the cluster (do not change other parameters)

Click now on the cluster link

Dataproc

Clusters **CREATE CLUSTER** **REFRESH** **START** **STOP** **DELETE** **REGIONS**

Filter Search clusters, press Enter

Name	Status	Region	Zone	Total worker nodes	Scheduled deletion
cluster-3a15	Running	europe-west1	europe-west1-b	2	Off

Jobs on clusters

- Clusters
- Jobs
- Workflows

Now click on the SSH button, this will enable an SSH connection with the Master

MONITORING JOBS **VM INSTANCES** CONFIGURATION WEB INTERFACES

Filter Filter instances

Name	Role	SSH
cluster-3a15-m	Master	SSH
cluster-3a15-w-0	Worker	
cluster-3a15-w-1	Worker	

You have now the possibility of using HDFS, MapReduce, Spark and Hive.

For Spark in Python just use pyspark

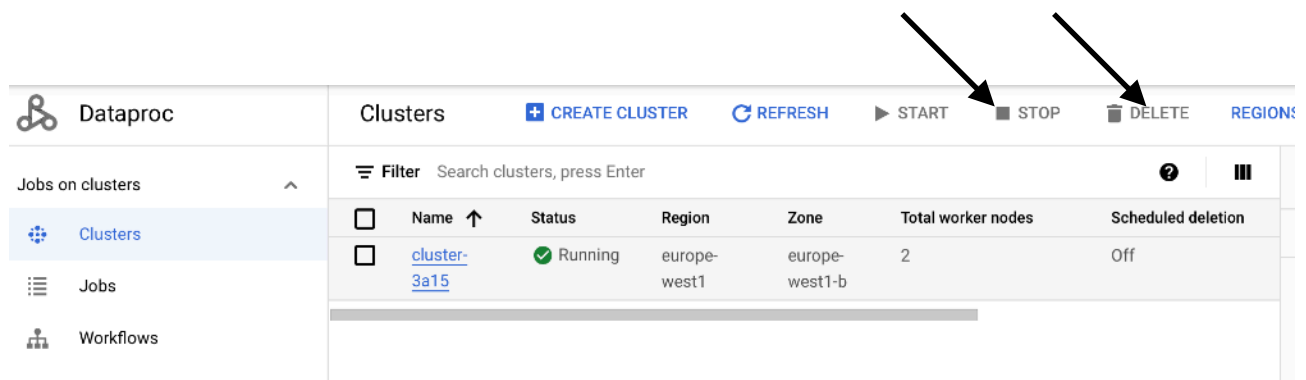
Concerning HDFS, with

you get information about all living DataNodes (including their addresses).

IMPORTANT: once finished with your session, use the exit command for exiting the SSH terminal.

-->Then go to the Google Cloud interface for the cluster VMs and stop the cluster

IMPORTANT: This will stop the cluster by stopping money consumption, just click on start to start again when needed. Note that there will be some minimal money consumption, so if you think that for several days/weeks you will not use the cluster, just DELETE it and re-CREATE when needed.



The screenshot displays the Google Cloud Dataproc Clusters management interface. On the left, a sidebar shows navigation options: Clusters (selected), Jobs, and Workflows. The main panel shows a table of clusters. The table has columns for Name, Status, Region, Zone, Total worker nodes, and Scheduled deletion. A single cluster is listed with the name 'cluster-3a15', status 'Running', region 'europe-west1', zone 'europe-west1-b', 2 worker nodes, and no scheduled deletion. Above the table, there are buttons for 'CREATE CLUSTER', 'REFRESH', 'START', 'STOP', and 'DELETE'. Two arrows point to the 'STOP' and 'DELETE' buttons.

Name	Status	Region	Zone	Total worker nodes	Scheduled deletion
cluster-3a15	Running	europe-west1	europe-west1-b	2	Off

Big data processing on Google Cloud

Run WordCount on the Google Cloud Cluster.

Once the SSH connection/terminal available (see previous pictures), proceed with the following steps

1. Download the mapper, reducer and a text file in you home directory, with the following commands

```
> wget https://www.dropbox.com/s/471k0286292ifho/mapper.py
```

```
> wget https://www.dropbox.com/s/e8s6f7rsiwts84m/reducer.py
```

Concerning the input text file:

```
> wget https://www.dropbox.com/s/hj8khqc94vsrzw9/shake.txt
```

2. Give the execution right to the .py files.

```
chmod +x *.py
```

3. Now create a wc directory in the HDFS /user directory.

```
hdfs dfs -mkdir /user/wc
```

4. Now you need to create a 'input' subdirectory containing the input file for job

```
hdfs dfs -mkdir /user/wc/input
```

5. Finally, you need to transfer the input file from the local file system to the HDFS file system so that MapReduce can process it.

```
hdfs dfs -put shake.txt /user/wc/input/.
```

In order to check the status of shake.txt inside HDFS use

```
hdfs fsck /user/wc/input/shake.txt -files -blocks -locations
```

Big data processing on Google Cloud

In particular this will show number of blocks and respective DataNodes storing them.

6. You are ready now to launch our MapReduce job. Copy this command into the SSH cluster terminal. You will get several statistics upon job completion. The job results are in the 'output' HDFS directory. Each reduce task has produced its own file result.

```
hadoop jar /usr/lib/hadoop/hadoop-streaming-3.2.2.jar \  
-input /user/wc/input \  
-output /user/wc/output \  
-file ~/mapper.py \  
-mapper ~/mapper.py \  
-file ~/reducer.py \  
-reducer ~/reducer.py
```

Note that we are assuming that Python programs are in your home (~) directory in the local file system (not in HDFS), otherwise you have to change this path.

To see the output files (one per reduce task):

```
hdfs dfs -ls /user/wc/output
```

To see the content of these files:

```
hdfs dfs -cat /user/wc/output/*
```

If you want to use the combiner then use this variant (note that you re-use the reducer as a combiner, and that this holds for the word-count problem, but not necessarily for other problems, e.g. computing averages). Note that you have to use another output file.

```
hadoop jar /usr/lib/hadoop/hadoop-streaming-3.2.2.jar \  
-input /user/wc/input \  
-output /user/wc/output1 \  
-file ~/mapper.py \  
-mapper ~/mapper.py \  
-file ~/reducer.py \  
-reducer ~/reducer.py \  
-combiner ~/reducer.py
```

Of course if your combiner is not the reducer, then you have to specify both -file and -combiner parameters with the same path for the .py file including the combiner.

This version is to set a particular number of reduce tasks (e.g., 3 reduce tasks)

```
hadoop jar /usr/lib/hadoop/hadoop-streaming-3.2.2.jar \  
-input /user/wc/input \  
-output /user/wc/output2 \  
-file ~/mapper.py \  
-mapper ~/mapper.py \  
-file ~/reducer.py \  
-reducer ~/reducer.py \  
-jobconf mapred.reduce.tasks=3
```

Big data processing on Google Cloud

7. You can use the `-ls` and `-cat` HDFS command to, respectively, list the content of the output directory, and display one of the output file.