# Calculating Pi with PySpark

https://hc.labnet.sfbu.edu/~henry/npu/classes/learning_spark/key_value_pair/slide/exercise_key_value_pair.html
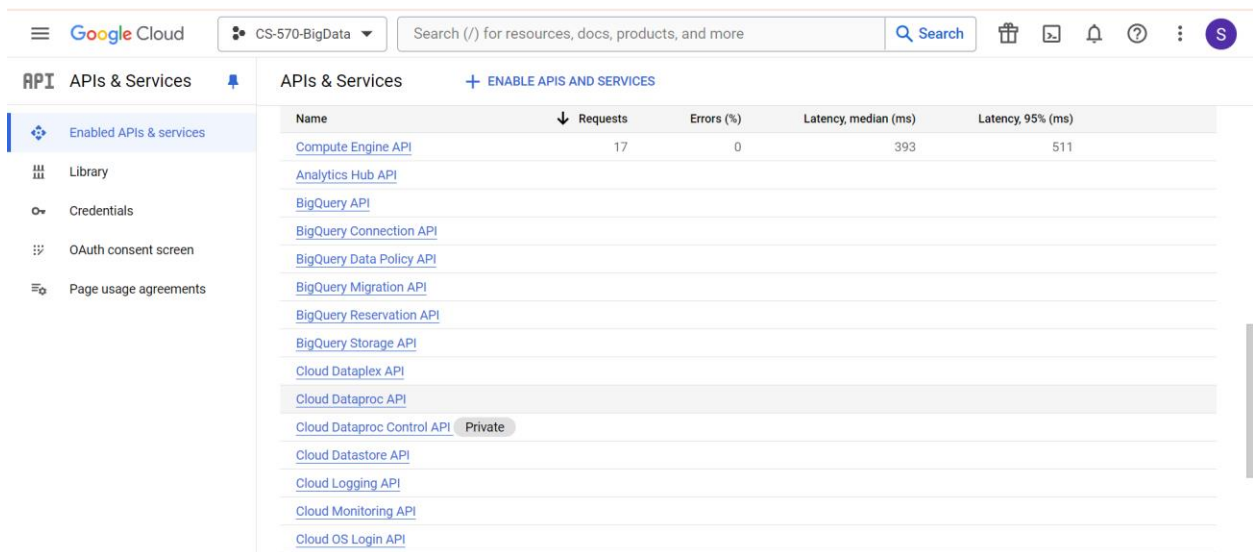
Q24 ==> Project: Calculating Pi

## Note:

- If you have an existing cluster go to step 3
- If you have an existing cluster and bucket ready go to step 4

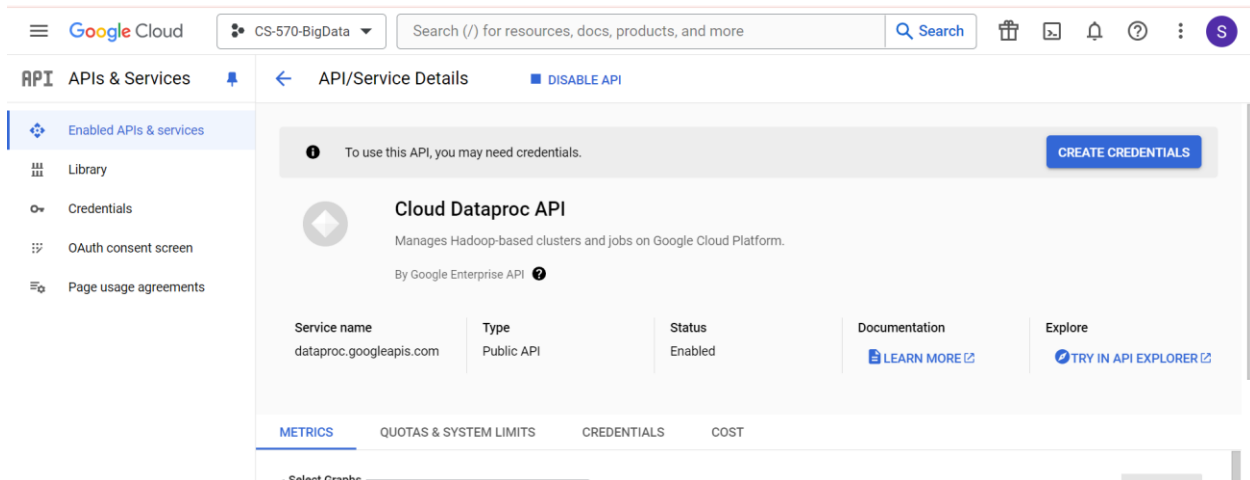## Step 1: Go to your existing vm-project

- Go to navigation menu then API and Services



- Enable the Cloud Dataproc API

## Step 2: Create Cluster

- Go to Dataproc on navigation menu
- Create a Dataproc cluster

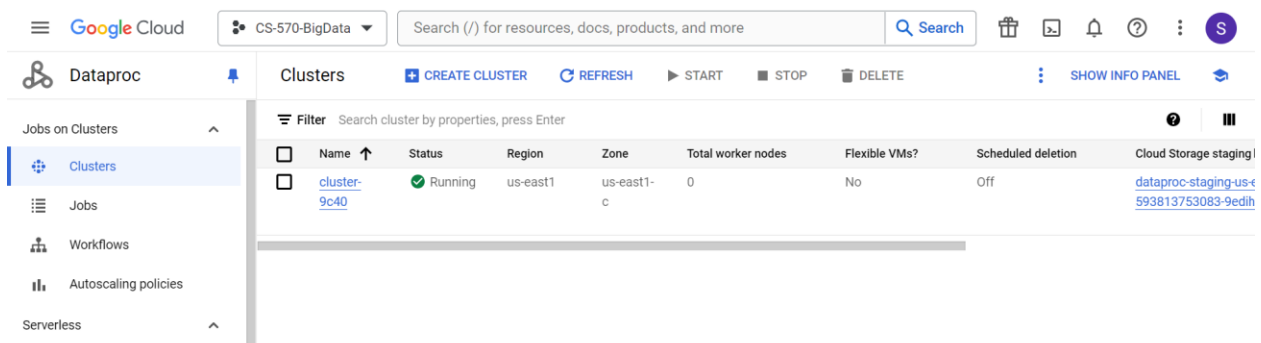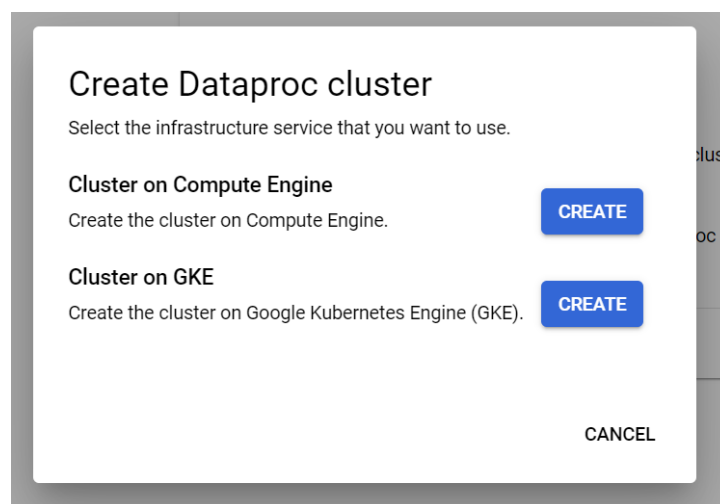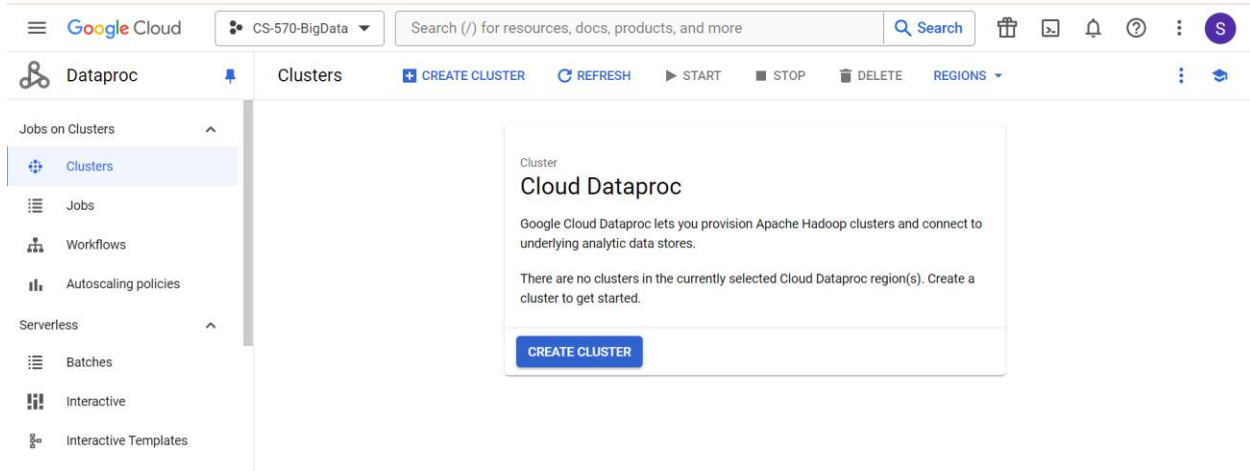Now our cluster is created we proceed with logging in our ssh-browser the Dataproc master server.

Click on the cluster then go to VM INSTANCES

Click on the ssh of the cluster to start working on the ssh-browser



# Step 3: Create Bucket

Create bucket to store input and output files

Configure the bucket name, location and storage



# Step 4: Create Pi_calculate code file

Now go to the cloud shell and create the pi_calculate.py file.

Start your Free Trial with $300 in credit. Don't worry—you won't be charged if you run out of credits. Learn more ⮺

DISMISS    START FREE

≡  Google Cloud    ⁍ CS-570-BigData ▼    Search (/) for resources, docs, products, and more    🔍 Search    ⌷ 🔔 ⑦ ⋮ S

⚴ Dataproc 📌    Clusters    ⊞ CREATE CLUSTER    ⟳ REFRESH    ▶ START    ■ STOP    🗑 DELETE    ⋮    SHOW INFO PANEL ◈

Jobs on Clusters    ⌃    ⇶ Filter  Search cluster by properties, press Enter    ❷  ▥

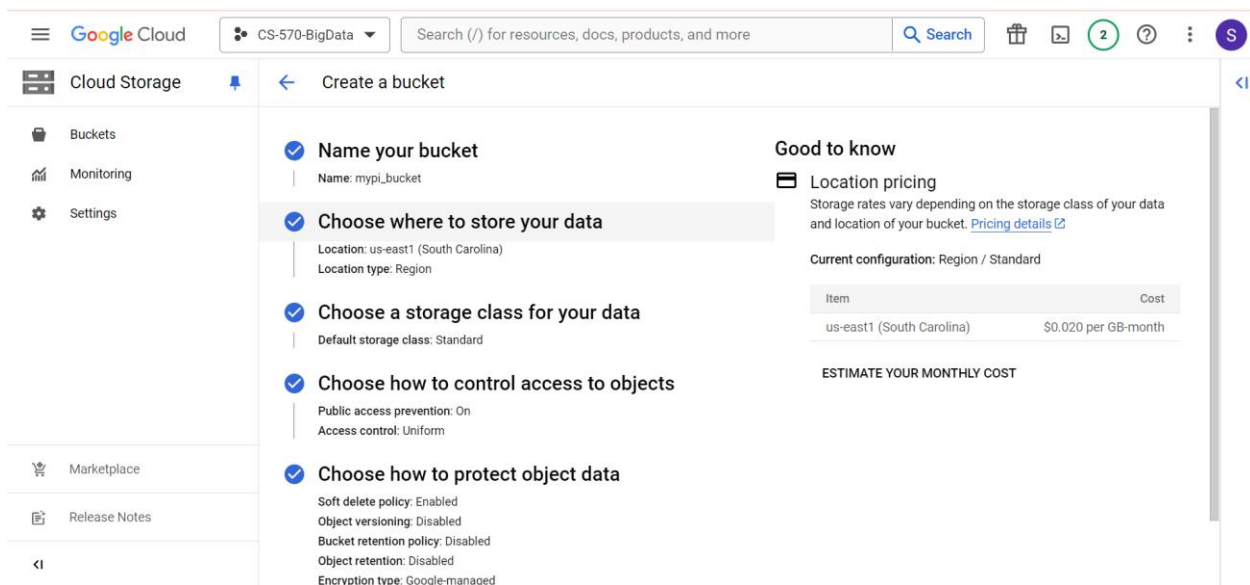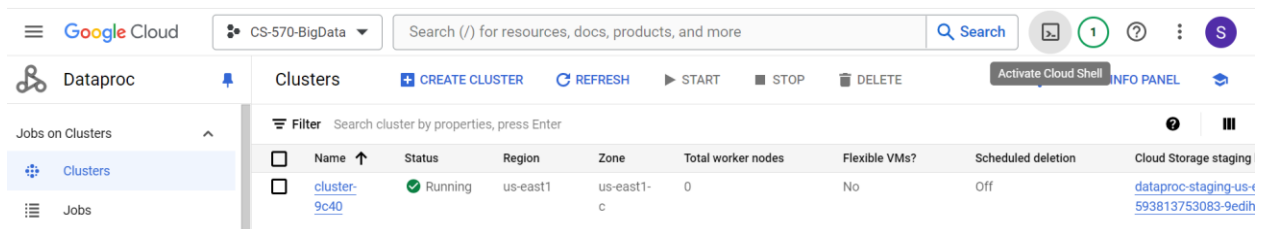| | Name ↑ | Status | Region | Zone ↑ | Total worker nodes | Flexible VMs? | Scheduled deletion | Cloud Storage staging |
|---|---|---|---|---|---|---|---|---|
📑 Release Notes | ☐ | cluster-9c40 | ✅ Running | us-east1 | us-east1-c | 0 | No | Off | dataproc-staging-us-e 593813753083-9edih

◁|

CLOUD SHELL
≥ Terminal    (cs-570-bigdata) ✕    + ▾         ✎ Open Editor    ⌨ ⚙ ⊡ 🖳 ⋮  _  ⇕  ⬀  ✕

```
Welcome to Cloud Shell! Type "help" to get started.
Your Cloud Platform project in this session is set to cs-570-bigdata.
Use "gcloud config set project [PROJECT_ID]" to change to a different project.
syohanne998@cloudshell:~ (cs-570-bigdata)$ ▮
```

Go to open Editor and write the following code
............................................................................................................................
.

```python
import argparse
import logging
from operator import add
from random import random
from pyspark.sql import SparkSession

logger = logging.getLogger(__name__)
logging.basicConfig(level=logging.INFO, format='%(levelname)s:
%(message)s')

def calculate_pi(partitions, output_uri):
    def calculate_hit(_):
        x = random() * 2 - 1
        y = random() * 2 - 1
        return 1 if x ** 2 + y ** 2 < 1 else 0
    tries = 100000 * partitions
    logger.info(
        "Calculating pi with a total of %s tries in %s partitions.",
tries, partitions)
    with SparkSession.builder.appName("My PyPi").getOrCreate() as spark:
        hits = spark.sparkContext.parallelize(range(tries), partitions)\
            .map(calculate_hit)\
```

```python
            .reduce(add)
        pi = 4.0 * hits / tries
        logger.info("%s tries and %s hits gives pi estimate of %s.",
tries, hits, pi)
        if output_uri is not None:
            df = spark.createDataFrame(
                [(tries, hits, pi)], ['tries', 'hits', 'pi'])
            df.write.mode('overwrite').json(output_uri)
if __name__ == "__main__":
    parser = argparse.ArgumentParser()
    parsers.add_argument(
        '--partitions', default=2, type=int,
        help="The number of parallel partitions to use when calculating
pi.")
    parsers.add_argument(
        '--output_uri', help="The URI where output is saved, typically an
S3 bucket.")
    args = parsers.parse_args()

    calculate_pi(args.partitions, args.output_uri)
```

........................................................................

.

▶ Open Terminal | 👁 | ▢ | ⋮ | ⤢ | ⟗ | ✕

🐍 calculate_pi.py ✕                                                   ▷ ∨ ▢ ⋯

🐍 calculate_pi.py > ⊙ calculate_pi

```python
1   import argparse
2   import logging
3   from operator import add
4   from random import random
5   from pyspark.sql import SparkSession
6
7
8   logger = logging.getLogger(__name__)
9   logging.basicConfig(level=logging.INFO, format='%(levelname)s: %(message)s')
10
11
12  def calculate_pi(partitions, output_uri):
13
14      def calculate_hit(_):
15          x = random() * 2 - 1
16          y = random() * 2 - 1
17          return 1 if x ** 2 + y ** 2 < 1 else 0
18
19      tries = 100000 * partitions
20
21      logger.info(
22          "Calculating pi with a total of %s tries in %s partitions.", tries, partitions)
```

# Step 5: Run the code

Now click on the cluster and go to vm instances
Then click on the SHH to open SHH-browser

gcloud dataproc jobs submit pyspark gs://mypi_bucket/input/calculate_pi.py \
           --cluster cluster-9c40 \
           --region us-east1 \
           -- --partitions 4 --output_uri gs://mypi_bucket/pi-calc-output



Now go to job details and turn the line warp on

Viewing the output of pi value

```
syohanne998@cluster-9c40-m:~$ gsutil cat gs://mypi-bucket/pi-calc-output/*
{"tries":400000,"hits":314428,"pi":3.14428}
syohanne998@cluster-9c40-m:~$
```

# Step 6: Close or delete the cluster and the bucket
Remove bucket and close cluster
I have no use of the bucket and I have stopped my cluster instead of deleting it but you can also delete the cluster if you won't use.