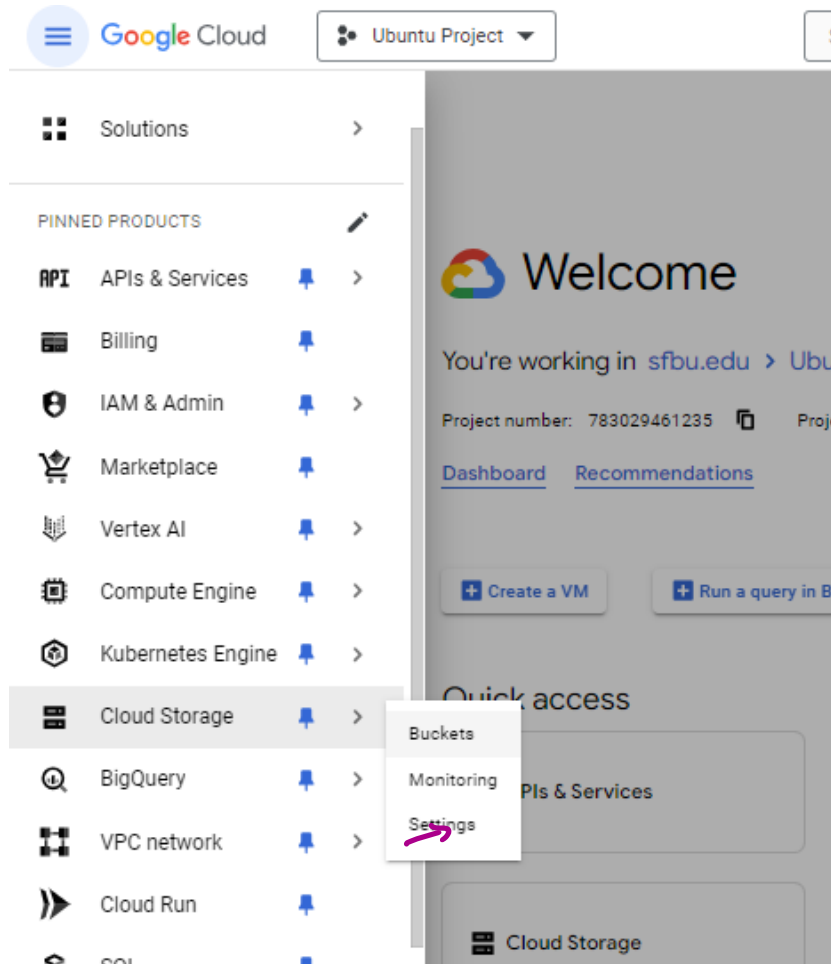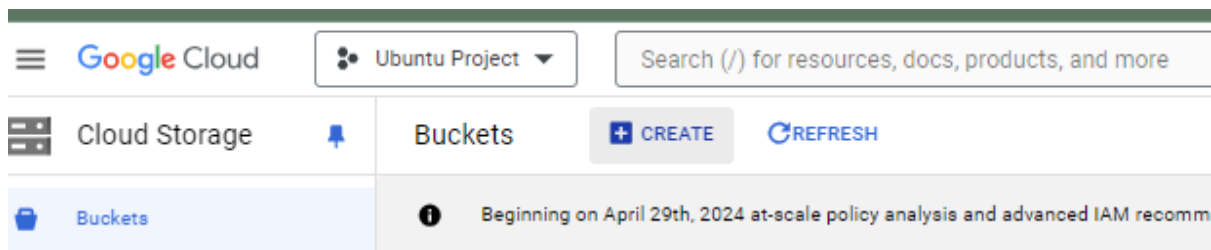Name: Sened Desalegn

# Pyspark on GCP

## *Create a Bucket*

1. To create a bucket in Google Cloud Storage, navigate to the Cloud Console, select or create a project, then go to Storage > Browser, click "Create Bucket", name it uniquely, choose a location, and click "Create".



2. Click Create bucket.

3. On the Create a Bucket page, enter your bucket information. To go to the next step, click Continue.



4. For Name your bucket, enter a name that meets the bucket name requirements. Choose the cheaper region. And rest things can be default values.

After you click "Create", you might face a pop up message like this:
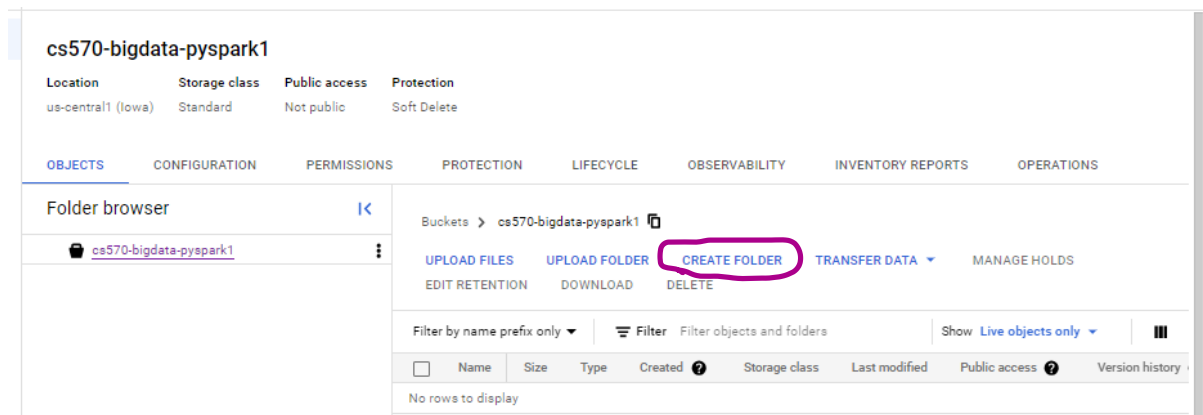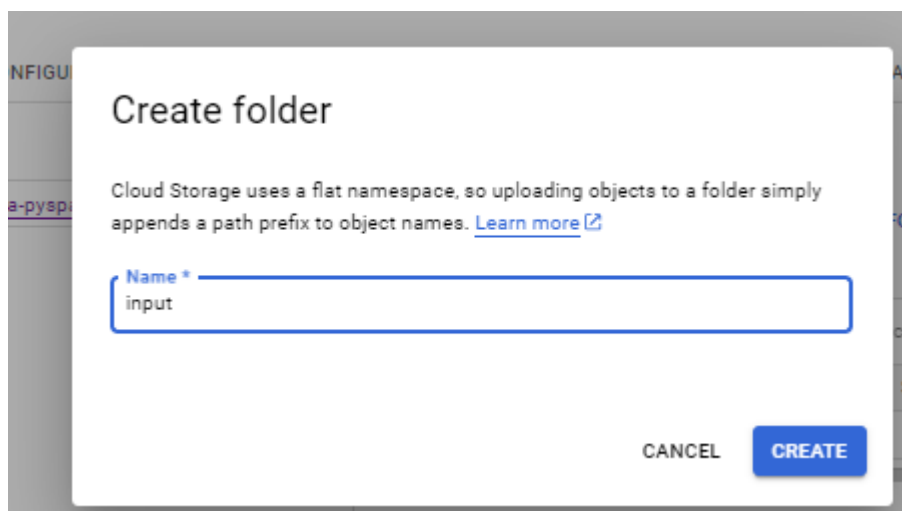
Just click "Confirm".



5. Now, we create a folder in the bucket and add sample.txt, you can follow these steps:

- ✓ In the Google Cloud console, go to the Cloud Storage Buckets page.
- ✓ Click on the name of the bucket you created

Name: Sened Desalegn


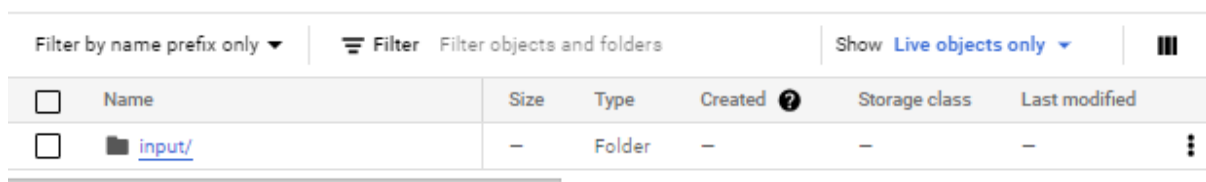
✓ Click on the Create folder button.
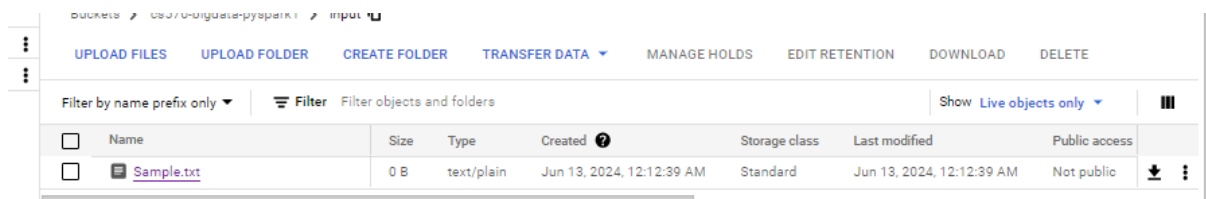✓ Enter a name for your folder and click Create.



6. Click on the name of the folder you just created.



7. Click on the Upload files button.

After I have uploaded the sample.txt file I have created.



8. Select your sample.txt file and click Open. The sample file can be any text file. I created a sample file containing information about Generative AI.

Name: Sened Desalegn

If you click on the authenticated URL, you will get what is inside the sample.txt file.



## Create a Dataproc cluster:

To create a Dataproc cluster, follow these steps:

1. Navigate to the APIs & Services page in the Google Cloud Console.

2. Click on "Enable APIs and Services"

Name: Sened Desalegn



3. Search for "Cloud Dataproc API" and select it.

4. Click Enable.



To create a Dataproc cluster in Google Cloud Platform, you can use the Google Cloud Console. Here are the steps:

1. In the Google Cloud console, go to the Dataproc Clusters page.

2. Click Create Cluster.



3. In the Create Dataproc cluster dialog, click Create in the Cluster on Compute engine row.



4. In the Cluster Name field, enter a name for your cluster.

5. In the Region and Zone lists, select the same region as the bucket and zone. Choose the cluster type as a single node. And the rest can be the default.

Name: Sened Desalegn

6. For all the other options, use the default settings. 7. To create the cluster, click Create.



Note: If you get error for subnetting issues, first you have to navigate to VPC networks, and under the VPC network, there is default and under default you have to select the subnet based on the zone you have and edit the network.

## To create and save the spark job python file:

1. In the top right corner of the console, click the Activate Cloud Shell button.



2. Once the Cloud Shell is activated, click on the Open Editor button in the top right corner of the Cloud Shell window.



3. Click on the new file icon beside your username to create a new file

Here is the python code used:

import pyspark

import sys

Name: Sened Desalegn

if len(sys.argv) != 3:

raise Exception("Exactly 2 arguments are required: <inputUri> <outputUri>")

inputUri = sys.argv[1]

outputUri = sys.argv[2]

sc = pyspark.SparkContext()

lines = sc.textFile(inputUri)

words = lines.flatMap(lambda line: line.split())

wordCounts = words.map(lambda word: (word, 1)).reduceByKey(lambda count1,

count2: count1 + count2)

wordCounts.saveAsTextFile(outputUri)

```
EXPLORER                    ...    words.py 1  ×
 SDESALEG137                        words.py > ...
    kubectl                      1   import pyspark
  ! kubia-rc.yaml                2   import sys
    README-cloudshell.txt        3   if len(sys.argv) != 3:
    words.py          1          4   raise Exception("Exactly 2 arguments are required: <inputUri> <outputUri>")
                                 5   inputUri = sys.argv[1]
                                 6   outputUri = sys.argv[2]
                                 7   sc = pyspark.SparkContext()
                                 8   lines = sc.textFile(inputUri)
                                 9   words = lines.flatMap(lambda line: line.split())
                               ●10   wordCounts = words.map(lambda word: (word, 1)).reduceByKey(lambda count1,
                                11   count2: count1 + count2)
                                12   wordCounts.saveAsTextFile(outputUri)
```

Save the file. I saved it as words.py and close the shell.

## Running PySpark Job on Google Cloud Dataproc and Sorting Output Files:

Step 1: Click on "Activate Cloud Shell" like we did previously. Authenticate with Google Cloud Platform (GCP). If you encounter an authentication error, run the command:

    $ gcloud auth login

This command will open a web page where you can authenticate and obtain new credentials. Follow the on-screen instructions to complete the authentication process

Click on the click, authorize, and copy the code. Paste the same code in the authentication code.

Name: Sened Desalegn



Step 2: Submit the PySpark job to Dataproc

This command submits a PySpark job (words.py) to the Dataproc cluster named cluster-3b08, which is located in the us-central1 region. The job processes an input file (sample.txt) from a specified GCS bucket and writes the output to another GCS directory.

Change the cluster name and Cluster Region, by going back to your cluster.



$ gcloud dataproc jobs submit pyspark words.py \

--cluster=cluster-3b08 \

--region=us-central1 \

-- gs://cs570-bigdata-pyspark/input/sample.txt gs://cs570-bigdata-pyspark/output



To get the location, open the bucket we created earlier, and go to the sample.txt. (My bucket name is cs570-bigdata-pyspark)

Name: Sened Desalegn



Step 3: List the files in the output directory

```
sdesaleg137@cloudshell:~$ gsutil ls gs://cs570-bigdata-pyspark1/output/
gs://cs570-bigdata-pyspark1/output/
gs://cs570-bigdata-pyspark1/output/_SUCCESS
gs://cs570-bigdata-pyspark1/output/part-00000
gs://cs570-bigdata-pyspark1/output/part-00001
sdesaleg137@cloudshell:~$
```

This command lists the files present in the gs://cs570-bigdata-pyspark1/output/ GCS bucket directory.

Step 4: Copy the output files to the current directory

```
sdesaleg137@cloudshell:~$ gsutil cp gs://cs570-bigdata-pyspark1/output/* .
Copying gs://cs570-bigdata-pyspark1/output/_SUCCESS...
Copying gs://cs570-bigdata-pyspark1/output/part-00000...
Copying gs://cs570-bigdata-pyspark1/output/part-00001...
/ [3 files][  1.7 KiB/  1.7 KiB]
Operation completed over 3 objects/1.7 KiB.
sdesaleg137@cloudshell:~$
```

This command copies all the files from the gs://cs570-bigdata-pyspark1/output/ GCS bucket directory to the current directory in the Cloud Shell environment.

Step 5: Combine the contents of part-00000 and part-00001 files

        $ cat part-00001 >> part-00000

```
sdesaleg137@cloudshell:~$ cat part-00001 >> part-00000
sdesaleg137@cloudshell:~$
```

This command appends the contents of the part-00001 file to the end of the part-00000 file.

Step 6: Sort the combined file based on the second column

This command reads the contents of the part-00000 file and pipes (|) it as input to the sort command. The sort command then sorts the input based on the second column (-k 2). The sorted output will be displayed in the terminal for easy viewing.

Name: Sened Desalegn

```
sdesaleg137@cloudshell:~$ cat part-00000 | sort -k 2
('a', 1)
('A', 1)
('across', 1)
('adapt', 1)
('advanced', 1)
('amounts', 1)
('and/or', 1)
('are', 1)
('attempted', 1)
('automating', 1)
('automation.', 1)
('be', 1)
('broader', 1)
('by', 1)
('can', 1)
('capable', 1)
('capacity', 1)
('comparing', 1)
('complex', 1)
('considered', 1)
('content,', 1)
('creating', 1)
('creativity', 1)
('data', 1)
('data,', 1)
('decision-making.', 1)
('designed', 1)
```
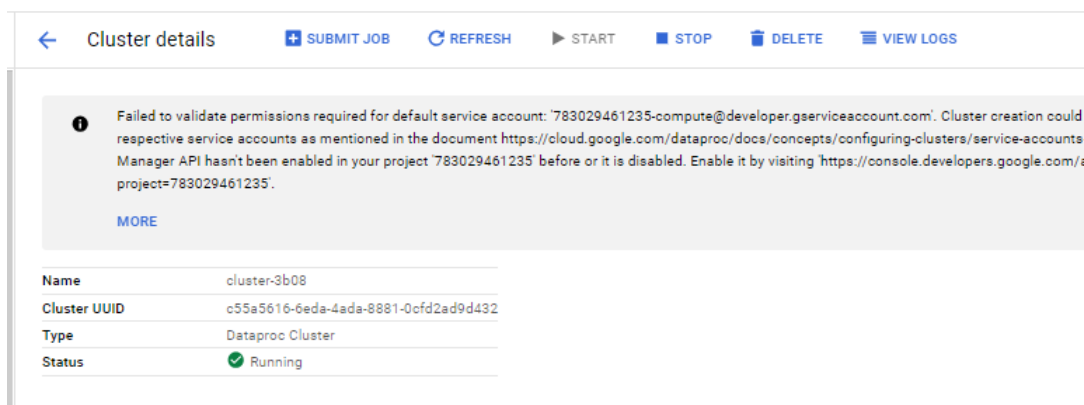
Name: Sened Desalegn

```
('traditionally', 1)
('Unlike', 1)
('was', 1)
('what', 1)
('while', 1)
('white-collar', 1)
('with', 1)
('automation,', 2)
('cognitive', 2)
('in', 2)
('influence', 2)
('is', 2)
('jobs', 2)
('or', 2)
('potential', 2)
('studies', 2)
('tasks', 2)
('traditional', 2)
('which', 2)
('occupations', 3)
('of', 4)
('the', 5)
('and', 6)
('AI', 8)
('to', 9)
sdesaleg137@cloudshell:~$
```

## Delete the Dataproc cluster and the bucket:

1. To delete a Dataproc cluster and the bucket we created on GCP, you can follow these

steps:

- ✓ Open the Google Cloud Console.
- ✓ Navigate to the Dataproc Clusters page in the console.
- ✓ Select the cluster you wish to delete.
- ✓ Click the 'Delete' button at the top of the page.



2. In the confirmation dialog box, click on 'confirm'.

3. Navigate to the Cloud Storage page.

Name: Sened Desalegn

4. Select the bucket you want to delete.

5. Click on the 'Delete' button at the top of the page.



6. In the confirmation dialog box, type 'DELETE'.

Reference

1. https://www.youtube.com/watch?v=_lwrfxE2RtE&ab_channel=SkillCurb
2. https://hc.labnet.sfbu.edu/~henry/npu/classes/cloud_computing/pyspark/hw/q2/2023_summer/CS570_week4_q1_19744_SriVardhan_Kotturu.pdf