**Name:** Harsh Chheda

**Roll Number:** 31031521005 / 22-15405

**Class:** Msc. Computer Science

**Subject**: Cloud Computing
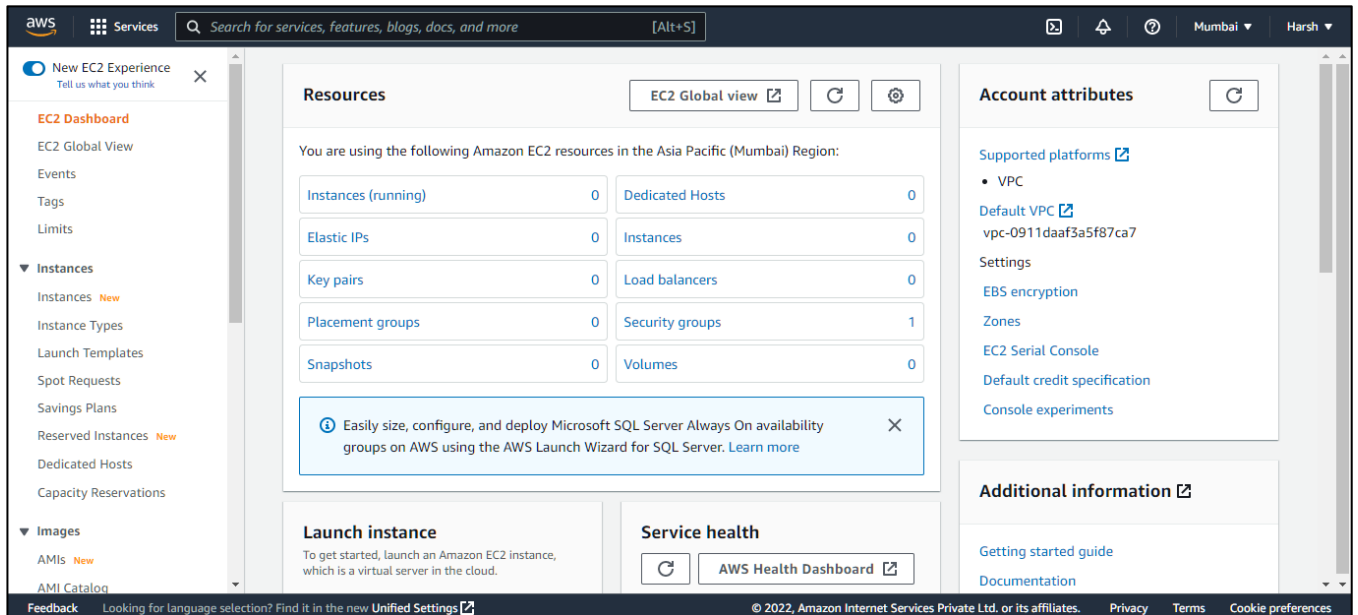
**Year:** 2022-23

# Practical 8

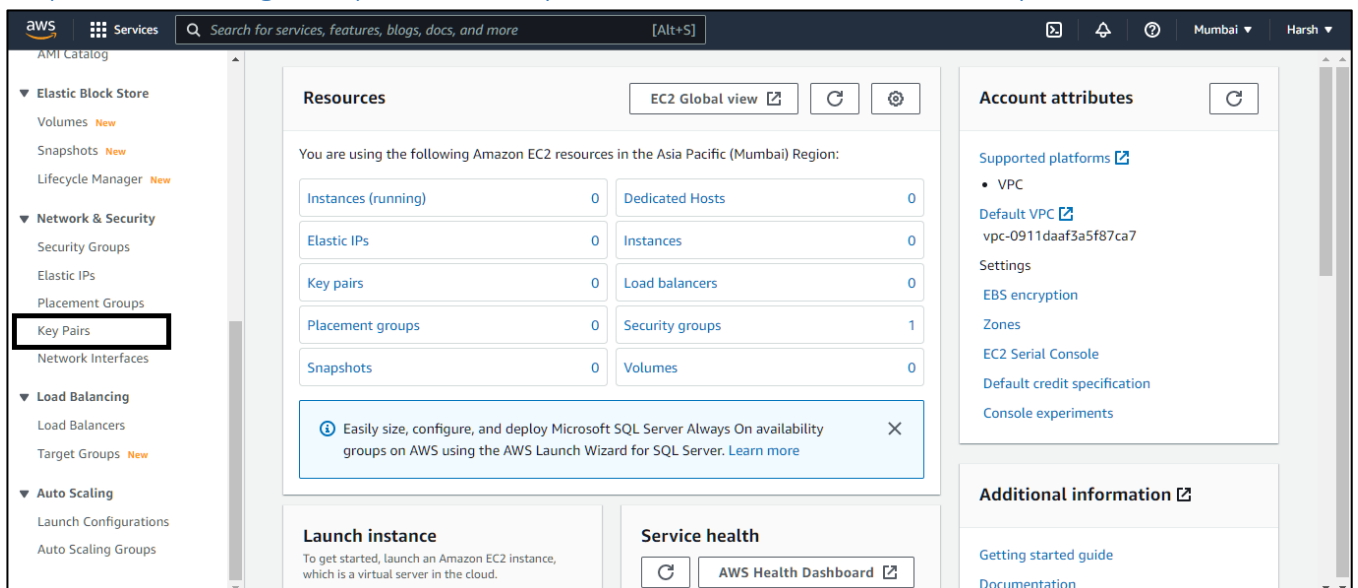**Aim:** Demonstration of data analytics in Cloud
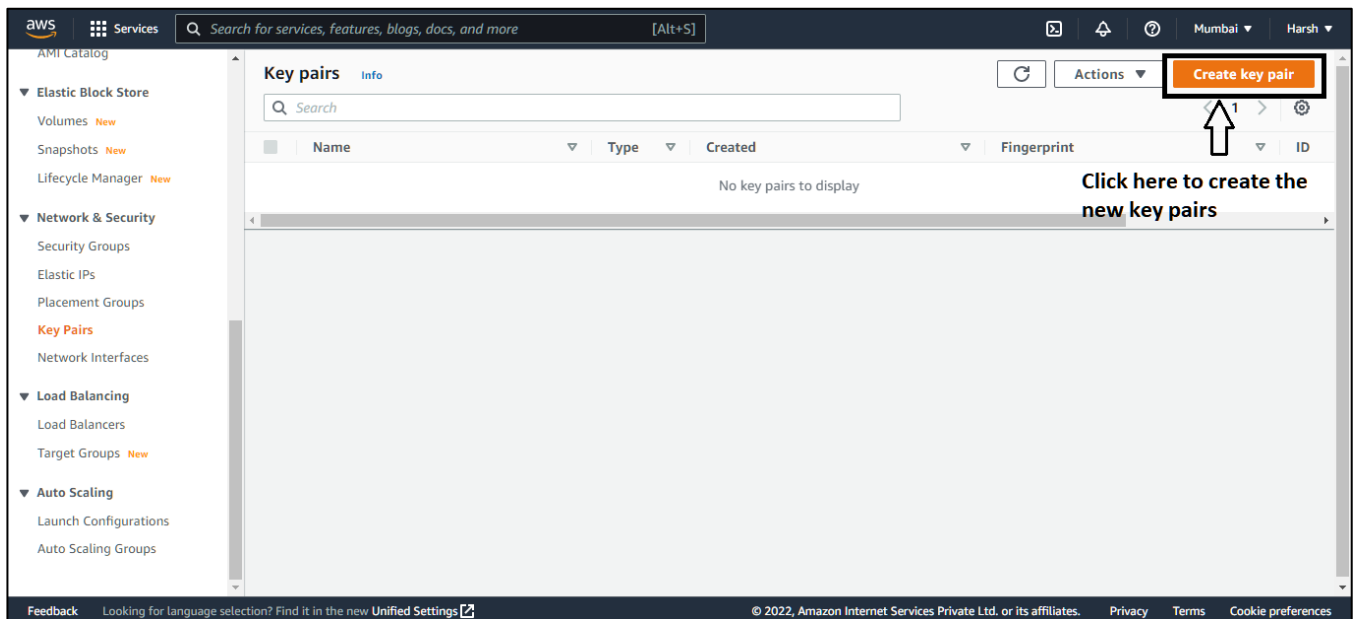
**Code:**

## Setting up EC2 Key-Pairs

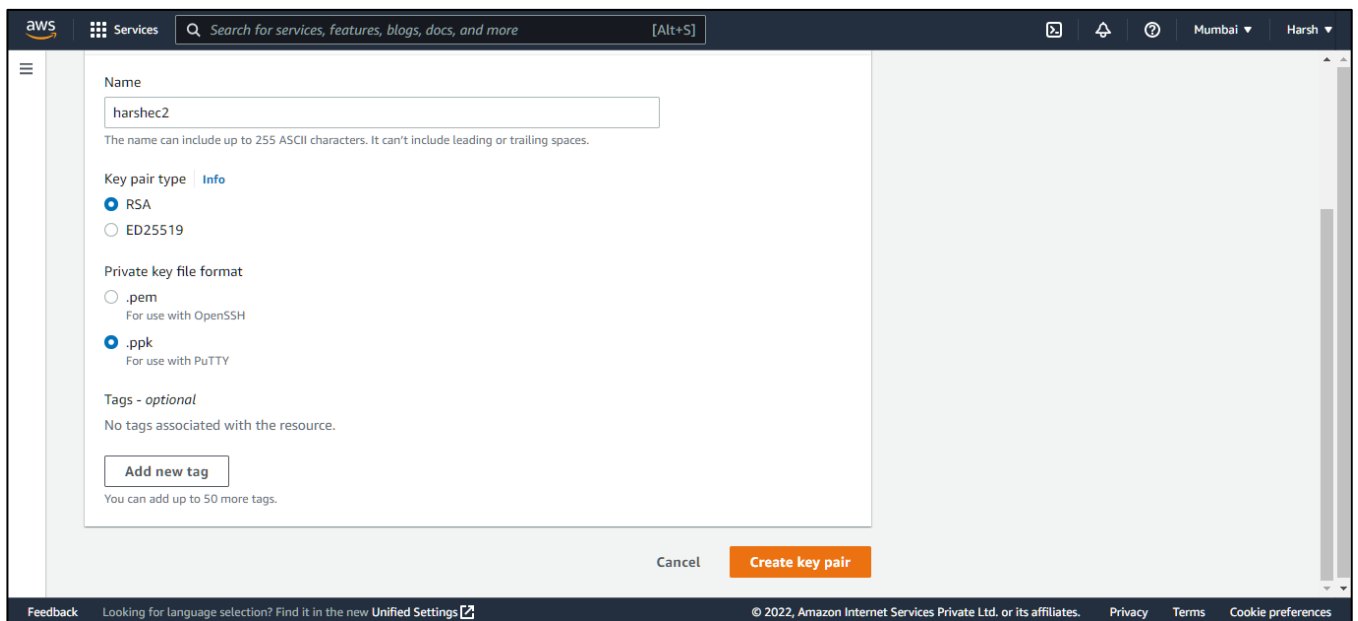### Step 1: Go to Amazon EC2 Console



### Step 2: In the Navigation pane, click Key Pairs under Network and Security Section

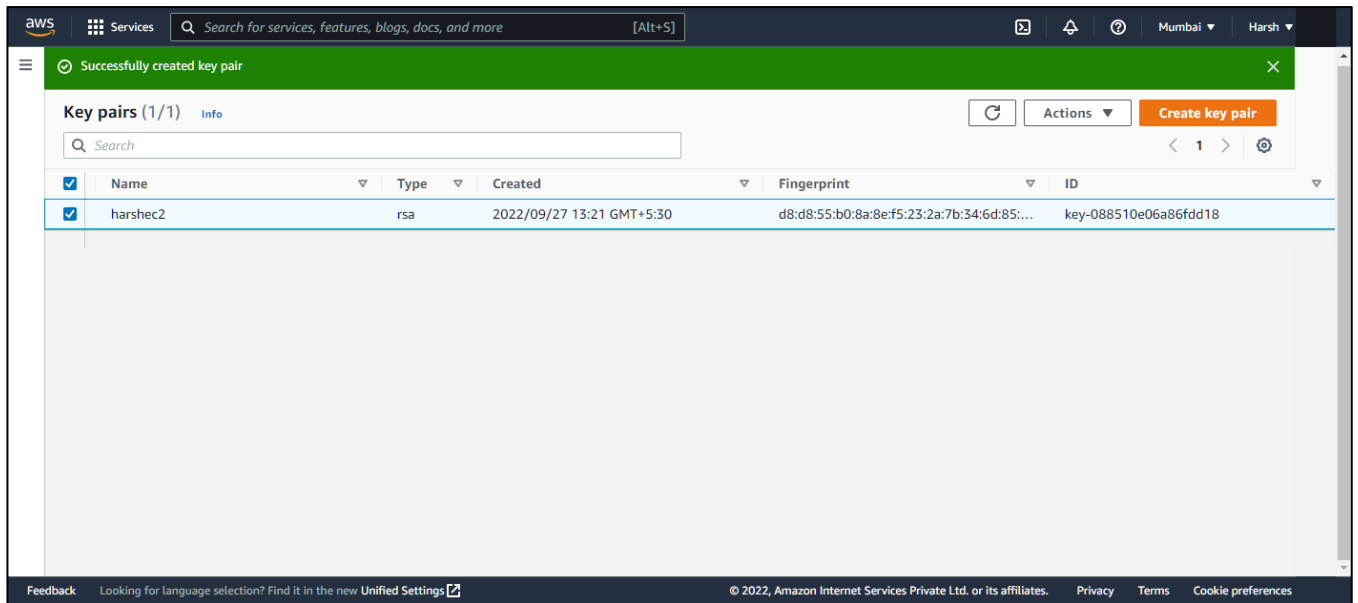## Step 3: On the Key Pairs page, click Create Key Pair



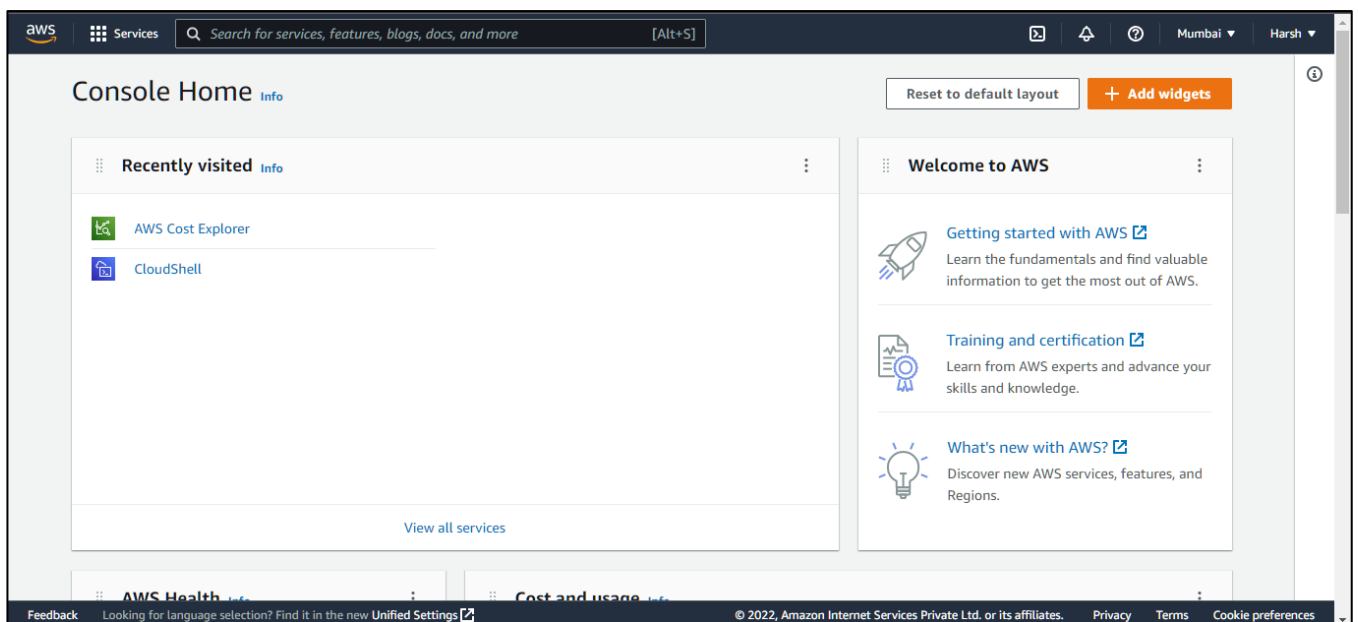## Step 4: In the Create Key Pair dialog box, enter a name for your key pair, such as, mykeypair



## Step 5: Click Create key Pair

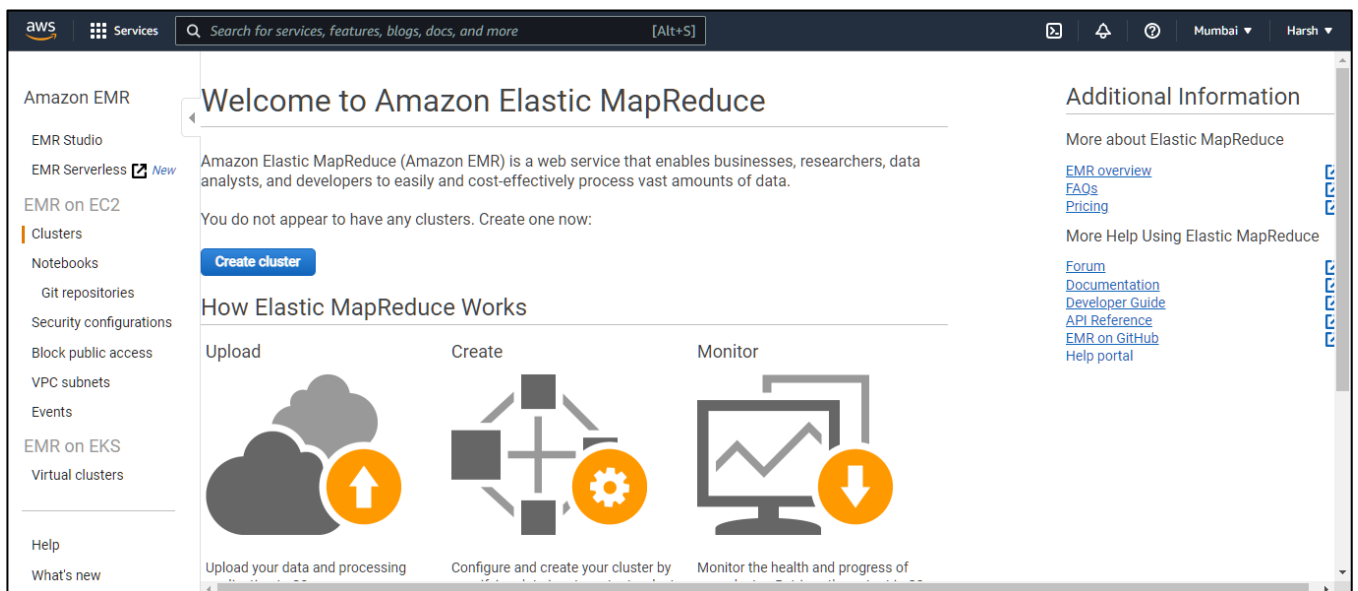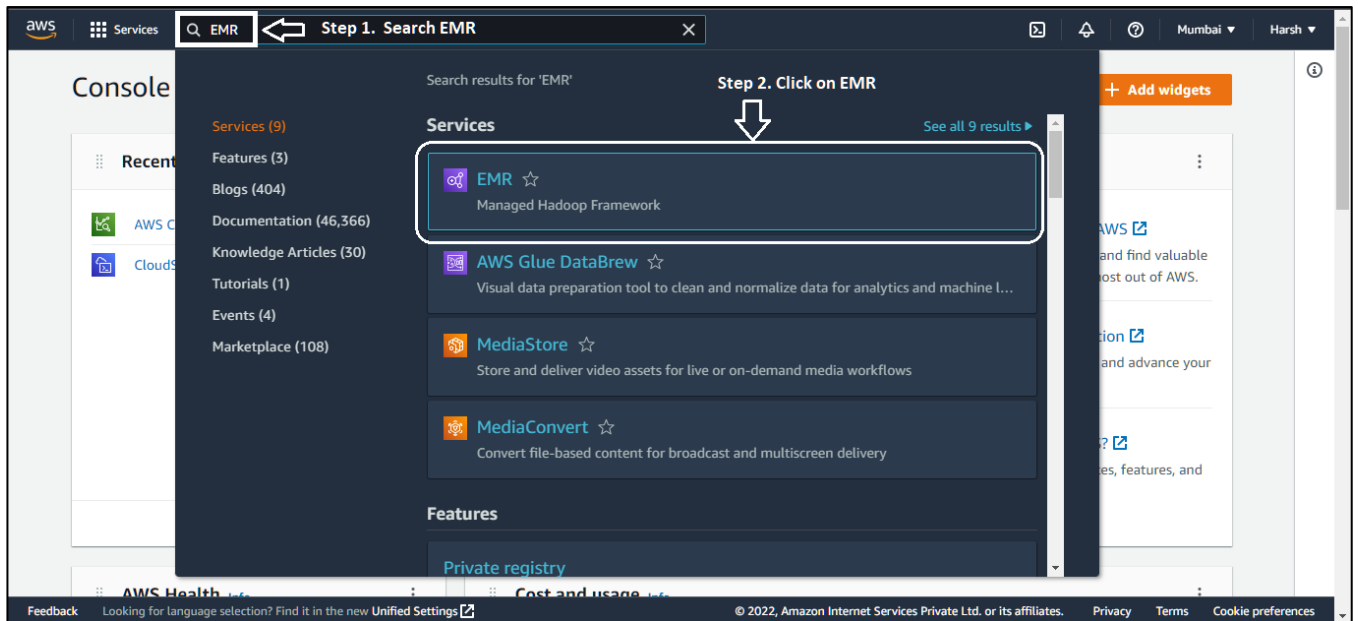## Step 6: Save the resulting PEM file in a safe location

# Setting up your environment on Amazon EMR

## Step 1: Create an AWS account and sign in to the console.



## Step 2: Search EMR in the Search Box

Cloud Computing    Msc. Computer Science    Roll Number: 31031521005    Name: Harsh Chheda

## Step 3: Creating new cluster

# Downloading Dataset

Step 1: **Click Here** to download dataset

# Setting up S3 Environment

Step 1: Search **S3** in the Search Box

## Step 2: To create new Bucket Click on Create Bucket

Step 3: Click on Create Bucket

Step 4: Once the bucket is created you will be able to see the bucket



# Uploading Dataset to the S3 Bucket

Step 1: Click and open the S3 bucket

Step 2: Click on Create Folder and create the new folder with the name data-source and click on Create folder.

## Step 3: Now open the Folder and upload the dataset.

Cloud Computing          Msc. Computer Science          Roll Number: 31031521005          Name: Harsh Chheda

# Code

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col


# FIND THE S3 URI IN THE S3 BUCKET
# PATH  S3_URI/FILENAME
S3_DATA_SOURCE_PATH="s3://harshc3294awsbucket/data-source/survey_results_public.csv"
S3_DATA_OUTPUT_PATH="s3://harshc3294awsbucket/data-output"

def main ():
    spark= SparkSession.builder.appName("HarshDemoApp").getOrCreate()
    all_data=spark.read.csv(S3_DATA_SOURCE_PATH,header=True)
    print("The total number of records int the source data : %s" % all_data.count())
    selected_data = all_data.where((col("Country")=="United States") &
(col("WorkWeekHrs")>45))
    print("The number of engineers who worked more than 45 hours a week in the US are:
%s" % selected_data.count())
    selected_data.write.mode("overwrite").parquet(S3_DATA_OUTPUT_PATH)
    print("Selected data was successfully saved to S3 %s"% S3_DATA_OUTPUT_PATH)

if __name__==    "__main__":
    main()
```

# Setting up the Security in EMR

## Step 1: Open EMR and click on the instance that is created. Scroll Down to **Security groups for Master**



## Step 2: Select the master node security group

Step 3. Click on Inbound Rules and click on edit inbound rules

## Step 4: Add new Rule for SSH and click on Save Rule.



# Running PYSPARK Cluster

Step 1: Open EMR and click on the instances that is created and click on **Connect to the Master Node Using SSH**


Step 2: Download PuTTY.exe to your computer from:
http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html


Step 3: Start PuTTY.

Step 4: In the Category list, click Session.

Step 5: In the Host Name field, type hadoop@ec2-15-206-168-25.ap-south-1.compute.amazonaws.com

Step 6: In the Category list, expand Connection > SSH, and then click Auth.

Step 7: For Private key file for authentication, click Browse and select the private key file (**harshec2.ppk**) used to launch the cluster.

Step 8: Click Open.

Step 9: Click Yes to dismiss the security alert.



Step 10: open vi main.py

Step 11: copy the code and press i in the terminal and paste the code

Step 12: ESC then :wq

Step 13: spark-submit main.py

## Output

```
hadoop@ip-172-31-27-104:~                                                                                  —  □  ×
E::::::::::::::::::E M:::::M          M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM          MMMMMMM RRRRRRR     RRRRRR

[hadoop@ip-172-31-27-104 ~]$ vi main.py
[hadoop@ip-172-31-27-104 ~]$ spark-submit main.py
22/09/27 09:23:57 INFO SparkContext: Running Spark version 2.4.7-amzn-1
22/09/27 09:23:57 INFO SparkContext: Submitted application: HarshDemoApp
22/09/27 09:23:57 INFO SecurityManager: Changing view acls to: hadoop
22/09/27 09:23:57 INFO SecurityManager: Changing modify acls to: hadoop
22/09/27 09:23:57 INFO SecurityManager: Changing view acls groups to:
22/09/27 09:23:57 INFO SecurityManager: Changing modify acls groups to:
22/09/27 09:23:57 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users  with view permissions: Set(hadoop); groups with view permissi
ons: Set(); users  with modify permissions: Set(hadoop); groups with modify permissions: Set()
22/09/27 09:23:58 INFO Utils: Successfully started service 'sparkDriver' on port 36871.
22/09/27 09:23:58 INFO SparkEnv: Registering MapOutputTracker
22/09/27 09:23:58 INFO SparkEnv: Registering BlockManagerMaster
22/09/27 09:23:58 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
22/09/27 09:23:58 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
22/09/27 09:23:58 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-dc043fd2-2bd0-4c86-a96a-81c2f0510e7f
22/09/27 09:23:58 INFO MemoryStore: MemoryStore started with capacity 912.3 MB
22/09/27 09:23:58 INFO SparkEnv: Registering OutputCommitCoordinator
22/09/27 09:23:58 INFO Utils: Successfully started service 'SparkUI' on port 4040.
22/09/27 09:23:58 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-27-104.ap-south-1.compute.internal:4040
22/09/27 09:23:58 INFO Utils: Using initial executors = 50, max of spark.dynamicAllocation.initialExecutors, spark.dynamicAllocation.minExecutors and spark.executor.ins
tances
22/09/27 09:23:59 INFO RMProxy: Connecting to ResourceManager at ip-172-31-27-104.ap-south-1.compute.internal/172.31.27.104:8032
22/09/27 09:23:59 INFO Client: Requesting a new application from cluster with 2 NodeManagers
22/09/27 09:23:59 INFO Configuration: resource-types.xml not found
22/09/27 09:23:59 INFO ResourceUtils: Unable to find 'resource-types.xml'.
22/09/27 09:23:59 INFO ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
22/09/27 09:23:59 INFO ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
22/09/27 09:23:59 INFO Client: Verifying our application has not requested more than the maximum memory capability of the cluster (12288 MB per container)
22/09/27 09:23:59 INFO Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
22/09/27 09:23:59 INFO Client: Setting up container launch context for our AM
22/09/27 09:23:59 INFO Client: Setting up the launch environment for our AM container
22/09/27 09:23:59 INFO Client: Preparing resources for our AM container
22/09/27 09:23:59 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
22/09/27 09:24:01 INFO Client: Uploading resource file:/mnt/tmp/spark-f9ae28be-37b8-49ef-82b6-b22e195f664a/__spark_libs__2302063699033026815.zip -> hdfs://ip-172-31-27-
104.ap-south-1.compute.internal:8020/user/hadoop/.sparkStaging/application_1664270493255_0001/__spark_libs__2302063699033026815.zip
22/09/27 09:24:03 INFO Client: Uploading resource file:/usr/lib/spark/python/lib/pyspark.zip -> hdfs://ip-172-31-27-104.ap-south-1.compute.internal:8020/user/hadoop/.sp
arkStaging/application_1664270493255_0001/pyspark.zip
22/09/27 09:24:03 INFO Client: Uploading resource file:/usr/lib/spark/python/lib/py4j-0.10.7-src.zip -> hdfs://ip-172-31-27-104.ap-south-1.compute.internal:8020/user/ha
doop/.sparkStaging/application_1664270493255_0001/py4j-0.10.7-src.zip
22/09/27 09:24:03 INFO Client: Uploading resource file:/mnt/tmp/spark-f9ae28be-37b8-49ef-82b6-b22e195f664a/__spark_conf__3160276972070887987.zip -> hdfs://ip-172-31-27-
```

```
22/09/27 09:41:47 INFO DAGScheduler: ResultStage 3 (count at NativeMethodA
22/09/27 09:41:47 INFO DAGScheduler: Job 2 finished: count at NativeMethodA
The total number of records int the source data : 64461
22/09/27 09:41:47 INFO FileSourceStrategy: Pruning directories with:
22/09/27 09:41:47 INFO FileSourceStrategy: Post-Scan Filters: isnotnull(Cou
> 45)
```

```
22/09/27 09:41:48 INFO DAGScheduler: ResultStage 6 (count at NativeMethodAccessorImpl.j
22/09/27 09:41:48 INFO DAGScheduler: Job 4 finished: count at NativeMethodAccessorImpl.
The number of engineers who worked more than 45 hours a week in the US are: 1527
22/09/27 09:41:48 INFO FileSourceStrategy: Pruning directories with:
22/09/27 09:41:48 INFO FileSourceStrategy: Post-Scan Filters: isnotnull(Country#18),isn
> 45)
```

```
22/09/27 09:41:52 INFO FileFormatWriter: Finished processing stats for write job 71a55dc8-d856-4399-
Selected data was successfully saved to S3 s3://harshc3294awsbucket/data-output
22/09/27 09:41:52 INFO SparkContext: Invoking stop() from shutdown hook
22/09/27 09:41:52 INFO SparkUI: Stopped Spark web UI at http://ip-172-31-27-104.ap-south-1.compute.i
```