# Contents

# 1    pin up a new EMR cluster

Go to aws.amazon.com and sign in to the console



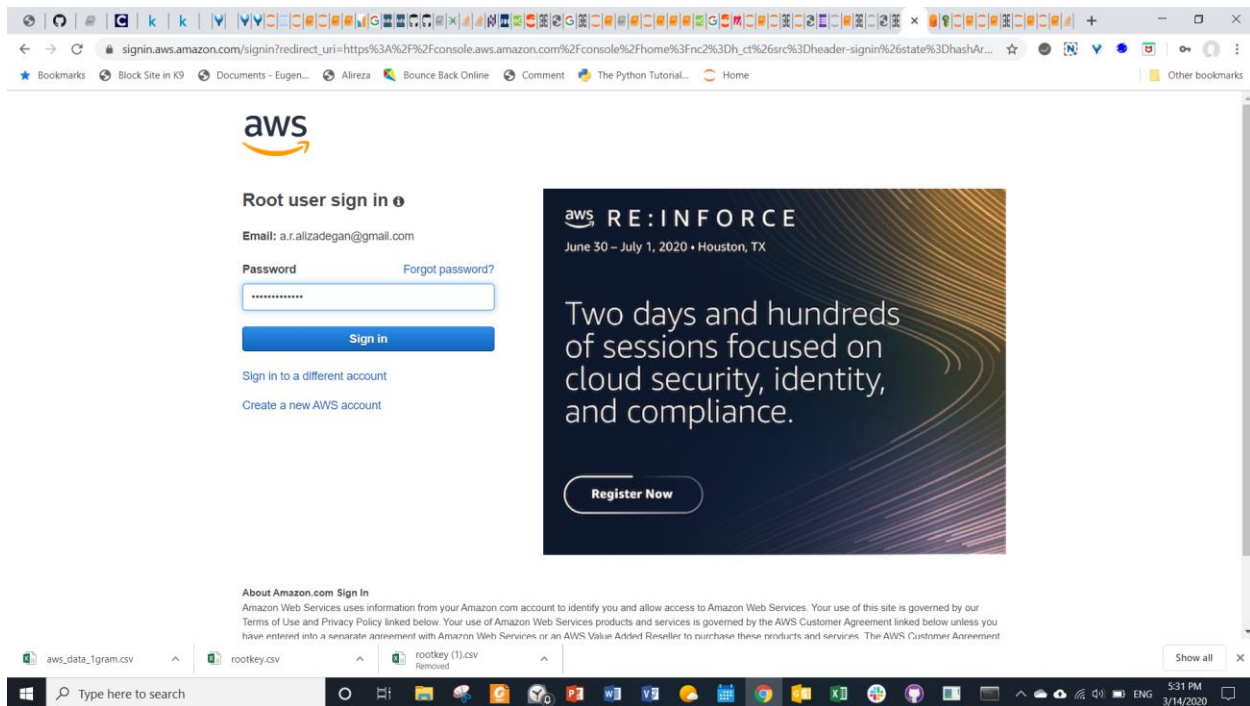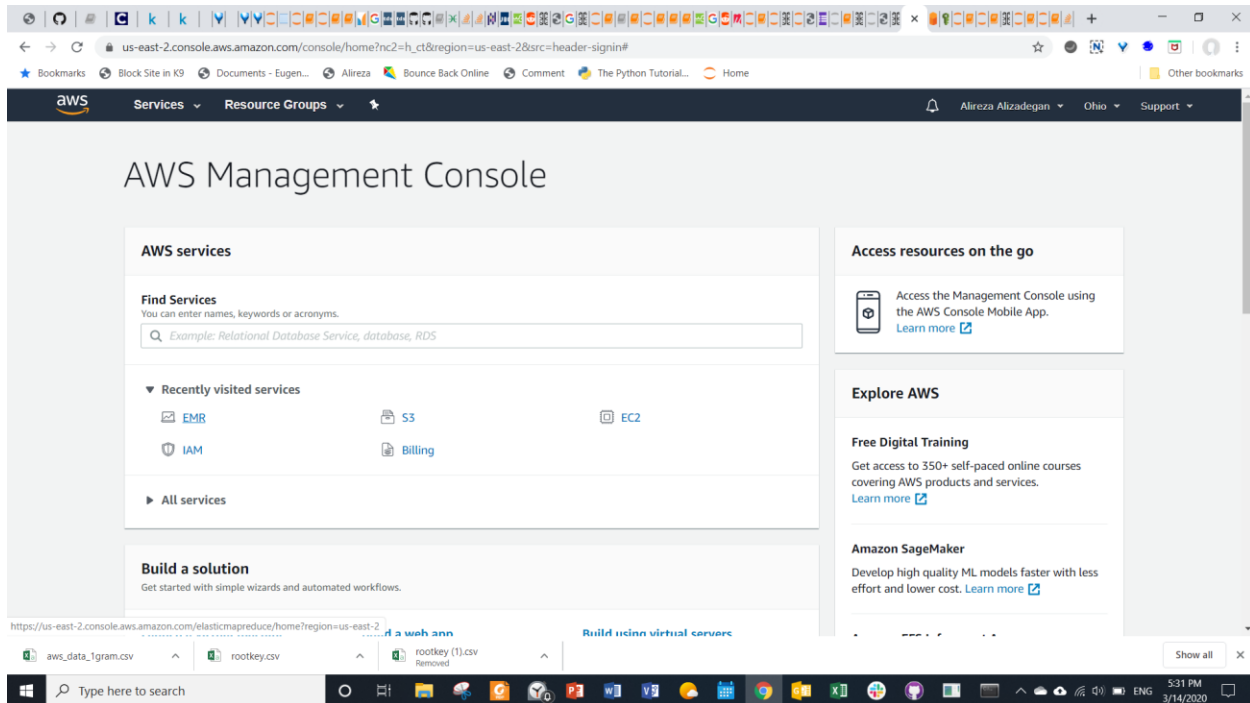Enter root user email address and click next

Enter password and sign in



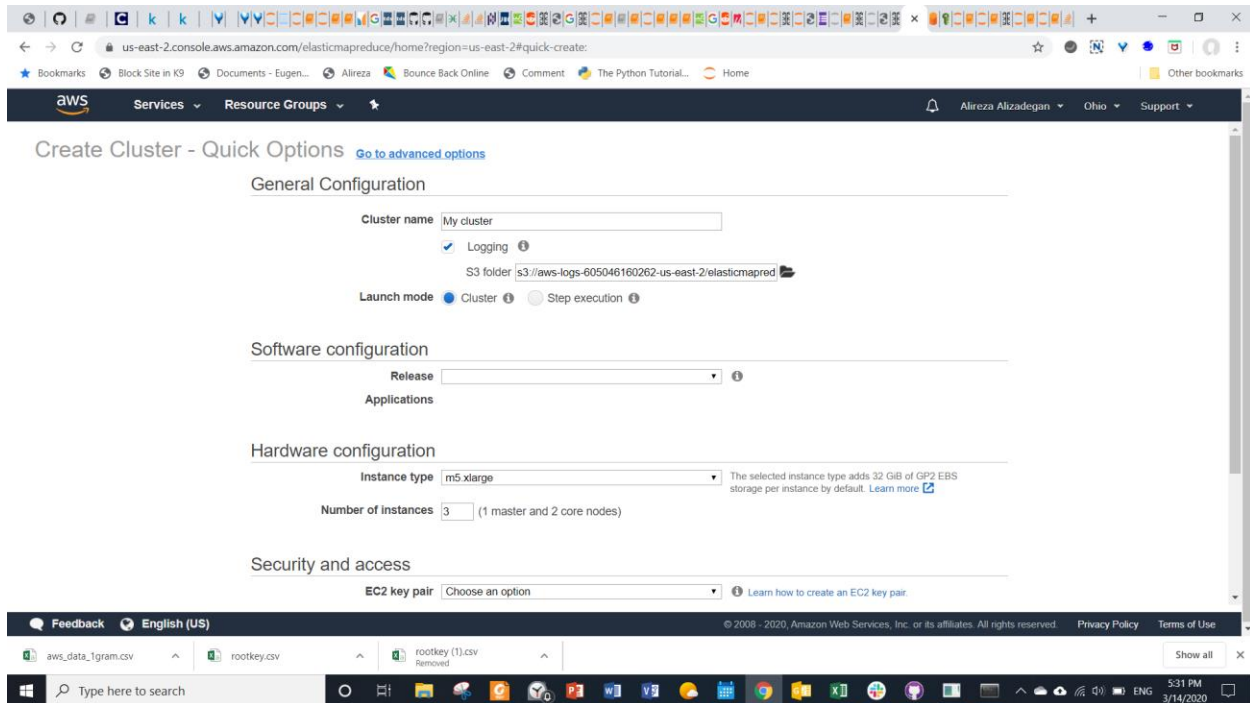Choose EMR from AWS services in their management console

Click create cluster in EMR clusters page



Go to advanced options in the options page

Select Hadoop, jupyterhub, hive, and spark in software & steps step



Do not make any changes in hardware step

Name the cluster in the general cluster settings step



Select the EC2 key pair in security options section of security step and click create cluster

Starting status shows up in green font in front of the cluster name as below



Wait until status changes to running as below then click on SSH link in master public DNS

# 2 Connect to master node of cluster

Add location of the private key to provided command in step 2 of instructions



Paste command in git bash and enter Hadoop environment as below

# 3 Copy data from S3 bucket to Hadoop directory

Create a directory in Hadoop environment



Copy data from brainstation public repository to directory by pasting following command in git bash

hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram

# 4 Analyze the data with spark

Click enable web connections link in front of connections



Add location of the private key to provided command in step 2 of instructions and paste it in git bash as below

Click the activated JupyterHub link in connections as below



Enter 'jovyan' for username and 'jupyter' for password and sign in

Create new pyspark3 notebook and run commands as attached jupyter notebook



Verify that results of analytics in spark are saved to Hadoop by listing content

# 5   Merge data into master node and copy to S3 bucket

Collect the data to master node as a CSV file



Copy the file from master node to personal AWS S3 bucket 'lastassignmentbucket'

```
hadoop@ip-172-31-2-171:~

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
20 package(s) needed for security, out of 37 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMM        MMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::E M:::::::M        M:::::::M R:::::::::::::::R
EE::::EEEEEEEEE::::E M::::::::M      M::::::::M R::::::RRRRR:::::R
  E::::E       EEEEE M:::::::::M    M:::::::::M RR::::R      R::::R
  E::::E             M::::::M:::M  M:::M::::::M   R:::R      R::::R
  E:::::EEEEEEEEEE    M::::::M M:::M M:::M::::::M   R:::RRRRRR::::R
  E::::::::::::::E    M::::::M  M:::M:::M M::::::M   R:::::::::::RR
  E:::::EEEEEEEEEE    M::::::M   M:::::M  M::::::M   R:::RRRRRR::::R
  E::::E             M::::::M    M:::M   M::::::M   R:::R      R::::R
  E::::E       EEEEE M::::::M     MMM    M::::::M   R:::R      R::::R
EE::::EEEEEEEE::::E M::::::M            M::::::M   R:::R      R::::R
E::::::::::::::::::E M::::::M            M::::::M RR::::R      R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM            MMMMMMM RRRRRR      RRRRRR

[hadoop@ip-172-31-2-171 ~]$ hadoop fs -mkdir /user/hadoop/eng_1M_1gram
[hadoop@ip-172-31-2-171 ~]$ hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram
20/03/15 00:42:56 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, skipCRC=false, blocking=true, numListstatusThr
eads=0, maxMaps=20, mapBandwidth=100, sslConfigurationFile='null', copyStrategy='uniformsize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=
[s3://brainstation-dsft/eng_1M_1gram.csv], targetPath=/user/hadoop/eng_1M_1gram, targetPathExists=true, filtersFile='null'}
20/03/15 00:42:56 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-2-171.us-east-2.compute.internal/172.31.2.171:8032
20/03/15 00:42:58 INFO executor.GlobalS3Executor: Bucket brainstation-dsft is in the ca-central-1 region. Please configure the proper region to avoid multiple unnecessary redirects
20/03/15 00:42:59 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
20/03/15 00:42:59 INFO tools.SimpleCopyListing: Build file listing completed.
20/03/15 00:42:59 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
20/03/15 00:42:59 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
20/03/15 00:42:59 INFO tools.DistCp: Number of paths in the copy list: 1
20/03/15 00:42:59 INFO tools.DistCp: Number of paths in the copy list: 1
20/03/15 00:43:00 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-2-171.us-east-2.compute.internal/172.31.2.171:8032
20/03/15 00:43:00 INFO mapreduce.JobSubmitter: number of splits:1
[hadoop@ip-172-31-2-171 ~]$ hadoop fs -ls
Found 2 items
drwxr-xr-x   - hadoop hadoop          0 2020-03-15 00:44 eng_1M_1gram
drwxr-xr-x   - livy   hadoop          0 2020-03-15 01:00 eng_data_1gram
[hadoop@ip-172-31-2-171 ~]$ hadoop fs -getmerge /user/hadoop/eng_data_1gram collected_data_1gram.csv
[hadoop@ip-172-31-2-171 ~]$ aws s3 cp collected_data_1gram.csv s3://lastassignmentbucket/aws_data_1gram.csv
upload: ./collected_data_1gram.csv to s3://lastassignmentbucket/aws_data_1gram.csv
[hadoop@ip-172-31-2-171 ~]$
```