# CS 643, CLOUD COMPUTING
# PROGRAMMING ASSIGNMENT 2

GitHub Link: https://github.com/keerthikalla/kk224-programming-assignment-2

Docker Repository: https://hub.docker.com/repository/docker/keerthikalla123/winequalitypred

**Input for model training:**

Create an S3 bucket to upload the training and validation datasets

**Setting up 4 parallel EC2 instances:**

1. Launch an EMR cluster with following specifications:



2. Go to security groups of EMR-Master to edit the inbound rules to allow ssh to be done.

3. Click on edit inbound rules and add the following rule and click save.

| - | SSH | ▼ | TCP | 22 | Anywhe... ▼ | Q | | Delete |
|---|---|---|---|---|---|---|---|---|
| | | | | | | 0.0.0.0/0 ✕ | | |

Add rule

Cancel    Preview changes    **Save rules**

4. Connect to the Master Node using SSH as follows:

Cluster: winepredcluster    Waiting  Cluster ready after last step completed.

Summary | Application user interfaces | Monitoring | Hardware | Configurations | Events | Steps | Bootstrap actions

Summary

**SSH**                                                                ✕

**Connect to the Master Node Using SSH**

You can connect to the Amazon EMR master node using SSH to run interactive queries, examine log files, submit Linux commands, and so on.
Learn more ⧉.

| Windows | Mac / Linux |

1. Open a terminal window. On Mac OS X, choose Applications > Utilities > Terminal. On other Linux distributions, terminal is typically found at Applications > Accessories > Terminal.
2. To establish a connection to the master node, type the following command. Replace ~/EC2-A.pem with the location and filename of the private key file (.pem) used to launch the cluster.

    ssh -i ~/EC2-A.pem hadoop@ec2-52-90-87-86.compute-1.amazonaws.com

3. Type yes to dismiss the security warning.

                                                                Close

Auto-termination: Terminate if idle for 1 hour

Security and access
        Key name: EC2-A
    EC2 instance profile: EMR_EC2_DefaultRole

hadoop@ip-172-31-81-21:~

```
C:\Users\nnade>ssh -i D:/EC2-A.pem hadoop@ec2-34-230-38-110.compute-1.amazonaws.com
The authenticity of host 'ec2-34-230-38-110.compute-1.amazonaws.com (34.230.38.110)' can't be established.
ECDSA key fingerprint is SHA256:JlVsThzdT4kCdYZCXp7rnuq+iHFQ5IY8dfZ3/16Vv7Q.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-34-230-38-110.compute-1.amazonaws.com,34.230.38.110' (ECDSA) to the list of known hosts.
Last login: Thu Dec  8 18:03:01 2022

       __|  __|_  )
       _|  (     /   Amazon Linux 2 AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
37 package(s) needed for security, out of 50 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEEE MMMMMMMM           MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M         M:::::::M R::::::::::::::::R
EE:::::EEEEEEEEE:::E M::::::::M       M::::::::M R:::::RRRRRR:::::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::::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 RRRRRRR      RRRRRR

[hadoop@ip-172-31-81-21 ~]$
```

5. Create a .py file using the command "nano train.py".

Then write your code, save, and close by using the following command: Shift O + Enter + Shift X.

6. Install all necessary libraries using the "pip install <libraryname>" command.

7. Run the code using '"spark-submit train.py"



8. Go to the application interface and select spark history server.



9. Check status of your job.

**Spark Jobs** <sup>(?)</sup>

**User:** hadoop
**Total Uptime:** 20 s
**Scheduling Mode:** FIFO
**Completed Jobs:** 37

▸ Event Timeline

▾ **Completed Jobs (37)**

| Job Id ▾ | Description | Submitted | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|---|---|---|---|---|---|
| 36 | parquet at treeEnsembleModels.scala:446<br>parquet at treeEnsembleModels.scala:446 | 2022/12/08 18:46:56 | 0.9 s | 1/1 | 1/1 |
| 35 | runJob at SparkHadoopWriter.scala:78<br>runJob at SparkHadoopWriter.scala:78 | 2022/12/08 18:46:55 | 0.2 s | 1/1 | 1/1 |
| 34 | count at /home/hadoop/winetrain.py:94<br>count at /home/hadoop/winetrain.py:94 | 2022/12/08 18:46:55 | 14 ms | 1/1 | 1/1 |
| 33 | count at /home/hadoop/winetrain.py:94<br>count at /home/hadoop/winetrain.py:94 | 2022/12/08 18:46:54 | 0.3 s | 1/1 | 1/1 |
| 32 | toPandas at /home/hadoop/winetrain.py:82<br>toPandas at /home/hadoop/winetrain.py:82 | 2022/12/08 18:46:54 | 0.1 s | 1/1 | 1/1 |
| 31 | showString at NativeMethodAccessorImpl.java:0<br>showString at NativeMethodAccessorImpl.java:0 | 2022/12/08 18:46:54 | 0.1 s | 1/1 | 1/1 |
| 30 | runJob at PythonRDD.scala:153<br>runJob at PythonRDD.scala:153 | 2022/12/08 18:46:54 | 0.1 s | 1/1 | 1/1 |
| 29 | runJob at PythonRDD.scala:153<br>runJob at PythonRDD.scala:153 | 2022/12/08 18:46:53 | 0.4 s | 1/1 | 1/1 |
| 28 | collectAsMap at RandomForest.scala:567<br>collectAsMap at RandomForest.scala:567 | 2022/12/08 18:46:53 | 27 ms | 2/2 | 2/2 |
| 27 | collectAsMap at RandomForest.scala:567<br>collectAsMap at RandomForest.scala:567 | 2022/12/08 18:46:53 | 22 ms | 2/2 | 2/2 |
| 26 | collectAsMap at RandomForest.scala:567<br>collectAsMap at RandomForest.scala:567 | 2022/12/08 18:46:53 | 32 ms | 2/2 | 2/2 |
| 25 | collectAsMap at RandomForest.scala:567<br>collectAsMap at RandomForest.scala:567 | 2022/12/08 18:46:53 | 32 ms | 2/2 | 2/2 |
| 24 | collectAsMap at RandomForest.scala:567<br>collectAsMap at RandomForest.scala:567 | 2022/12/08 18:46:53 | 47 ms | 2/2 | 2/2 |

## **Developing a Spark Application (Model Implementation):**

## 1. Launch an EC2 instance as follows



## 2. Connect to your new instance from the local terminal

```
D:\>ssh -i "EC2-A" ec2-54-167-35-118.compute-1.amazonaws.com
Warning: Identity file EC2-A not accessible: No such file or directory.
The authenticity of host 'ec2-54-167-35-118.compute-1.amazonaws.com (54.167.35.118)' can't be established.
ECDSA key fingerprint is SHA256:dvwcc0bkqcWRmz/6Gg07+9JvA4rREDb98VTcoRF+Rdk.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-54-167-35-118.compute-1.amazonaws.com,54.167.35.118' (ECDSA) to the list of known hosts.
keerthi@ec2-54-167-35-118.compute-1.amazonaws.com: Permission denied (publickey,gssapi-keyex,gssapi-with-mic).

D:\>ssh -i EC2-A.pem  ec2-54-167-35-118.compute-1.amazonaws.com
keerthi@ec2-54-167-35-118.compute-1.amazonaws.com: Permission denied (publickey,gssapi-keyex,gssapi-with-mic).

D:\>ssh -i EC2-A.pem  ec2user@ec2-54-167-35-118.compute-1.amazonaws.com
ec2user@ec2-54-167-35-118.compute-1.amazonaws.com: Permission denied (publickey,gssapi-keyex,gssapi-with-mic).

D:\>ssh -i EC2-A.pem  ec2-user@ec2-54-167-35-118.compute-1.amazonaws.com

       __|  __|_  )
       _|  (     /   Amazon Linux 2 AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
19 package(s) needed for security, out of 31 available
```

## 3. Install Scala



## 4. Install Spark

```
[ec2-user@ip-172-31-22-133 ~]$ wget https://archive.apache.org/dist/spark/spark-2.4.5/spark-2.4.5-bin-hadoop2.7.tgz
--2022-12-08 20:12:01--  https://archive.apache.org/dist/spark/spark-2.4.5/spark-2.4.5-bin-hadoop2.7.tgz
Resolving archive.apache.org (archive.apache.org)... 138.201.131.134, 2a01:4f8:172:2ec5::2
Connecting to archive.apache.org (archive.apache.org)|138.201.131.134|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 232530699 (222M) [application/x-gzip]
Saving to: 'spark-2.4.5-bin-hadoop2.7.tgz'

100%[==============================================================================>] 232,530,699 16.2MB/s   in 14s

2022-12-08 20:12:16 (15.6 MB/s) - 'spark-2.4.5-bin-hadoop2.7.tgz' saved [232530699/232530699]

[ec2-user@ip-172-31-22-133 ~]$ sudo tar xvf spark-2.4.5-bin-hadoop2.7.tgz -C /opt
spark-2.4.5-bin-hadoop2.7/
spark-2.4.5-bin-hadoop2.7/licenses/
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-jtransforms.html
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-zstd.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-zstd-jni.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-xmlenc.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-vis.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-spire.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-sorttable.js.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-slf4j.txt
spark-2.4.5-bin-hadoop2.7/licenses/LICENSE-scopt.txt
```

5. Check Java -version. In case it doesn't exist, follow the instructions provided in this link.

https://techviewleo.com/install-java-openjdk-on-amazon-linux-system/

6. Check python version

```
[ec2-user@ip-172-31-22-133 ~]$ python --version
Python 2.7.18
[ec2-user@ip-172-31-22-133 ~]$ sudo yum -y install python-pip
Loaded plugins: extras_suggestions, langpacks, priorities, update-motd
amzn2-core                                                                    | 3.7 kB  00:00:00
Resolving Dependencies
--> Running transaction check
---> Package python2-pip.noarch 0:20.2.2-1.amzn2.0.3 will be installed
--> Finished Dependency Resolution

Dependencies Resolved

=================================================================================================== Package
                      Repository               Size
===================================================================================================Installing:
 python2-pip              noarch          20.2.2-1.amzn2.0.3            amzn2-core              2.0 M

Transaction Summary
===================================================================================================Install  1 Package

Total download size: 2.0 M
Installed size: 9.5 M
Downloading packages:
python2-pip-20.2.2-1.amzn2.0.3.noarch.rpm                                     | 2.0 MB  00:00:00
Running transaction check
Running transaction test
Transaction test succeeded
Running transaction
  Installing : python2-pip-20.2.2-1.amzn2.0.3.noarch                                              1/1
```

7. Create a .py file using the command "nano test.py".

Then write your code, save, and close by using the following command: Shift O + Enter + Shift X.

8. Install all necessary libraries using the "pip install <libraryname>" command.

9. Run the code using '"spark-submit test.py"

```
F1- score:  0.55
[[ 0   0   1   0   0   0]
 [ 0   0   2   0   0   0]
 [ 0   0  48  14   4   0]
 [ 0   0  17  34  14   0]
 [ 0   0   5  11   6   0]
 [ 0   0   0   3   1   0]]
/home/hadoop/.local/lib/python3.7/site-packages/sklearn/metrics/_classi
eing set to 0.0 in labels with no predicted samples. Use `zero_division
  _warn_prf(average, modifier, msg_start, len(result))
/home/hadoop/.local/lib/python3.7/site-packages/sklearn/metrics/_classi
eing set to 0.0 in labels with no predicted samples. Use `zero_division
  _warn_prf(average, modifier, msg_start, len(result))
/home/hadoop/.local/lib/python3.7/site-packages/sklearn/metrics/_classi
eing set to 0.0 in labels with no predicted samples. Use `zero_division
  _warn_prf(average, modifier, msg_start, len(result))
            precision    recall  f1-score   support

       3.0       0.00      0.00      0.00         1
       4.0       0.00      0.00      0.00         2
       5.0       0.66      0.73      0.69        66
       6.0       0.55      0.52      0.54        65
       7.0       0.24      0.27      0.26        22
       8.0       0.00      0.00      0.00         4

   accuracy                           0.55       160
  macro avg       0.24      0.25      0.25       160
weighted avg       0.53      0.55      0.54       160

Accuracy 0.55
```

**Predicting using Docker Images:**

Launch your ec2-instance and then step-up docker as follows:

1. Go to your Docker repository

2. Pull the image to Docker hub repository by using the following command

"docker pull keerthikalla123/winequalitypred:tag"

3. Run the image using the following command

"docker run -p 4000:80 keerthikalla123/winequalitypred:tag"

Accuracy and F1 score will be displayed accordingly.

```
[ec2-user@ip-172-31-94-198 ~]$ history
    1  sudo yum update
    2  sudo yum install docker
    3  sudo service docker start
    4  sudo usermod -a -G docker ec2-user
    5  exit
    6  docker info
```