**Name: Swathi Muralinathan**

**Email Address: sm3345@njit.edu**

**Subject : CS643853-Cloud Computing**

**Assignment:**

**Module 03 Assignment 03: Programming**

**Assignment 2**

**Wine Quality Prediction AWS Spark Application:**

**Pa2Winepred:** This project involves the development of a Python application utilizing the PySpark interface.

The application is deployed on an Amazon Web Services (AWS) Elastic MapReduce (EMR) cluster. The primary objective is to parallelly train a machine learning model on EC2 instances for predicting wine quality using publicly available data. Subsequently, the trained model is employed to predict the quality of wine. Docker is utilized to create a container image for the trained machine learning model, streamlining the deployment process.

**Link for GitHub:**

https://github.com/Swmural/Wine

**Link for Docker:**

Docker

**Steps for the Execution for Wine Quality Prediction AWS Spark :**

1.Create a Key-pair for the EMR Cluster :go to EC2/Network/Key-pairs

 Use the format of .pem and download the keypair
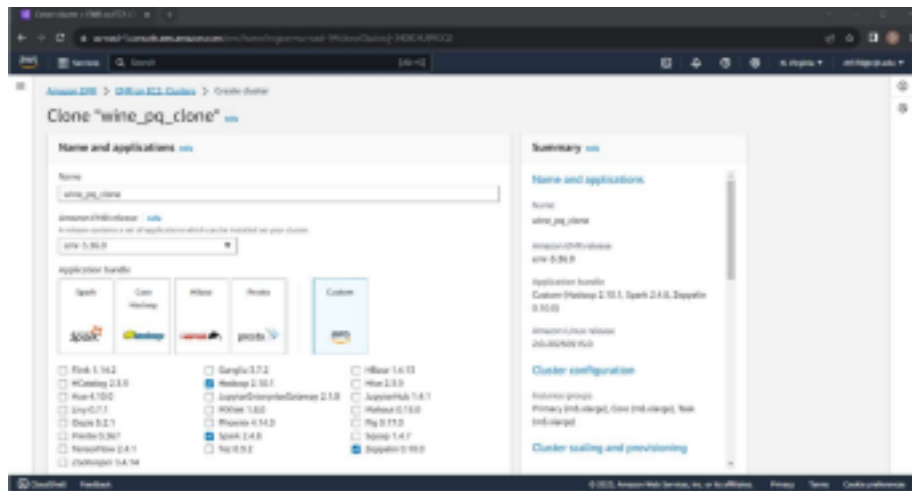
 Created key pair as: pa2assmahi.pem

2.Create an S3 bucket

 Created an S3 bucket in aws: pa2winebucket1

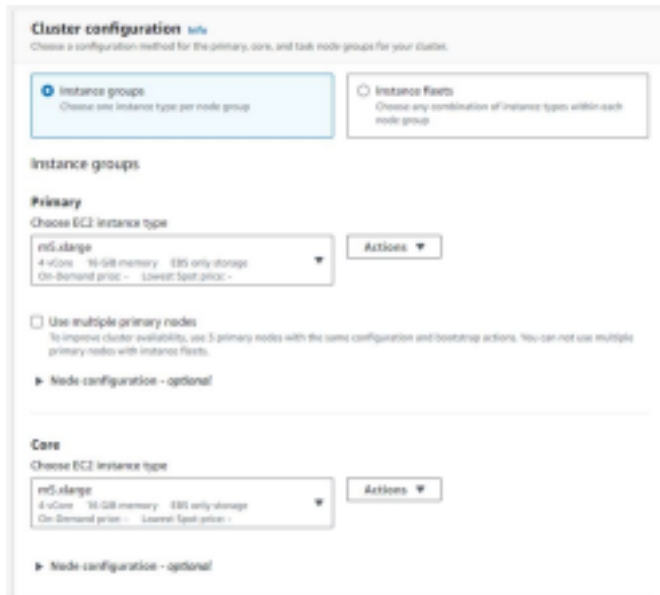3.Then go to EMR console and create EMR cluster

4. Creating the spark in the AWS instance by using EMR console:

 Creating the spark cluster by using the EMR console, and create the 4

instances:   Name and application

Note: Here it says Clone "wine_pq_clone" as I have cloned the previous configuration instead of creating from scratch to save time.

**Cluster Configuration:**

## Cluster Scaling and provisioning:



## Networking & Cluster Termination:

**Security Configuration and EC2 Key pair & Identity and access management(IAM) role:**



We can follow above steps for creating EMR cluster for the instances

5.Now we are training ML model into spark cluster with ec2 instances in

parallel:

1.Now the cluster will accept the tasks to run the ML model

Need to connect the Master instance in the Terminal:

ssh -i "pa2assmahi.pem"

ec2-user@ec2-44-201-107-82.compute-1.amazonaws.com  and it is

successfully login.

2.After the login of Master instance then change the root by using
Sudo su

3.

Submit the task by the command:

spark-submit s3://pa2winebucket1/ winequilityprediction.py

4.Then you can find the trace status for the above tasks, The status is succeed then there is a creation of test.model in the s3 bucket s3://pa2winebucket1



6.Now we are running ML model using the Docker:

1.Create an docker account and sign up.

2. After the successful login then download and setup the docker in your local system 3.Install the docker

4. Login the docker in the power shell by the command

docker login

Pwd

5.After login you need to build the image:

docker build -t winequlpred .

6. The push and pull into the docker hub repository:

PUSH:

docker tag winequlpred dt37824/winequlpred

docker push dt37824/winequlpred

PULL:

docker pull dt37824/winequlpred



7. Store your test data file in a designated folder, referred to as "dir." Mount this directory with the Docker container, and execute the container using the following command.

docker run -v C:\Pa2Winepred\data\csv winequlpred testdata.csv

# deeps2201/wine_quality_prediction_swathi:final

DIGEST: sha256:a407be76a81a9bcdc84261efbee14e66d1f4c0a4d92ac495c1e90fe52c2232cc

| OS/ARCH | COMPRESSED SIZE ⓘ | LAST PUSHED | TYPE |
|---|---|---|---|
| linux/amd64 | 556.59 MB | an hour ago by deeps2201 | Image |

## IMAGE LAYERS ⓘ

| # | | |
|---|---|---|
| 1 | **ADD file ... in /** | **72.57 MB** |
| 2 | LABEL org.label-schema.schema-version=1.0 org.label… | 0 B |
| 3 | CMD ["/bin/bash"] | 0 B |
| 4 | RUN /bin/sh -c yum -y | 206.97 MB |
| 5 | RUN /bin/sh -c python -V | 93 B |
| 6 | RUN /bin/sh -c python3 -V | 93 B |

**Command**

ADD file:b3ebbe8bd304723d43b7b44a6d996

| # | | |
|---|---|---|
| 3 | CMD ["/bin/bash"] | 0 B |
| 4 | RUN /bin/sh -c yum -y | 206.97 MB |
| 5 | RUN /bin/sh -c python -V | 93 B |
| 6 | RUN /bin/sh -c python3 -V | 93 B |
| 7 | ENV PYSPARK_DRIVER_PYTHON=python3 | 0 B |
| 8 | ENV PYSPARK_PYTHON=python3 | 0 B |
| 9 | RUN /bin/sh -c pip3 install | 4.48 MB |
| 10 | RUN /bin/sh -c pip3 install | 53.91 MB |
| 11 | RUN /bin/sh -c wget --no-verbose | 218.64 MB |
| 12 | RUN /bin/sh -c ln -s | 176 B |
| 13 | RUN /bin/sh -c echo 'export | 341 B |
| 14 | RUN /bin/sh -c mkdir /code | 110 B |
| 15 | RUN /bin/sh -c mkdir /code/data | 124 B |
| 16 | RUN /bin/sh -c mkdir /code/data/csv | 145 B |
| 17 | RUN /bin/sh -c mkdir /code/data/model | 145 B |
| 18 | RUN /bin/sh -c mkdir /code/src | 124 B |

```
17    RUN /bin/sh -c mkdir /code/data/model                145 B

18    RUN /bin/sh -c mkdir /code/src                       124 B

19    RUN /bin/sh -c mkdir /code/data/testdata.model/      150 B

20    COPY src/test.py /code/src # buildkit                957 B

21    COPY data/model/testmodel.model/ /code/data/mod…    4.79 KB

22    COPY data/csv/ /code/data/csv # buildkit           20.78 KB

23    RUN /bin/sh -c rm /bin/sh                            160 B

24    RUN /bin/sh -c /bin/bash -c                           93 B

25    RUN /bin/sh -c /bin/sh -c                             93 B

26    WORKDIR /code/                                        32 B

27    ENTRYPOINT ["/opt/spark/bin/spark-submit" "src/test…   0 B
```

docker          Why          Products          Developers          Company

Conclusion: As shown in the image above, got an accuracy of ~98% while predicting the wine  quality.