# Programming Assignment 2

**Name:** Utakarsh Aggarwal

**Class:** CS643-852


**Docker Hub Link:**

https://hub.docker.com/repository/docker/utakarshagg/winequalitytesting/general

**Github Link:**

https://github.com/UtakarshAgg/Wine-Quality-Testing---Programming-Assignment-2


**Instructions:**

1. Launch AWS Academy Learner Lab
2. Create EMR cluster
   2.1. Add up to 4 tasks inside the cluster under *Cluster Configuration*
   2.2. Create or upload a keypair (.ppk) under *Security Configuration*
   2.3. Change IAM roles to default roles under *Access Management*

3. Create S3 Bucket
   3.1. No configuration change required
4. Upload all required files in S3 bucket
   4.1. Provided datasets
   4.2. Training and prediction code
   4.3. Docker file



5. Open the EMR cluster
   5.1. Find 'Connect to primary node using SSH'
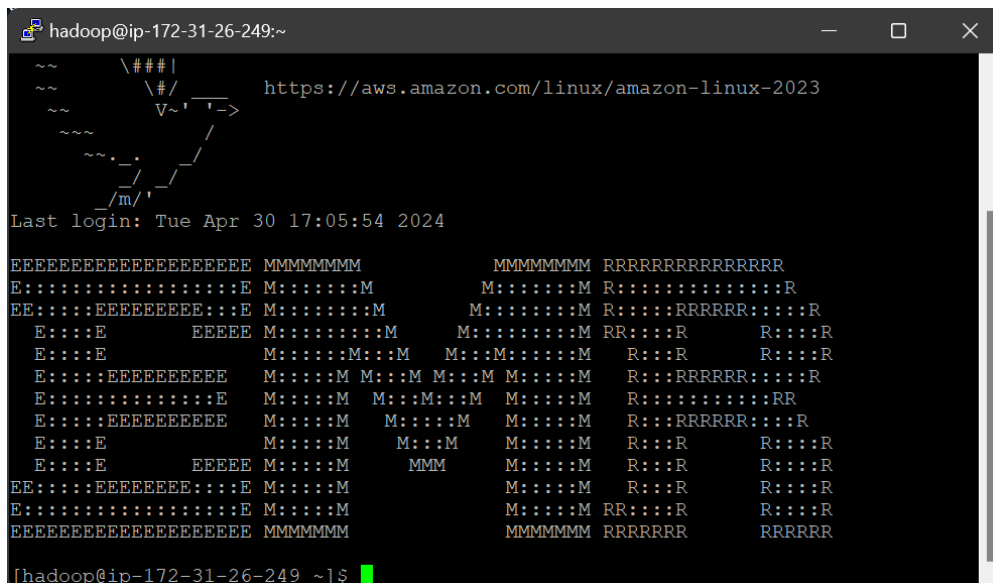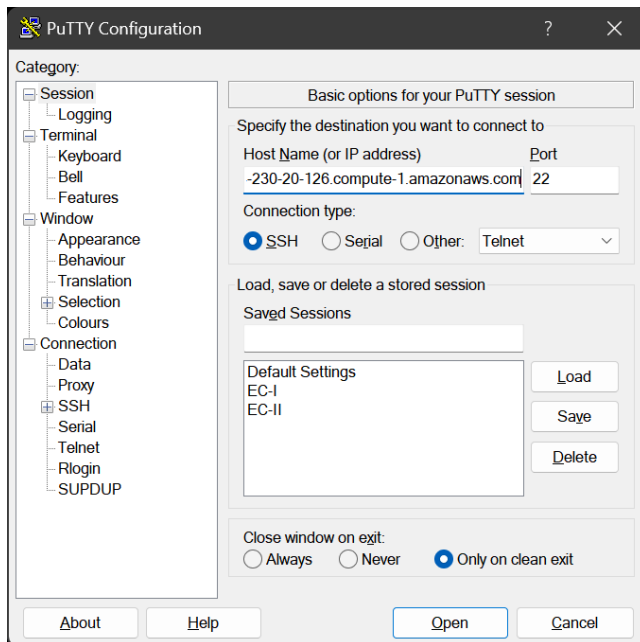   5.2. Copy the hostname mentioned under Windows tab

6. Open PuTTy
   6.1. Paste the copied hostname under *Session*
   6.2. Go to Connection >> SSH >> Auth >> Credentials
   6.3. Upload the keypair (.ppk)
   6.4. Launch the session

7. Run the commands below in the EMR session
   7.1. Switch to root user
      o sudo su
   7.2. Install Numpy
      o pip install numpy
   7.3. Sync all files in the instance
      o aws s3 sync s3://ua9-pa2-bucket/ .
      o You can verify the synced files using 'ls' command
   7.4. Run the training script
      o spark-submit ua9_training.py
   7.5. Copy the files in Hadoop File System
      o hadoop fs -copyFromLocal TrainingDataset.csv hdfs://ip-172-31-26-249.ec2.internal:8820/user/root/
      o hadoop fs -copyFromLocal ValidationDataset.csv hdfs://ip-172-31-26-249.ec2.internal:8820/user/root/
   7.6. Copy the model from HDFS to root directory
      o hadoop fs -ls hdfs://ip-172-31-26-249.ec2.internal:8820/user/root/
      o hadoop fs -get hdfs://ip-172-31-26-249.ec2.internal:8820/user/root/ua9-trainedmodel .
      o Check imported model in root directory using 'ls' command
   7.7. Run the model for prediction
      o spark-submit ua9_prediction.py
      o The test score and weighted F1 score is printed

```
Wine Prediction Model:
Test Accuracy = 0.96875
/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/context.py:158: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCre
ate() instead.
Prediction Model Weighted F1 Score = 0.9541901629072682
Exiting Spark Application
[root@ip-172-31-26-249 hadoop]# docker login
```

8. Running the prediction model using Docker

    8.1. Login into docker

        o   docker login

    8.2. Create the docker image using dockerfile

        o   docker build -t utakarshagg/winequalitytesting .



    8.3. Run the image

        o   docker run utakarshagg/winequalitytesting

    8.4. Push the image in docker hub

        o   docker push utakarshagg/winequalitytesting