

## CS643 Programming Assignment 2

Rudolph Paulin

12/11/2024

Github Link: <https://github.com/Bonli9/CS643-Wine-Prediction.git>

DockerHub Link: <https://hub.docker.com/repository/docker/rpaulin/wine-quality-prediction>

### Training Setup

Go to Amazon Console Home, Select S3 (Scalable Storage in the cloud) and click on "Create Bucket." Name the bucket "cs643rudolphpaulin," and uncheck the "Block all public access" option, keep the default settings for all other options and create the bucket. After the bucket is created, click on it and hit **Upload** to select "TrainingDataset.csv" and "ValidationDataset.csv" files.

### AWS EMR

Go to EMR on EC2 → Create Cluster → Give the cluster a name ("rp832-wineq-prediction") → Choose emr-7.5.0(Amazon EMR release). Then select the following applications: Zeppelin 0.11.1 ; Pig 0.17.0 ; and leave the other applications selected by default, as shown in the image below:

Name  
rp832-wineq-prediction

Amazon EMR release [Info](#)  
A release contains a set of applications which can be installed on your cluster.  
emr-7.5.0

Application bundle

Spark Interactive	Core Hadoop	Flink	HBase	Presto	Trino	Custom
<input type="checkbox"/> AmazonCloudWatchAgent 1.300032.2	<input type="checkbox"/> HCatalog 3.1.3	<input type="checkbox"/> Flink 1.19.1	<input type="checkbox"/> HBase 2.5.10	<input type="checkbox"/> Presto 0.287	<input type="checkbox"/> Trino 446	
<input type="checkbox"/> Hue 4.11.0	<input checked="" type="checkbox"/> Livy 0.8.0	<input checked="" type="checkbox"/> Hadoop 3.4.0	<input checked="" type="checkbox"/> Hive 3.1.3	<input checked="" type="checkbox"/> JupyterEnterpriseGateway 2.6.0	<input checked="" type="checkbox"/> JupyterHub 1.5.0	
<input checked="" type="checkbox"/> Pig 0.17.0	<input type="checkbox"/> TensorFlow 2.16.1	<input type="checkbox"/> Oozie 5.2.1	<input checked="" type="checkbox"/> Phoenix 5.2.0	<input checked="" type="checkbox"/> Spark 3.5.2	<input checked="" type="checkbox"/> Trino 446	
<input checked="" type="checkbox"/> Zeppelin 0.11.1	<input type="checkbox"/> ZooKeeper 3.9.2					

Choose a primary instance type m5.xlarge and 3 core groups. Then scroll down to select by defaults:

EMR\_DefaultRole→EMR\_EC2\_DefaultRole→EMR\_AutoScaling\_DefaultRole

Then hit "create cluster" button

While waiting for the cluster to initialize, go to the left panel and click on "Workspaces (Notebooks)" then select "Create Notebook." Name the notebook "wine-training." Under "Cluster," click "Choose" and pick the cluster that was set up earlier ("CS643rudolphpaulin"). Leave the other settings as defaults and click "Create Workspace" and name it **wine-training**. Wait for the status to update from "Starting" to "Ready." Once ready, click on the workspace and Jupyter application will open. In Jupyter, select the "wine-training.ipynb" notebook, then from the toolbar, click on Kernel > Change Kernel > PySpark.

Download the code from "training.py" (see Github repository) and paste it into the first cell of the notebook. After that, click "Run." Once executed, the results for both the LogisticRegression and RandomForestClassifier models will appear, showing the training outputs from all four nodes.

The screenshot shows the AWS EMR Studio interface. On the left, the 'Compute' panel is visible, showing the 'EMR on EC2 cluster' selected. The main area displays the 'wine-training.ipynb' notebook. The notebook content includes:

```
Data loaded from S3 bucket.
-----fixed acidity----- quality-----
0      8.9 ...      6
1      7.6 ...      5
2      7.9 ...      5
3      8.5 ...      5
4      6.9 ...      6

[5 rows x 12 columns]
Data has been formatted.
fixed acidity  volatile acidity  citric acid  ...  sulphates  alcohol  label
0      8.9      0.22      0.48  ...      0.53      9.4      6
1      7.6      0.39      0.31  ...      0.65      9.7      5
2      7.9      0.43      0.21  ...      0.91      9.5      5
3      8.5      0.49      0.11  ...      0.53      9.4      5
4      6.9      0.40      0.14  ...      0.63      9.7      6

[5 rows x 12 columns]
F1 Score for LogisticRegression Model: 0.5729445029855991
F1 Score for RandomForestClassifier Model: 0.5149515912576688
Since the Logistic Regression model has the superior F1 score, it will be selected for the prediction application.
```

## Prediction Configuration

In the AWS Console, go to Services > EC2 > Launch Instances. Choose "Amazon Linux 2 AMI..." Select the t2.large instance. Leave all other settings as default.

Generate a new key pair and name it "pa2.pem." Click "Download key pair," then select "Launch Instances." After that, click "View Instances." Initially, the EC2 instance status will likely display as "Pending." While it transitions to "Running," open a terminal and move the downloaded .pem file to your home directory. Next, execute the following command to adjust the file permissions for the .pem file:

- First, set the correct permissions for your .pem file by running the following command:

```
$ chmod 400 pa2.pem
```

- Once the EC2 instance is running, connect to it using the following SSH command:  
\$ ssh -i ~/pa2.pem ec2-user@<YOUR\_INSTANCE\_PUBLIC\_DNS>
- Next, open a terminal on your local machine to copy the data files (TrainingDataset.csv and TestDataset.csv):

```
$ scp -i ~/pa2.pem TrainingDataset.csv ec2user@<YOUR_INSTANCE_PUBLIC_DNS>:~/
```

```
$ scp -i ~/pa2.pem TestDataset.csv ec2-user@<YOUR_INSTANCE_PUBLIC_DNS>:~/
```

- Afterward, SSH into your EC2 instance again. You should see both data files in the home directory. Move them to the appropriate folder on your EC2 instance by running:

```
$ sudo mkdir /app
```

```
$ sudo cp TrainingDataset.csv TestDataset.csv /wineapp/
```

### Without Docker

After connecting to your EC2 instance via SSH, use the following commands to install and configure the Java Development Kit (JDK):

```
$ sudo yum install java-17-amazon-corretto
```

```
$ export JAVA_HOME=/usr/lib/jvm/java-17-amazon-corretto
```

To prepare your environment for Spark and Python, perform these steps:

1. Fetch Apache Spark 3.5.3 and unpack it into a dedicated directory:

```
$ wget https://archive.apache.org/dist/spark/spark-3.5.3/spark-3.5.3-bin-hadoop3.tgz -P ~/server
```

```
$ cd ~/server
```

```
$ sudo tar -xvzf spark-3.5.3-bin-hadoop3.tgz
```

2. Download the Anaconda installer and execute the installation:

```
$ curl -O https://repo.anaconda.com/archive/Anaconda3-2023.11-Linux-x86_64.sh
```

```
$ mv Anaconda3-2023.11-Linux-x86_64.sh /tmp
```

```
$ cd /tmp
```

```
$ bash Anaconda3-2023.11-Linux-x86_64.sh
```

**Then set up the environment variables.**

Reconnect to your EC2 instance and execute these commands to complete the setup:

1. Install additional Python dependencies:

```
$ pip install quinn
```

2. Set a password for Jupyter Notebook:

```
$ jupyter notebook password
```

3. Start the Jupyter Notebook server:

```
$ jupyter notebook
```

---

## Access Jupyter Notebook on Your Local Machine

To access the notebook from your browser:

1. Open a terminal on your local machine and create an SSH tunnel to forward traffic from your EC2 instance to your local machine:

```
$ ssh -i ~/pa2.pem -N -f -L localhost:8888:localhost:8888 ec2-user@<YOUR_INSTANCE_PUBLIC_DNS>
```

2. Open your web browser and go to <http://localhost:8888>.
3. Log in using the password you set during the configuration process.

See the image below



Start the Docker service: `$ sudo service docker start`

```
Login Succeeded
[ec2-user@ip-172-31-28-215 ~]$ docker tag wine-quality-prediction rpaulin/wine-quality-prediction:latest
[ec2-user@ip-172-31-28-215 ~]$ docker push rpaulin/wine-quality-prediction:latest
The push refers to repository [docker.io/rpaulin/wine-quality-prediction]
f22c2baf81f9: Pushed
e77f0fdcb0dd: Pushed
6ee9de778e45: Pushed
c50473b30ba4: Pushed
2cb18b10fe84: Pushed
16ebdaf10048: Pushing [=====>] 1.566GB
```