

Cloud Computing - Programming Assignment 2

Chaitanya Karnati

UCID: ck338

Wine Quality Prediction AWS Spark Application:

pa2assignment:

This project entails creating a Python application that leverages the PySpark interface. The application is deployed on an Amazon Web Services (AWS) Elastic MapReduce (EMR) cluster. The main goal is to parallelize the training of a machine learning model on EC2 instances to predict wine quality using publicly accessible data. Following the training phase, the model is then used to predict the quality of wine. Docker is employed to produce a container image for the trained machine learning model, simplifying the deployment process.

GitHub Link: <https://github.com/ChaitanyaKarnati/Cloud-Computing-PA2>

Docker Link: <https://hub.docker.com/repository/docker/chaitanyakarnati/winequlpred/general>

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

1. Generate an EMR Cluster Key Pair:
 - Navigate to EC2/Network/Key-pairs in the AWS console.
 - Use the .pem format and download the key pair.
 - Example key pair: pa2assignment.pem
2. Set Up an S3 Bucket:
 - Create an S3 bucket in AWS, e.g., pa2assignwineprediction
3. EMR Cluster Creation:
 - Access the EMR console.
 - Create a new EMR cluster.
4. Spark Cluster Configuration:
 - Configure the Spark cluster using the EMR console.
 - Specify the cluster details, including:
 - Name: chaitanyapa2
 - Amazon EMR release: emr-5.33.0
 - Application bundle: Hadoop 2.10.1, Spark 2.4.7, Zippeline 0.9.0, and Yarn.
5. Instance Configuration:
 - Create four instances for the Spark cluster.

These steps collectively establish the necessary infrastructure and configurations for running the AWS Spark application dedicated to predicting wine quality.

6. Currently, we are concurrently training an ML model on a Spark cluster using EC2 instances. Subsequently, the cluster is prepared to execute tasks related to running the ML model. To establish a connection with the Master instance through the Terminal:
 - **ssh -i "pa2assignment.pem" ec2-user@ec2-54-242-27-160.compute-1.amazonaws.com**
7. Upon successful login, After accessing the Master instance, switch to the root user by executing:
 - **sudo su**
8. Submit the task by the command:
 - **spark-submit s3://pa2assignwineprediction/wine_prediction.py**
9. Subsequently, check the trace status for the aforementioned tasks. If the status indicates success, it implies the successful creation of the "test.model" in the S3 bucket at the location **s3://pa2assignwineprediction**.

Docker steps:

1. Presently, we are executing the ML model using Docker:
 - Begin by creating a Docker account and completing the signup process.
 - Following a successful login, proceed to download and set up Docker on your local system.
 - Install Docker on your system.
 - In the PowerShell, log in to Docker using the following command:

docker login
docker login -u chaitanyakarnati
pwd
 - After login you need to build the image: **docker build -t winequlpre .**
2. The push and pull into the docker hub repository:
 - PUSH:

docker tag winequlpre chaitanyakarnati/winequlpre
docker push chaitanyakarnati/winequlpre
 - PULL:

docker pull chaitanyakarnati/winequlpre
3. Place your test data file in a specified folder named "dir." Mount this directory to the Docker container and run the container using the following command:
 - **docker run -v /Users/chaitanyakarnati/Documents/Wineprediction/data/csv winequlpred testdata.csv**

Running the machine learning model without utilizing Docker:

1. Clone the Repository:

Begin by cloning this repository onto your local machine.

2. Set Up Local Spark Environment:

Ensure that you have a local Spark environment ready for running this application. If you haven't set up Spark yet, you can follow the instructions provided in the [official Spark documentation](<https://spark.apache.org/docs/latest>).

3. Navigate to the Python File Folder:

Access the 'python file' directory within the cloned repository.

4. Prepare Test Data:

Store your test data in the '/Users/chaitanyakarnati/Documents/Wineprediction/data/csv' folder.

Conclusion:

As depicted in the image below, an accuracy of approximately 98% was achieved while predicting wine quality.

```
2022/07/20 17:02 INFO SharedDataset: Warehouse path is 'file:/cdd00/spark-warehouse...'
---Input file for test data is---
data/csv/testdata.csv

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|fixed acidity|volatile acidity|citric acid|residual sugar|chlorides|free sulfur dioxide|total sulfur dioxide|density| pH|sulphates|alcohol|quality|      features|label|      rawPrediction|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      8.9|      0.22|      0.48|      1.8|      0.077|      29.0|      60.0| 0.9968|3.39|      0.53|      9.4|      6.0|[8.9,0.22,0.48,1,...]| 1.0|[3.48851027289842...]| [0.0697
7020545796...| 1.0|
|      7.6|      0.39|      0.31|      2.3|      0.082|      23.0|      71.0| 0.9982|3.52|      0.65|      9.7|      5.0|[7.6,0.39,0.31,2,...]| 0.0|[48.1243835079459...]| [0.9624
8767015891...| 0.0|
|      7.9|      0.43|      0.21|      1.6|      0.106|      10.0|      37.0| 0.9966|3.17|      0.91|      9.5|      5.0|[7.9,0.43,0.21,1,...]| 0.0|[48.1539002576703...]| [0.9630
7800515340...| 0.0|
|      8.5|      0.49|      0.11|      2.3|      0.084|      9.0|      67.0| 0.9968|3.17|      0.53|      9.4|      5.0|[8.5,0.49,0.11,2,...]| 0.0|[47.6785761357096...]| [0.9535
7152271419...| 0.0|
|      6.9|      0.4|      0.14|      2.4|      0.085|      21.0|      40.0| 0.9968|3.43|      0.63|      9.7|      6.0|[6.9,0.4,0.14,2.4,...]| 1.0|[1.82872254349015...]| [0.0365
7445086980...| 1.0|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

None
Test Accuracy of wine prediction model = 0.983580922595778
Weighted f1 score of wine prediction model = 0.9776578095108527
(base) chaitanyakarnati@Chaitanyas-MBP Wineprediction %
```