## Description

In this assignment, the task is to develop a parallel machine learning model in the Amazon AWS cloud environment. In this project, I have trained the "winequality" prediction model based on given training data and then utilized this model to test my model for the given validation dataset. I have utilized the Random Forest regression model for training and testing wine quality prediction applications and then calculated the accuracy of the model to determine how well the model performed. The task was to deploy our model with and without docker. The following are steps in the development of the Wine Quality prediction system.

Docker Link: <a href="https://hub.docker.com/r/system25/winequalityprediction">https://hub.docker.com/r/system25/winequalityprediction</a>

## How to set up the cloud environment and run the model training and the application prediction without docker

Following are the steps to a setup cloud environment:

- 1. Navigate to the Amazon AWS website i.e. "https://aws.amazon.com/
- 2. Create an Amazon AWS account and log in to the console by providing login credentials.
- 3. From the search for EC2 service and navigate to the EC2 dashboard.
- 4. Click on Launch Instance and follow all the steps by providing all the mandatory details.
- 5. Choose a number of instances. In this case, I have selected 5 instances because we are required to run our application on multiple instances.
- 6. Click on the launch button to launch your EC2 instance.
- 7. From the command line ssh into the ec2 machine from our local system
- 8. Copy application files and folders from the local system into the EC2 instance.
- 9. Run the following command to build the model
  - 1. Spark-submit test.py

# How to set up the cloud environment and run the model training and the application prediction with docker

Following are the steps to run the application in docker in Amazon AWS:

- Login to your Amazon AWS account.
- SSH into EC2 instance from the local machine.
- Install docker into the EC2 instance by using the following commands:
  - o sudo yum install docker
  - sudo service docker start
  - o sudo usermod -a -G docker ec2-user
- Navigate to the directory with application files.
- Run the following command to build the image and run the container:
  - docker build -t winequalityprediction
  - o docker run -p 80:80 winequalityprediction

### Screenshots:

#### Docker build in ec2 instance

```
ec2-user@ip-172-31-33-167:~/winequalityprediction
Adding debian:Hellenic_Academic_and_Research_Institutions_ECC_RootCA_2015.pem
Adding debian:Hellenic_Academic_and_Research_Institutions_RootCA_2011.pem
Adding debian:Hellenic_Academic_and_Research_Institutions_RootCA_2015.pem
Adding debian:Hongkong_Post_Root_CA_1.pem
Adding debian:Hongkong_Post_Root_CA_3.pem
Adding debian:ISRG_Root_X1.pem
Adding debian:IdenTrust_Commercial_Root_CA_1.pem
Adding debian:IdenTrust_Public_Sector_Root_CA_1.pem
Adding debian:Izenpe.com.pem
Adding debian:Microsec_e-Szigno_Root_CA_2009.pem
Adding debian:Microsoft_ECC_Root_Certificate_Authority_2017.pem
Adding debian:Microsoft_RSA_Root_Certificate_Authority_2017.pem
Adding debian:NAVER_Global_Root_Certification_Authority.pem
Adding debian:NetLock_Arany_=Class_Gold=_Főtanúsítvány.pem
Adding debian:Network_Solutions_Certificate_Authority.pem
Adding debian:OISTE_WISeKey_Global_Root_GB_CA.pem
Adding debian:OISTE_WISeKey_Global_Root_GC_CA.pem
Adding debian:QuoVadis_Root_CA.pem
Adding debian:QuoVadis_Root_CA_1_G3.pem
Adding debian:QuoVadis_Root_CA_2.pem
Adding debian:QuoVadis_Root_CA_2_G3.pem
Adding debian:QuoVadis_Root_CA_3.pem
Adding debian:QuoVadis_Root_CA_3_G3.pem
 Adding debian:SSL.com_EV_Root_Certification_Authority_ECC.pem
```

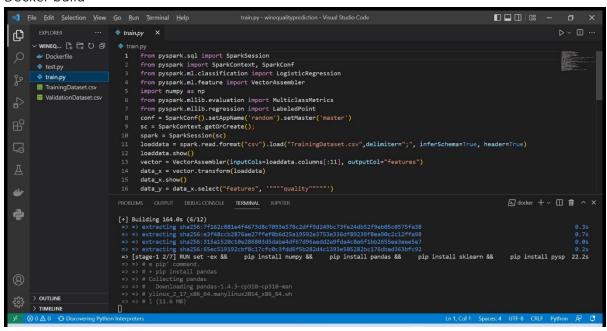
#### Docker Run

```
ec2-user@ip-172-31-38-174:~/winequalityprediction
            at org.sparkproject.guava.cache.LocalCache$Segment.loadSync(LocalCache.java:2379)
             .. 58 more
22/08/05 11:01:56 INFO FileScanRDD: Reading File path: file:///winequalitypredict/ValidationDataset.csv, range: 0-8762,
 partition values: [empty row]
22/08/05 11:01:56 INFO BlockManagerInfo: Removed broadcast_136_piece0 on a3524101733a:44055 in memory (size: 27.0 KiB,
 ree: 413.7 MiB)
22/08/05 11:01:56 INFO PythonRunner: Times: total = 210, boot = 9, init = 200, finish = 1
22/08/05 11:01:56 INFO Executor: Finished task 0.0 in stage 68.0 (TID 68). 2718 bytes result sent to driver
22/08/05 11:01:56 INFO TaskSetManager: Finished task 0.0 in stage 68.0 (TID 68) in 285 ms on a3524101733a (executor driv
22/08/05 11:01:56 INFO TaskSchedulerImpl: Removed TaskSet 68.0, whose tasks have all completed, from pool 22/08/05 11:01:56 INFO DAGScheduler: ResultStage 68 (toPandas at /test.py:30) finished in 0.295 s 22/08/05 11:01:56 INFO DAGScheduler: Job 68 is finished. Cancelling potential speculative or zombie tasks for this job 22/08/05 11:01:56 INFO TaskSchedulerImpl: Killing all running tasks in stage 68: Stage finished 22/08/05 11:01:56 INFO DAGScheduler: Job 68 finished: toPandas at /test.py:30, took 0.300492 s
               1:01:56 INFO SparkContext: Invoking stop() from shutdown hook
22/08/05 11:01:56 INFO SparkUI: Stopped Spark web UI at http://a3524101733a:4040
22/08/05 11:01:56 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
22/08/05 11:01:56 INFO MemoryStore: MemoryStore cleared
22/08/05 11:01:56 INFO BlockManager: BlockManager stopped
22/08/05 11:01:56 INFO BlockManagerMaster: BlockManagerMaster stopped
22/08/05 11:01:56 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
22/08/05 11:01:56 INFO SparkContext: Successfully stopped SparkContext
22/08/05 11:01:56 INFO ShutdownHookManager: Shutdown hook called
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-a43b3d49-371f-4c6b-99f5-06c180614b58
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-a43b3d49-371f-4c6b-99f5-06c180614b58/pyspark-8
3735332-3b37-4bf8-aa27-3ec37c8614e4
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-9f13ab34-0f04-47c3-a6fe-00e0cfbf2ee7
[ec2-user@ip-172-31-38-174 winequalityprediction]$
```

#### Docker run

```
ec2-user@ip-172-31-38-174:~/winequalityprediction
                                                                                                                                                                                        \Box
                                                                                                                                                                                                   ×
             at org.sparkproject.guava.cache.LocalCache$Segment.loadSync(LocalCache.java:2379)
             ... 58 more
22/08/05 11:01:56 INFO FileScanRDD: Reading File path: file:///winequalitypredict/ValidationDataset.csv, range: 0-8762,
partition values: [empty row]
22/08/05 11:01:56 INFO BlockManagerInfo: Removed broadcast_136_piece0 on a3524101733a:44055 in memory (size: 27.0 KiB, f
22/08/03 11:01:35 In 0 Block analyse Info. Nembers of State 210, 100 to 20, 101:01:35 Info Block analyse Info. Nembers 313.7 MiB)
22/08/05 11:01:56 INFO PythonRunner: Times: total = 210, boot = 9, init = 200, finish = 1
22/08/05 11:01:56 INFO Executor: Finished task 0.0 in stage 68.0 (TID 68). 2718 bytes result sent to driver
22/08/05 11:01:56 INFO TaskSetManager: Finished task 0.0 in stage 68.0 (TID 68) in 285 ms on a3524101733a (executor driv
 er) (1/1)
22/08/05 11:01:56 INFO TaskSchedulerImpl: Removed TaskSet 68.0, whose tasks have all completed, from pool
22/08/05 11:01:56 INFO DAGScheduler: ResultStage 68 (toPandas at /test.py:30) finished in 0.295 s
22/08/05 11:01:56 INFO DAGScheduler: Job 68 is finished. Cancelling potential speculative or zombie tasks for this job 22/08/05 11:01:56 INFO TaskSchedulerImpl: Killing all running tasks in stage 68: Stage finished 22/08/05 11:01:56 INFO DAGScheduler: Job 68 finished: toPandas at /test.py:30, took 0.300492 s
F1- score: 0.55625
22/08/05 11:01:56 INFO SparkContext: Invoking stop() from shutdown hook
22/08/05 11:01:56 INFO SparkUI: Stopped Spark web UI at http://a3524101733a:4040
22/08/05 11:01:56 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
22/08/05 11:01:56 INFO MemoryStore: MemoryStore cleared
22/08/05 11:01:56 INFO BlockManager: BlockManager stopped
22/08/05 11:01:56 INFO BlockManagerMaster: BlockManagerMaster stopped
22/08/05 11:01:56 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
22/08/05 11:01:56 INFO SparkContext: Successfully stopped SparkContext
22/08/05 11:01:56 INFO ShutdownHookManager: Shutdown hook called
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-a43b3d49-371f-4c6b-99f5-06c180614b58
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-a43b3d49-371f-4c6b-99f5-06c180614b58/pyspark-8
3735332-3b37-4bf8-aa27-3ec37c8614e4
22/08/05 11:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-9f13ab34-0f04-47c3-a6fe-00e0cfbf2ee7
[ec2-user@ip-172-31-38-174 winequalityprediction]$
```

#### Docker build



## **Image Creation**

### Model

8/04 21:16:49 INFO CodeGenerator: Code generated in 43.385794 ms									
	idity""" """v "pH""" """sul	phates""" "	"""alcoho	1"""" """	"quality""			fur dioxide""" """total sulfu	r dioxide"""
	7.4		0.7		0.0	1.9	+ 0.076	11	34
.9978	3.51 7.8	0.56	0.88	9.4	0.0	5 [7.4,0.7,0.0,1.9,  2.6	0.098	25	67
.9968	3.2	0.68	0.761	9.8	0.04	5 [7.8,0.88,0.0,2.6		451	Eal
0.997	7.8  3.26	0.65	0.76	9.8	0.04	2.3  5 [7.8,0.76,0.04,2	0.092  	15	54
0.998	11.2  3.16	0.58	0.28	9.8	0.56	1.9  6 [11.2,0.28,0.56,1	0.075	17	69
0.556	7.4	0.56	0.7	9.0	0.0	1.9	0.076	11	34
.9978	3.51 7.4	0.56	0.66	9.4	0.0	5 [7.4,0.7,0.0,1.9,  1.8	0.075	13	40
.9978	3.51	0.56		9.4	0.01	5 [7.4,0.66,0.0,1.8			401
.9964	7.9  3.3	0.46	0.6	9.4	0.06	1.6  5 [7.9,0.6,0.06,1.6	0.069	15	59
	7.3		0.65		0.0	1.2	0.065	15	21
.9946	3.39 7.8	0.47	0.58	10.0	0.02	7 [7.3,0.65,0.0,1.2	0.073	9	18
.9968	3.36	0.57		9.5		7 [7.8,0.58,0.02,2			
.9978	7.5  3.35	0.8	0.5	10.5	0.36	6.1  5 [7.5,0.5,0.36,6.1	0.071  	17	102
	6.7		0.58		0.08	1.8	0.097	15	65
.9959	3.28 7.5	0.54	0.5	9.2	0.36	5 [6.7,0.58,0.08,1  6.1	0.071	17	102
.9978	3.35	0.8	0 6451	10.5	٠ م ا	5 [7.5,0.5,0.36,6.1			Fol
.9943	5.6  3.58	0.52	0.615	9.9	0.0	1.6    5 [5.6,0.615,0.0,1	0.089  	16	59
.9974	7.8  3.26	1.56	0.61	9.1	0.29	1.6  5 [7.8,0.61,0.29,1	0.114	9	29
	8.9		0.62		0.18	3.8	0.176	52	145
.9986	3.16 8 9	0.88	9 62	9.2	a 19l	5 [8.9,0.62,0.18,3		51	148
.9986	8.9  3.17	0.93	0.62	9.2	0.19	3.9    5 [8.9,0.62,0.19,3	0.17  	51	