

# CS643 - Programming Assignment 2

Docker link :

[https://hub.docker.com/repository/docker/mg946/wine\\_quality/general](https://hub.docker.com/repository/docker/mg946/wine_quality/general)

Github Repository link :

[https://github.com/MohanKumar946/Wine\\_Quality\\_CS643](https://github.com/MohanKumar946/Wine_Quality_CS643)

## Wine Quality Prediction Using Spark and Amazon Web Services (AWS)

- The procedure to use AWS services to train ML (Machine Learning) models on multiple parallel EC2(Elastic Compute Cloud) instances using EMR cluster.
- The ML program is written in python using Apache Spark MLlib libraries. The training and prediction programs are configured to run inside a container.
- First we set up an EMR on the AWS academy account. We can give any name to our EMR and add a total of 4 tasks on the EMR. Create a new key value pair and download the .pem file.
- Add appropriate permissions required for EMR,EC2 and autoscaling. Create the EMR cluster.

- Now start running the EMR cluster.
- Check all the EC2 instances that are running in AWS account. There will be 6 EC2 instances running. There will be one master EC2 instance and rest of them will be slave EC2 instances.
- If now if not already added add a new SSH inbound rule in the security of the master instance and add your IP address

Now create a new S3 bucket where we can add all our files that have to run on these EC2 instances.

- Create a new S3 bucket and give a unique name to it.
- Now upload all the files we need to the S3 bucket.
- Upload [training.py](#) , [prediction.py](#) , [TrainingDataset.csv](#) , [ValidationDataset.csv](#)
- Make sure the files are visible in the S3 bucket after they are uploaded.

Go to the EMR cluster created :

- Check for the option Connect to the Primary node using SSH
- We will get a SSH connection link for terminal

**ssh -i ~/EMR.pem**

[\*\*hadoop@ec2-34-202-160-201.compute-1.amazonaws.com\*\*](#)

- Open a terminal on local where the .pem file is present. Before running the above command we need to change permissions of the .pem file

- Run

**chmod 400 ./EMR.pem**

- Now run the SSH connection command

```
mohankumargm@Mohans-MBP IS601 TA % chmod 400 EMR.pem
mohankumargm@Mohans-MBP IS601 TA % ssh -i ~/EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com
Warning: Identity file /Users/mohankumargm/EMR.pem not accessible: No such file or directory.
The authenticity of host 'ec2-34-202-160-201.compute-1.amazonaws.com (34.202.160.201)' can't be established.
ED25519 key fingerprint is SHA256:gDFx3325+IRAKJ1CKtJiYnSyGvto8L1WvyoPWVa6r0c.
This key is not known by any other names.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-34-202-160-201.compute-1.amazonaws.com' (ED25519) to the list of known hosts.
hadoop@ec2-34-202-160-201.compute-1.amazonaws.com: Permission denied (publickey,gssapi-keyex,gssapi-with-mic).
mohankumargm@Mohans-MBP IS601 TA % ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com
#_
~\_ ##### Amazon Linux 2023
~~ \####|
~~ \|#
[ ~~ \#/ __ https://aws.amazon.com/linux/amazon-linux-2023
~~ V~' '-->
~~ /
~~ .-
~/-
~/-
~/m/`
```

```
EEEEEEEEEEEEEEEEEE MMMMMMM MBBBBBB RRRRRRRRRRRRRR
E:::::::::::E M:::::M M:::::M R:::::R
EE:::::EEEEE:::E M:::::M M:::::M R:::::RRRRRR:::::R
E:::E EEEEE M:::::M M:::::M RR:::R R:::::R
E:::E M:::::M::::M M:::::M R:::R R:::::R
E:::::EEEEE M:::::M M:::::M M:::::M R:::::RRRRRR:::::R
E:::::::::::E M:::::M M:::::M M:::::M R:::::RR
E:::::EEEEE M:::::M M:::::M M:::::M R:::RRRRR:::::R
E:::E M:::::M M:::::M M:::::M R:::R R:::::R
E:::E EEEEE M:::::M MMM M:::::M R:::R R:::::R
EE:::::EEEEE:::E M:::::M M:::::M R:::R R:::::R
E:::::::::::E M:::::M M:::::M RR:::R R:::::R
EEEEEEEEEEEEEEEEEE MMMMMMM MBBBBBB RRRRRRR
```

- Successful connection and the hadoop terminal will be connected
- Now connect the S3 bucket with the EMR cluster in the terminal to access its files.

**aws s3 sync s3://emr-ml-bucket/ .**

- 'ls' command and check if all the files are present.

Setup all python dependencies to run the python files :

- Install numpy dependency

**pip install numpy --user**

- Now we have to copy the files to Hadoop file system from the local for both Training.py and Validation.py

**hadoop fs -copyFromLocal TrainingDataset.csv  
hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/**

**hadoop]# hadoop fs -copyFromLocal ValidationDataset.csv  
hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/**

Once all the dependencies are installed we can run spark submit

- Run the command :

**spark-submit Training.py**

```
[root@ip-172-31-26-205 ~]# hadoop fs -copyFromLocal TrainingDataset.csv hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/
[root@ip-172-31-26-205 ~]# hadoop fs -copyFromLocal ValidationDataset.csv hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/
[root@ip-172-31-26-205 ~]# spark-submit Training.py
Apr 29, 2024 12:52:18 AM org.apache.spark.launcher.Log4jHotPatchOption staticJavaAgentOption
WARNING: spark.log4jHotPatch.enabled is set to true, but /usr/share/log4j-cve-2021-44228-hotpatch/jdk17/Log4jHotPatchFat.jar does not exist at the configured location

^[[A/usr/bin/python3: can't open file '/home/hadoop/Training.py': [Errno 2] No such file or directory
24/04/29 09:52:12 INFO ShutdownHookManager: Shutdown hook called
24/04/29 09:52:12 INFO ShutdownHookManager: Deleting directory /mnt/tmp/spark-fe1977e3-9680-4226-825b-dbe8ca82469
[root@ip-172-31-26-205 ~]# spark-submit Training.py
Apr 29, 2024 12:52:23 AM org.apache.spark.launcher.Log4jHotPatchOption staticJavaAgentOption
WARNING: spark.log4jHotPatch.enabled is set to true, but /usr/share/log4j-cve-2021-44228-hotpatch/jdk17/Log4jHotPatchFat.jar does not exist at the configured location

Starting Spark Application in EMR
24/04/29 09:52:26 INFO SparkContext: Running Spark version 3.5.0-amzn-1
24/04/29 09:52:26 INFO SparkContext: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 09:52:26 INFO SparkContext: Java version 17.0.10
=====
24/04/29 09:52:26 INFO ResourceUtils: =====
24/04/29 09:52:26 INFO ResourceUtils: No custom resources configured for spark.driver.
24/04/29 09:52:26 INFO ResourceUtils: =====
24/04/29 09:52:26 INFO SparkContext: Submitted application: Wine-Quality-Prediction-SPARK-ML
24/04/29 09:52:26 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 4, script: , vendor: , memory -> name: memory, amount: 9486, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpus -> name: cpus, amount: 1.0)
24/04/29 09:52:26 INFO ResourceProfile: limiting resource is cpus at 4 tasks per executor
24/04/29 09:52:26 INFO SecurityManager: Setting spark.driver.securityManager.principalId: 0
24/04/29 09:52:26 INFO SecurityManager: Changing view acls to: root
24/04/29 09:52:26 INFO SecurityManager: Changing view acl groups to: root
24/04/29 09:52:26 INFO SecurityManager: Changing modify acl groups to:
24/04/29 09:52:26 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: root; groups with view permissions: EMPTY; users with modify permissions: root; groups with modify permissions: EMPTY
24/04/29 09:52:27 INFO Utils: Successfully started service 'sparkDriver' on port 37565.
24/04/29 09:52:27 INFO SparkEnv: Registering MapOutputTracker
```

```

24/04/29 00:52:38 INFO ServerInfo: Adding filter to /spis: org.apache.hadoop.yarn.server.webproxy.ammfilter.AmIpFilter
24/04/29 00:52:38 INFO ServerInfo: Adding filter to /jobs/job/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:52:38 INFO ServerInfo: Adding filter to /stages/stage/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:52:38 INFO ServerInfo: Adding filter to /metrics/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:52:38 INFO YarnClientSchedulerBackend: SchedulerBackend is ready for scheduling beginning after reached minRegisteredResourcesRatio: 0.0
Reading training CSV file from TrainingDataset.csv
Creating VectorAssembler
Creating StringIndexer
Caching data for faster access
Creating RandomForestClassifier
Creating Pipeline for training
Retraining model on multiple parameters using CrossValidator
Fitting CrossValidator to the training data
Saving the best model to new param 'model'
Saving the best model to S3

```

- Once the training.py is completed and a trained model is stored in the S3 bucket. Use ‘ls’ command and find the trained model

```

72-31-26-205 hadoop]# ls
TrainingDataset.csv ValidationDataset.csv prediction.py trainedmodel training.py
● 72-31-26-205 hadoop]# spark-submit prediction.py

```

- Now run spark command for predicting the data from prediction.py

## spark-submit prediction.py

```

ApplicationMaster host: N/A
ApplicationMaster RPC port: -1
queue: default
start time: 1714352352256
final status: UNDEFINED
tracking URL: http://ip-172-31-26-205.ec2.internal:20888/proxy/application_1714350124576_0004/
user: root
24/04/29 00:59:16 INFO YarnClientSchedulerBackend: Add WebUI Filter. org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter, Map[PROXY_HOSTS -> ip-172-31-26-205.ec2.internal, PROXY_PORTS -> 20888], org.apache.hadoop.yarn.server.webproxy.ammfilter.AmIpFilter
24/04/29 00:59:17 INFO YarnSchedulerBackend$YarnSchedulerEndpoint: ApplicationMaster registered as NettyRpcEndpointRef(spark-client://YarnAM)
24/04/29 00:59:17 INFO Client: Application report for application_1714350124576_0004 (state: RUNNING)
24/04/29 00:59:17 INFO Client:
  client token: N/A
  diagnostics: N/A
  ApplicationMaster host: 172.31.23.66
  ApplicationMaster RPC port: -1
  queue: default
  start time: 1714352352256
  final status: UNDEFINED
  tracking URL: http://ip-172-31-26-205.ec2.internal:20888/proxy/application_1714350124576_0004/
  user: root
24/04/29 00:59:17 INFO YarnClientSchedulerBackend: Application application_1714350124576_0004 has started running.
24/04/29 00:59:17 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 37383.
24/04/29 00:59:17 INFO NettyBlockTransferService: Server created on ip-172-31-26-205.ec2.internal:37383
24/04/29 00:59:17 INFO BlockManager: External shuffle service port = 7337
24/04/29 00:59:17 INFO BlockManagerMaster: Registering BlockManagerId(driver, ip-172-31-26-205.ec2.internal:37383, None)
24/04/29 00:59:17 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383, None)
24/04/29 00:59:17 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383, None)
24/04/29 00:59:17 INFO SingleEventLogFileWriter: Logging events to file: /var/log/spark/apps/application_1714350124576_0004.inprogress
24/04/29 00:59:17 INFO Utils: Using 50 preallocated executors (minExecutors: 0). See spark.dynamicAllocation.preallocateExecutors to 'false' disable executor preallocation.
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/job/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/stage/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /environment/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage/rdd/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /static: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /api: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/job/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/stage/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /metrics/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO YarnClientSchedulerBackend: SchedulerBackend is ready for scheduling beginning after reached minRegisteredResourcesRatio: 0.0
Test Accuracy of wine prediction model = 0.96875
/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/context.py:158: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Weighted F1 Score of wine prediction model = 0.9541901629872682

```

- We can see a accuracies shown :

**Test Accuracy of wine prediction model = 0.96875**  
**Weighted F1 Score of wine prediction model =**  
**0.9541901629072682**

Once our prediction is done we can make our Docker image.

- First create a docker account
- Now in the same root terminal use command :

## docker login

- Enter your docker Username and Password credentials.
- Now we run a command to build the docker file

## docker build -t mg946/wine\_quality .

```
[root@mg946-172-31-26-285 hadoop]# docker login
Log in with your Docker ID or email address to push and pull images from Docker Hub. If you don't have a Docker ID, head over to https://hub.docker.com/ to create one.
You can log in with your password or a Personal Access Token (PAT). Using a limited-scope PAT grants better security and is required for organizations using SSD. Learn more at https://docs.docker.com/go/acces
kens/
Username: mg946@mjt.edu
Password:
WARNING! Your password will be stored unencrypted in /root/.dockercfg.json.
Configure a credential helper to remove this warning. See
https://docs.docker.com/engine/reference/commandline/login/#credentials-store

Login Succeeded
[root@mg946-172-31-26-285 hadoop]# docker build -t mg946/wine_quality .
[+] Building 95.3s (17/17) FINISHED                                            docker:default
--> [internal] load build definition from Dockerfile                         0.0s
--> [internal] load metadata for docker.io/library/openjdk:8-jre-slim          0.0s
--> [internal] load metadata for docker.io/library/openjdk:8-jre-slim           0.4s
--> [auth] library/openjdk:pull token for registry-1.docker.io                 0.0s
--> [internal] load .dockercfg                                                 0.0s
--> [internal] load index manifest                                             0.0s
--> [ 1/11] FROM docker.io/library/openjdk:8-jre-slim@sha256:53186129237fb8bc0a12dd36da6761f4c7a2e20233c20d4eb8d497e4045a4f5 2.8s
--> [ 2/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 3/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 4/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 5/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 6/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 7/11] RUN curl -sS https://repo.maven.apache.org/maven2/com/amazonaws/aws-sdk-java/1.11.500/aws-xml-databind-1.11.500.jar > aws-xml-databind-1.11.500.jar      0.0s
--> [ 8/11] COPY ValidationDataset.csv /opt/                                      0.0s
--> [ 9/11] COPY trainedmodel.pkl /opt/trainedmodel/                           0.0s
--> [10/11] COPY ValidationDataset.csv /opt/                                      0.0s
--> [11/11] COPY ValidationDataset.csv /opt/                                      0.0s
--> exporting to image                                                       4.5s
--> exporting layers                                                        4.5s
--> exporting annotations                                                    0.0s
--> naming to docker.io/mg946/wine_quality                                     0.0s
--> rm -rf /tmp/miniconda.sh /tmp/miniconda.sh                                0.0s
```

- We check if the Docker file has been created properly by running the file :

## docker run mg946/wine\_quality

```
IS601TA ~ root@ip-172-31-26-205:/home/hadoop ~ ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com ~ 212x67
  .-- naming to docker.io/mg946/wine_quality
[root@ip-172-31-26-205 ~]# docker run mg946/wine_quality
/bin/bash: /bin/ls: command not found
Starting Spark Application
24/04/29 01:11:13 INFO SparkContext: Running Spark version 3.5.0
24/04/29 01:11:13 INFO SparkContext: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:13 INFO SparkContext: Java version 1.8.0_342
24/04/29 01:11:13 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO ResourceUtils: custom resources configured for spark.driver.
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO SparkContext: Submitted application: Wine-Quality-Prediction-SPARK-ML
24/04/29 01:11:13 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map[cores -> name: cores, amount: 1, script: , vendor: , memory -> name: memory, amount: 1024, script: , vendor: , offHeap -> name: offHeap, amount: 1, vendor: ] task resources: Map[cpus -> name: cpus, amount: 1, vendor: ]
24/04/29 01:11:13 INFO ResourceProfile: Limiting resource is cpu
24/04/29 01:11:13 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/04/29 01:11:13 INFO SecurityManager: Changing view acls to: root
24/04/29 01:11:13 INFO SecurityManager: Changing modify acls to: root
24/04/29 01:11:13 INFO SecurityManager: Changing view acls groups to:
24/04/29 01:11:13 INFO SecurityManager: Changing modify acls groups to:
24/04/29 01:11:13 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: root; groups with view permissions: EMPTY; users with modify permissions: root; groups with modify permissions: root; groups with modify permissions: EMPTY
24/04/29 01:11:14 INFO Utils: Successfully started service 'sparkDriver' on port 37995.
24/04/29 01:11:14 INFO SparkEnv: Registering MapOutputTracker
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMasterEndpoint
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/04/29 01:11:14 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-cde8d74c-cfce-47af-a347-88b8e4c0940
24/04/29 01:11:14 INFO MemoryStore: MemoryStore started with capacity 366.3 MB
24/04/29 01:11:14 INFO SparkEnv: Registering OutputCommitter
24/04/29 01:11:14 INFO SparkEnv: Registering ExecutorMetricsReporter
24/04/29 01:11:14 INFO SparkUI: Starting SparkUI on port 4040.
24/04/29 01:11:14 INFO Executor: Starting executor ID driver on host ee5c07e4c40f
24/04/29 01:11:14 INFO Executor: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:14 INFO Executor: Java version 1.8.0_342
24/04/29 01:11:14 INFO Executor: Registered executor with 0.0 GB free memory
24/04/29 01:11:14 INFO Utils: Successfully started service 'SparkUI' on port 4040.
24/04/29 01:11:14 INFO Executor: Created & updated repl class loader org.apache.spark.util.MutableURLClassLoader@79116642d for default.
24/04/29 01:11:14 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 39683.
24/04/29 01:11:14 INFO BlockManager: Using org.apache.spark.storage.RandomizedBlockReplicationPolicy for block replication policy
24/04/29 01:11:14 INFO BlockManagerMaster: Registering block manager ee5c07e4c40f:39683 with 0.0 MB RAM, BlockManagerId(driver, ee5c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: Registering block manager ee5c07e4c40f:39683 with 0.0 MB RAM, BlockManagerId(driver, ee5c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, ee5c07e4c40f, 39683, None)
Test Accuracy of wine prediction model = 0.96875
/opt/spark/python/lib/pyspark.zip/pyspark/context.py:158: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Weighted F1 Score of wine prediction model = 0.9641981629872682
Exiting Spark Application
```

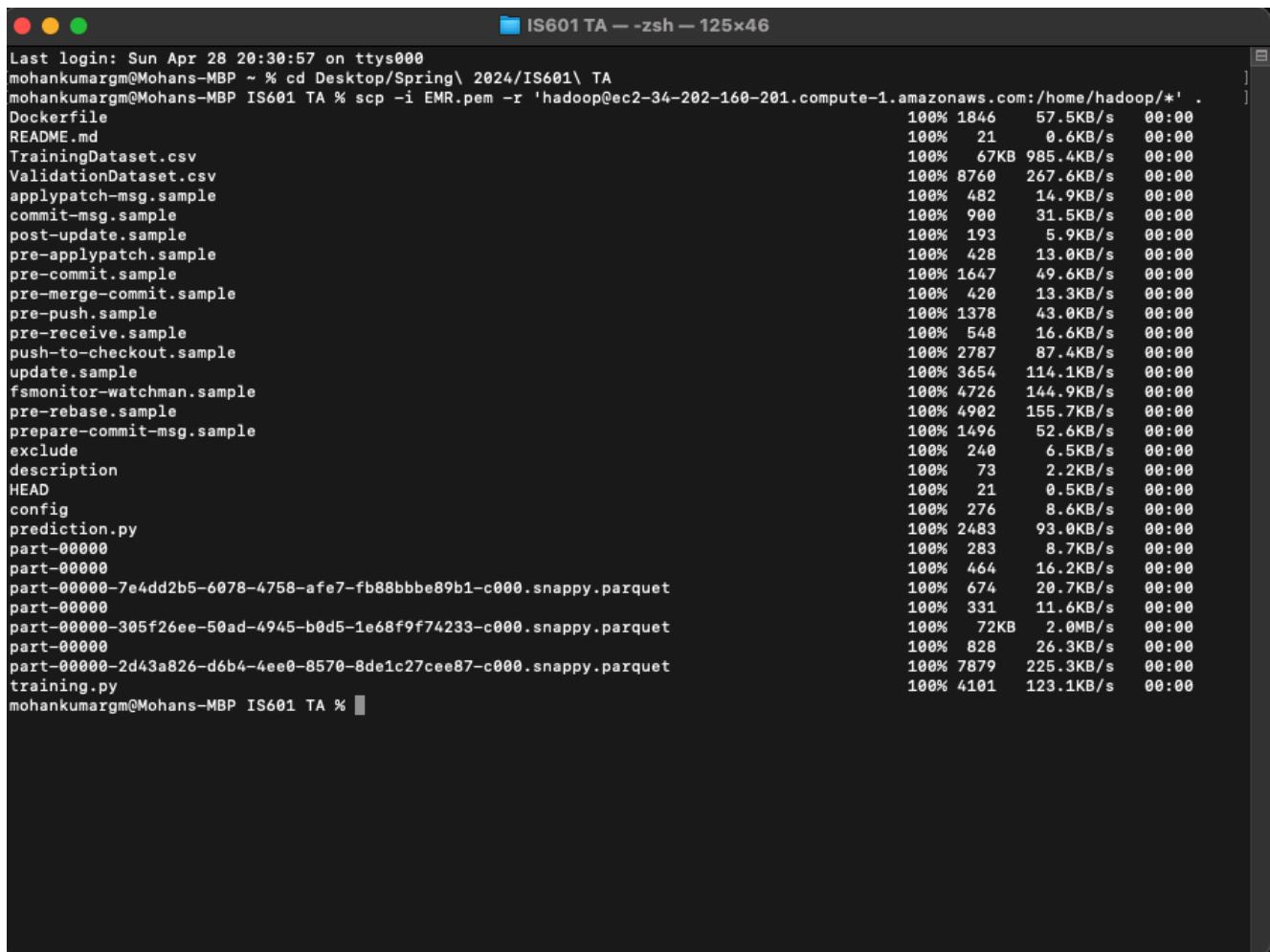
- Once verified push the docker image

## **docker push mg946/wine\_quality**

```
[root@ip-172-31-26-205 hadoop]# docker push mg946/wine_quality
Using default tag: latest
The push refers to repository [docker.io/mg946/wine_quality]
28c3e7ac5714: Pushed
9f9d12ca39f9: Pushed
06c4012ca3b1: Pushed
b72c5a50bd95: Pushed
f64b460f220c: Pushed
3563f7b387d3: Pushed
5f70bf18a086: Pushed
4bf6a6a5ac55: Pushed
c28332e68d2f: Pushed
7aa75242fa43: Pushed
b66078cf4b41: Mounted from library/openjdk
cd5a0a9f1e01: Mounted from library/openjdk
eafe6e032dbd: Mounted from library/openjdk
92a4e8a3140f: Mounted from library/openjdk
latest: digest: sha256:277e9efe2b389111a0d0df75f35739e2145161336f3c641c2fd61b0373c4850e size: 3262
[root@ip-172-31-26-205 hadoop]# sudo apt-get update
sudo: apt-get: command not found
[root@ip-172-31-26-205 hadoop]# sudo yum update -y
Last metadata expiration check: 1:06:16 ago on Mon Apr 29 00:18:56 2024.
Dependencies resolved.
Nothing to do.
Complete!
```

To get all the files present in the Trained Model on local we can run this command :

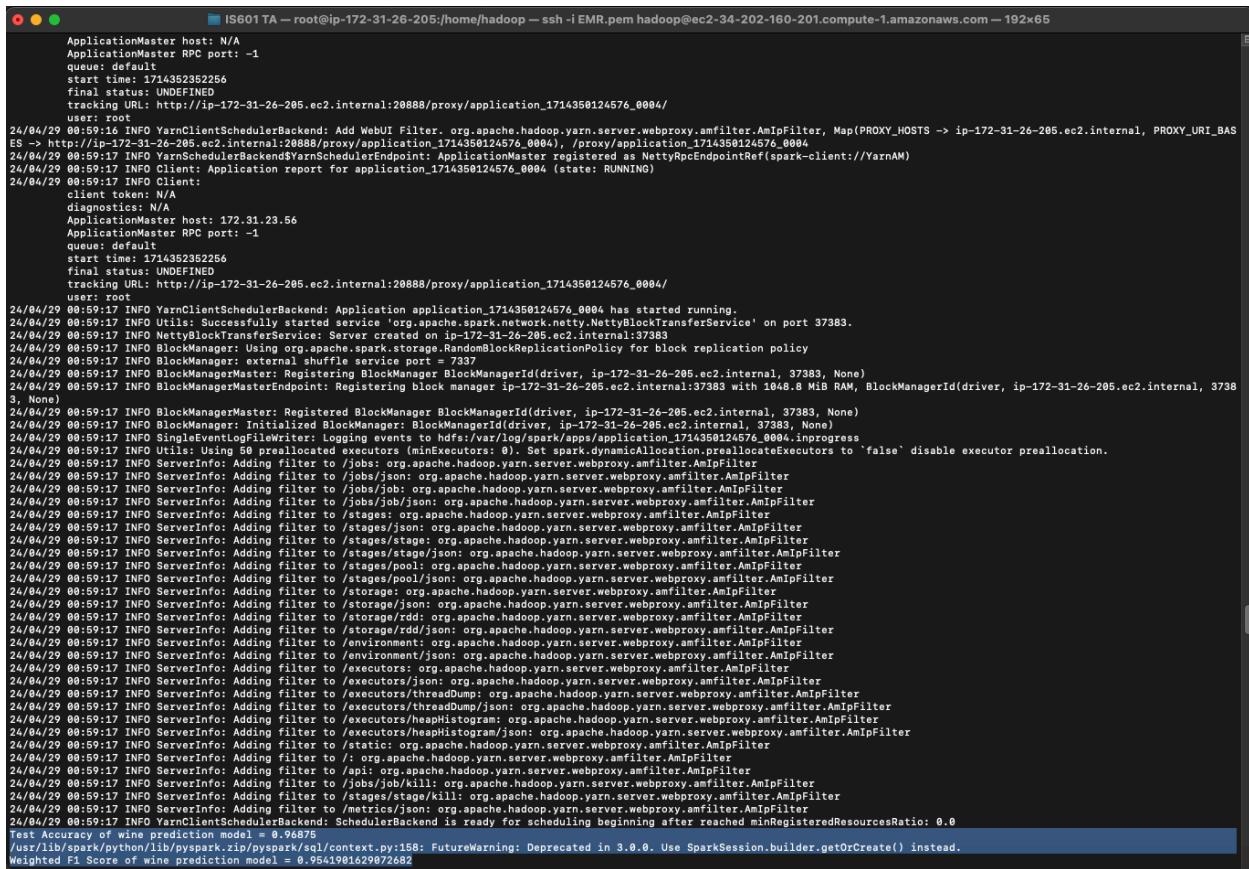
```
scp -i EMR.pem -r  
'hadoop@ec2-34-202-160-201.compute-1.amazonaws.com  
:/home/hadoop/*' .
```



A terminal window titled "IS601 TA — zsh — 125x46" showing the progress of an SCP command. The command is "scp -i EMR.pem -r 'hadoop@ec2-34-202-160-201.compute-1.amazonaws.com:/home/hadoop/\*' .". The terminal lists various files being transferred, including Dockerfile, README.md, TrainingDataset.csv, ValidationDataset.csv, and several sample files like applypatch-msg.sample, commit-msg.sample, post-update.sample, etc. The transfer speeds range from 0.6KB/s to 225.3KB/s, and the process is 100% complete.

```
Last login: Sun Apr 28 20:30:57 on ttys000  
mohankumargm@Mohans-MBP ~ % cd Desktop/Spring\ 2024/IS601\ TA  
mohankumargm@Mohans-MBP IS601 TA % scp -i EMR.pem -r 'hadoop@ec2-34-202-160-201.compute-1.amazonaws.com:/home/hadoop/*' .  
Dockerfile  
README.md  
TrainingDataset.csv  
ValidationDataset.csv  
applypatch-msg.sample  
commit-msg.sample  
post-update.sample  
pre-applypatch.sample  
pre-commit.sample  
pre-merge-commit.sample  
pre-push.sample  
pre-receive.sample  
push-to-checkout.sample  
update.sample  
fsmonitor-watchman.sample  
pre-rebase.sample  
prepare-commit-msg.sample  
exclude  
description  
HEAD  
config  
prediction.py  
part-00000  
part-00000  
part-00000-7e4dd2b5-6078-4758-afe7-fb88bbbe89b1-c000.snappy.parquet  
part-00000  
part-00000-305f26ee-50ad-4945-b0d5-1e68f9f74233-c000.snappy.parquet  
part-00000  
part-00000-2d43a826-d6b4-4ee0-8570-8de1c27cee87-c000.snappy.parquet  
training.py  
mohankumargm@Mohans-MBP IS601 TA %
```

# Important Screenshots :



The screenshot shows a terminal window titled "IS601 TA — root@ip-172-31-26-205:/home/hadoop — ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com — 192x65". The window displays a log of Hadoop application logs. The logs show the application starting up, registering with YARN, and initializing various services like BlockManager and AmIpFilter. The log ends with a test accuracy report and a warning about deprecated metrics.

```
ApplicationMaster host: N/A
ApplicationMaster RPC port: -1
queue: default
start time: 1714350124576
final status: UNDEFINED
tracking URL: http://ip-172-31-26-205.ec2.internal:20888/proxy/application_1714350124576_0004

24/04/29 00:59:16 INFO YarnClientSchedulerBackend: Add WebUI Filter. org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter, Map(PROXY_HOSTS -> ip-172-31-26-205.ec2.internal, PROXY_URI_BASES -> http://ip-172-31-26-205.ec2.internal:20888/proxy/application_1714350124576_0004)
24/04/29 00:59:17 INFO YarnSchedulerBackend$YarnSchedulerEndpoint: ApplicationMaster registered as NettyRpcEndpointRef[spark-client://YarnAM]
24/04/29 00:59:17 INFO Client: Application report for application_1714350124576_0004 (state: RUNNING)
24/04/29 00:59:17 INFO Client:
  client token: N/A
  diagnostics: N/A
  ApplicationMaster host: 172.31.23.56
  ApplicationMaster RPC port: -1
  queue: default
  start time: 1714350124576
  final status: UNDEFINED
  tracking URL: http://ip-172-31-26-205.ec2.internal:20888/proxy/application_1714350124576_0004/
  user: root

24/04/29 00:59:17 INFO YarnClientSchedulerBackend: Application application_1714350124576_0004 has started running.
24/04/29 00:59:17 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 37383.
24/04/29 00:59:17 INFO NettyBlockTransferService: Server created on ip-172-31-26-205.ec2.internal:37383
24/04/29 00:59:17 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block replication policy
24/04/29 00:59:17 INFO BlockManager: external shuffle service port = 7337
24/04/29 00:59:17 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383, None)
24/04/29 00:59:17 INFO BlockManagerMasterEndpoint: Registering master ip-172-31-26-205.ec2.internal:37383 with 1048.8 MiB RAM, BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383, None)
24/04/29 00:59:17 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383, None)
24/04/29 00:59:17 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, ip-172-31-26-205.ec2.internal, 37383)
24/04/29 00:59:17 INFO SingleEventLogFileWriter: Logging events to hdfs://var/log/spark/apps/application_1714350124576_0004.inprogress
24/04/29 00:59:17 INFO Utils: Using 50 preallocated executors (minExecutors: 0). Set spark.dynamicAllocation.preallocateExecutors to 'false' disable executor preallocation.
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/job: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/job/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/stage: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/stage/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/pool: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/pool/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage/rdd: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /storage/rdd/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /environment: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /environment/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/threadDump: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/threadDump/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/heaphistogram: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /executors/heaphistogram/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /static: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /api: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /jobs/job/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /stages/stage/kill: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: Adding filter to /metrics/json: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
24/04/29 00:59:17 INFO ServerInfo: BlockManager$BlockManagerBackend: BlockManager$BlockManagerBackend is ready for scheduling beginning after reached minRegisteredResourcesRatio: 0.0
Test Accuracy of wine prediction model = 0.96875
/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/context.py:158: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Weighted F1 Score of wine prediction model = 0.9541991429872682
```

```

IS601 TA — root@ip-172-31-26-205:/home/hadoop — ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com — 192x65
... 82 more

[[root@ip-172-31-26-205 hadoop]# docker login
Log in with your Docker ID or email address to push and pull images from Docker Hub. If you don't have a Docker ID, head over to https://hub.docker.com/ to create one.
You can log in with your password or a Personal Access Token (PAT). Using a limited-scope PAT grants better security and is required for organizations using SSO. Learn more at https://docs.docker.com/go/access-tokens/
Username: mg946@njit.edu
>Password:
WARNING! Your password will be stored unencrypted in /root/.docker/config.json.
Configure a credential helper to remove this warning. See
https://docs.docker.com/engine/reference/commandline/login/#credentials-store

Login Succeeded
[[root@ip-172-31-26-205 hadoop]# docker build -t mg946/wine_quality .
[+] Building 95.3s (17/17) FINISHED
--> [internal] load build definition from Dockerfile
--> [internal] transfer dockerfile: 1.94kB
--> [internal] load metadata for docker.io/library/openjdk:8-jre-slim
--> [auth] library/openjdk:pull token for registry-1.docker.io
--> [internal] load .dockercfgignore
--> [internal] transfer context: 2B
--> [ 1/1] FROM docker.io/library/openjdk:8-jre-slim@sha256:5318612937fb8bc0a12dd36da6761f4c7a2a20233c20d4eb8d497e4045a4f5
--> > resolve docker.io/library/openjdk:8-jre-slim@sha256:5318612937fb8bc0a12dd36da6761f4c7a2a20233c20d4eb8d497e4045a4f5
--> > sha256:1efc2761ff9f52c05d6e726cf7e6ccafdf8b9e0eed816bdd2f2860ad 31.37MB
--> > sha256:293d160227a67363547c03b377a593b7e76339b5860MB / 58MB
--> > sha256:5121a1e5557501555551578377a593b7e76339b5860MB / 58MB
--> > sha256:c02421c7a4bfc037d7fb87893c5fe145a2929fd39b5ee655701bf6c34a072d 41.70MB / 41.70MB
--> > sha256:53186129377bb8bc0a12dd36da6761f4c7a2a20233c20d4eb8d497e4045a4f5 549B / 549B
--> > sha256:285c61a1e5eeb7b3709729b69558670148c5fd6e6b704fae7d370042c51438 1.16KB / 1.16KB
--> > sha256:85b121a1ffdddfc7b66d18140b0b285ad1257bd11a676ddc7b108a3c0636c8 7.47KB / 7.47KB
--> > extracting sha256:1efc2761ff9f52c05d6e726cf7e6ccafdf8b9e0eed816bdd2f2860ad7
--> > extracting sha256:a272793da48276873890ac21b3:991053a7e864791aa8f82c39b7863c988093b
--> > extracting sha256:1a2deacc943157b2a50186781672aae0b145a6a4d88694b50f01d8f59fa7e
--> > extracting sha256:d2421c74abfc037d7fb87893c5fe145a329fd39b5ee655701bf6c34a072d
--> [internal] load build context: 99.59kB
--> [internal] transfer context: 99.59kB
--> [ 2/11] RUN apt-get update && apt-get install -y curl bzip2 wget unzip --no-install-recommends && rm -rf /var/lib/spt/lists/*
--> [ 3/11] RUN curl -s -L "https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh" --output /tmp/miniconda.sh && bash /tmp/miniconda.sh -b -f -p "/opt/miniconda"
--> [ 4/11] RUN pip install --no-cache pyspark==3.5.0 numpy pandas awscil
--> [ 5/11] WORKDIR /opt
--> [ 6/11] RUN wget --no-verbose -O apache-spark.tgz "https://archive.apache.org/dist/spark/spark-3.5.0/spark-3.5.0-bin-hadoop3.tgz" && tar -xf apache-spark.tgz && rm apache-s
--> [ 7/11] RUN wget https://repo1.maven.org/maven2/com/amazonaws/aws-java-sdk/1.8.0/aws-java-sdk-1.8.0.jar -P /opt/spark/jars/
--> [ 8/11] RUN wget https://repo1.maven.org/maven2/org/apache/hadoop/hadoop-aws/3.0.0/hadoop-aws-3.0.0.jar -P /opt/spark/jars/
--> [ 9/11] COPY prediction.py /opt/
--> [10/11] COPY ValidationDataset.csv /opt/
--> [11/11] COPY trainedmodel /opt/trainedmodel/
--> exporting to image
--> reporting layers
--> writing image sha256:98465d99ce944a6f3d75ae80f2ed660fc8f643cd2bfcb0daacd1d208cbe73cad
--> naming to docker.io/mg946/wine_quality
[[root@ip-172-31-26-205 hadoop]# docker run mg946/wine_quality
/opt/spark/bin/load-spark-env.sh: line 68: ps: command not found
Starting Spark Application
24/04/29 01:11:13 INFO SparkContext: Running Spark version 3.5.0
24/04/29 01:11:13 INFO SparkContext: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:13 INFO SparkContext: Java version 1.8.0_342
24/04/29 01:11:13 INFO SparkContext: Environment variables: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO ResourceUtils: No custom resources configured for spark.driver
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO SparkContext: Submitted application: Wine-Quality-Prediction-SPARK-ML
24/04/29 01:11:13 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 1, script: , vendor: , memory -> name: memory, amount: 1024, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpu -> name: cpus, amount: 1.0)
24/04/29 01:11:13 INFO ResourceProfile: Limiting resource is cpu

```

```
IS601 TA -- root@ip-172-31-26-205:/home/hadoop -- ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com - 192x65
Starting Spark Application
24/04/29 01:11:13 INFO SparkContext: Running Spark version 3.5.0
24/04/29 01:11:13 INFO SparkContext: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:13 INFO SparkContext: Java version 1.8.0_342
24/04/29 01:11:13 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO ResourceUtils: No custom resources configured for spark.driver.
24/04/29 01:11:13 =====
24/04/29 01:11:13 INFO SparkContext: Submitted application: Wine-Quality-Prediction-SPARK-ML
24/04/29 01:11:13 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores => name: cores, amount: 1, script: , vendor: , memory => name: memory, amount: 1024, scri
pt: vendor: , offHeap => name: offheap, amount: 0, script: , vendor: ), task resources: Map(cpu => name: cpus, amount: 1.0)
24/04/29 01:11:13 INFO ResourceProfileManager: Limiting resource is cpu
24/04/29 01:11:13 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/04/29 01:11:13 INFO SecurityManager: Changing view acls to root
24/04/29 01:11:13 INFO SecurityManager: Changing view acls to root
24/04/29 01:11:13 INFO SecurityManager: Changing view acls groups to:
24/04/29 01:11:13 INFO SecurityManager: Changing modify acls groups to:
24/04/29 01:11:13 INFO SecurityManager: authentication disabled; ui acls disabled; users with view permissions: root; groups with view permissions: EMPTY; users with modify pe
rmissions: root; groups with modify permissions: EMPTY
24/04/29 01:11:14 INFO Utils: Successfully started service 'sparkDriver' on port 37995.
24/04/29 01:11:14 INFO SparkEnv: Registering MapOutputTracker
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMaster
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/04/29 01:11:14 INFO BlockManager: Created local directory at /tmp/blockmgr-cde8d74c-cfce-a7af-a347-08b8e4c09440
24/04/29 01:11:14 INFO MemoryStore: MemoryStore started with capacity 366.3 MiB
24/04/29 01:11:14 INFO SparkEnv: Registering OutputCommitCoordinator
24/04/29 01:11:14 INFO JettyUtils: Start Jetty 0.8.0.v4848 for SparkUI
24/04/29 01:11:14 INFO Utils: Successfully started service 'SparkUI' on port 4040.
24/04/29 01:11:14 INFO Executor: Starting executor ID driver on host ee5c07e4c40f
24/04/29 01:11:14 INFO Executor: OS info Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:14 INFO Executor: Java version 1.8.0_342
24/04/29 01:11:14 INFO Executor: Starting executor with user classpath (userClassPathFirst = false): ''
24/04/29 01:11:14 INFO Executor: Created or updated repl class loader org.apache.spark.util.MutableURLClassLoader@7916642d for default.
24/04/29 01:11:14 INFO Executor: Registered executor ID driver with master 'spark://ip-172-31-26-205:37995'
24/04/29 01:11:14 INFO NettyBlockTransferService: Server created on ee5c07e4c40f:39683
24/04/29 01:11:14 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block replication policy
24/04/29 01:11:14 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, ee5c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: Registering block manager ee5c07e4c40f:39683 with 366.3 MiB RAM, BlockManagerId(driver, ee5c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, ee5c07e4c40f, 39683, None)
Test Accuracy of wine prediction model = 0.96875
/opt/spark/python/lib/pyspark.zip/pyspark/sql/context.py:168: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Weighted F1 Score of wine prediction model = 0.9541901629072682
Exiting Spark Application
[root@ip-172-31-26-205 hadoop]# docker push mg946/wine_quality
Using default tag: latest
The push refers to repository [docker.io/mg946/wine_quality]
2bc3e7ac574: Pushed
979d12ca3979: Pushed
86c4012ca3b1: Pushed
b72c5a50bd95: Pushed
f64b460f228: Pushed
3563f7b37d3: Pushed
5f70bf18a086: Pushed
4bf6a6a5ac5: Pushed
c2683a2a243: Pushed
1a7521fa43: Pushed
b66978cf4b41: Mounted from library/openjdk
cd5a0a9fe1e81: Mounted from library/openjdk
eafe6e032dbd: Mounted from library/openjdk
92a0e8a3140f: Mounted from library/openjdk
latest: digest: sha256:277e9efe2b38911a0d0df75f35739e2145161336f3c641c2fd61b0373c4850e size: 3262
[root@ip-172-31-26-205 hadoop]#
```

```

-D <property>=<value>           define a value for a given property
-fs <file://>/<url>                  specify default filesystem URL to use, overrides 'fs.defaultFS' property from configurations.
-jt <local|resourceManager|port>    specify a ResourceManager
--files <file...>                  specify a comma-separated list of files to be copied to the map/reduce cluster
--libjars <jar1,...>                specify a comma-separated list of jar files to be included in the classpath
--archives <archive1,...>          specify a comma-separated list of archives to be unarchived on the compute machines

The general command line syntax is:
command [genericOptions] [commandOptions]

Usage: hadoop fs [generic options] -copyFromLocal [-f] [-p] [-l] [-d] [-t <thread count>] [-q <thread pool queue size>] <localSrc> ... <dst>
[root@ip-172-31-26-205 hadoop] hadoop fs -copyFromLocal TrainingDataset.csv hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/
[root@ip-172-31-26-205 hadoop] hadoop fs -copyFromLocal ValidationDataset.csv hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/
[root@ip-172-31-26-205 hadoop] spark-submit Training.py
Apr 29, 2024 12:52:10 AM org.apache.spark.launcher.Log4jHotPatchOption staticJavaAgentOption
WARNING: spark.log4jHotPatch.enabled is set to true, but /usr/share/log4j/cve-2021-44228-hotpatch/jdk17/Log4jHotPatchFat.jar does not exist at the configured location
^[[A/usr/bin/python3: can't open file '/home/hadoop/TrainingPy': [Errno 2] No such file or directory
24/04/29 08:52:12 INFO ShutdownHookManager: Shutdown hook called
24/04/29 08:52:12 INFO ShutdownHookManager: Stopping all local directory /mnt/tmp/spark-f01777e3-7688-4226-825b-db7fe8ca82469
[root@ip-172-31-26-205 hadoop] hadoop fs -copyFromLocal ValidationDataset.csv hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/
Apr 29, 2024 12:52:23 AM org.apache.spark.launcher.Log4jHotPatchOption staticJavaAgentOption
WARNING: spark.log4jHotPatch.enabled is set to true, but /usr/share/log4j/cve-2021-44228-hotpatch/jdk17/Log4jHotPatchFat.jar does not exist at the configured location

Starting Spark Application in EMR
24/04/29 08:52:24 INFO SparkContext: Running Spark version 3.5.0-amzn-1
24/04/29 08:52:24 INFO SparkContext: OS info: Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 08:52:24 INFO SparkContext: Java version 17.0.10
24/04/29 08:52:24 INFO ResourceUtils: =====
24/04/29 08:52:24 INFO ResourceUtils: custom resources configured for spark.driver.
24/04/29 08:52:24 INFO ResourceUtils: =====
24/04/29 08:52:24 INFO SparkContext: Submitted application: Wine-Quality-Prediction-SPARK-ML
24/04/29 08:52:24 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cpus -> name: cores, amount: 4, script: , vendor: , memory -> name: memory, amount: 9486, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpus -> name: cpus, amount: 1.0)
24/04/29 08:52:24 INFO ResourceProfile: Limiting resource is cpus at 4 tasks per executor
24/04/29 08:52:24 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/04/29 08:52:24 INFO ResourceProfileManager: Registering view acls to: root
24/04/29 08:52:24 INFO SecurityManager: Changing view acls groups to:
24/04/29 08:52:24 INFO SecurityManager: Changing modify acls groups to:
24/04/29 08:52:24 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: root; groups with view permissions: EMPTY; users with modify permissions: root; groups with modify permissions: EMPTY
24/04/29 08:52:24 INFO ResourceUtils: Successfully started service 'sparkDriver' on port 37565.
24/04/29 08:52:27 INFO SparkEnv: Registering MapOutputTracker
24/04/29 08:52:27 INFO SparkEnv: Registering BlockManagerMaster
24/04/29 08:52:27 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
24/04/29 08:52:27 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
24/04/29 08:52:27 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/04/29 08:52:27 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-5ec1acf-a-b7-e4-89b0-65a002cef7bb
24/04/29 08:52:27 INFO MemoryStore: MemoryStore started with capacity 1048.8 MiB
24/04/29 08:52:27 INFO SparkEnv: Registering OutputCommitCoordinator
24/04/29 08:52:27 INFO SubResultCacheManager: Sub-result caches are disabled.
24/04/29 08:52:27 INFO JettyUtils: Started Jetty on 0.0.0.0:4400 for SparkUI
24/04/29 08:52:27 INFO Client: Successfully started service 'sparkDriver' on port 4048.
24/04/29 08:52:27 INFO Utils: Using 50 preallocated executors (minExecutors: 0). Set spark.dynamicAllocation.preallocateExecutors to 'false' disable executor preallocation.
24/04/29 08:52:27 INFO DefaultHARMFailoverProxyProvider: Connecting to ResourceManager at ip-172-31-26-205.ec2.internal/172.31.26.205:8082
24/04/29 08:52:28 INFO Configuration: resource-types.xml not found
24/04/29 08:52:28 INFO ResourceUtils: Unable to find 'resource-types.xml'
24/04/29 08:52:28 INFO Client: spark.yarn.archive is set to null, more than the maximum memory capability of the cluster (12288 MB per container)
24/04/29 08:52:28 INFO Client: Will allow our AM container, with its dependencies, including 384 MB overhead
24/04/29 08:52:28 INFO Client: Setting up container launch context for our AM
24/04/29 08:52:28 INFO Client: Preparing resources for our AM container
24/04/29 08:52:28 INFO Client: Setting up the launch environment for our AM container
24/04/29 08:52:28 INFO Client: Neither spark.yarn.jar nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
24/04/29 08:52:28 INFO Client: Uploading resource file: /tmp/blockmgr-5ec1acf-a-b7-e4-89b0-65a002cef7bb/libsparkStaging-1714358124574_0002/.spark_libs_293978290278611398.jar
24/04/29 08:52:30 INFO Client: Uploading resource file: /etc/spark/conf.dist/hive-site.xml -> hdfs://ip-172-31-26-205.ec2.internal:8020/user/root/.sparkStaging/application_1714358124576_0002/hive-site.xml

```

```

○  IS601 TA - root@ip-172-31-26-205:/home/hadoop - ssh -i EMR.pem hadoop@ec2-34-202-160-201.compute-1.amazonaws.com - 212x67
→ → naming to docker.io/mp946/wine_quality
[root@ip-172-31-26-205 hadoop]# docker run mp946/wine_quality
/opt/spark/bin/load-spark-env.sh: line 68: ps: command not found
Starting Spark Application
24/04/29 01:11:13 INFO SparkContext: Running Spark version 3.5.0
24/04/29 01:11:13 INFO SparkContext: OS info: Linux, 6.1.84-99.169.amzn2023.x86_64, amd64
24/04/29 01:11:13 INFO SparkContext: Java version 18.0.42
24/04/29 01:11:13 WARN NativeCodeLoader: Unable to load Native-hadoop library for your platform... using builtin-java classes where applicable
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO ResourceUtils: No custom resources configured for spark.driver.
24/04/29 01:11:13 INFO ResourceUtils: =====
24/04/29 01:11:13 INFO SparkContext: Default ResourceProfile file: Wine-Quality-Prediction-SPARK-ML
24/04/29 01:11:13 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 1, script: , vendor: , memory -> name: memory, amount: 1024, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpus -> name: cpus, amount: 1.0)
24/04/29 01:11:13 INFO ResourceProfile: Limiting resource is cpus
24/04/29 01:11:13 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/04/29 01:11:13 INFO SecurityManager: Changing view acls to: root
24/04/29 01:11:13 INFO SecurityManager: Changing modify acls to: root
24/04/29 01:11:13 INFO SecurityManager: Changing view acls groups to:
24/04/29 01:11:13 INFO SecurityManager: Changing modify acls groups to:
24/04/29 01:11:13 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: root; groups with view permissions: EMPTY; users with modify permissions: root; groups with modify permissions: EMPTY
24/04/29 01:11:14 INFO Utils: Successfully started service 'sparkDriver' on port 37995.
24/04/29 01:11:14 INFO SparkEnv: Registering MapOutputTracker
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMaster
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
24/04/29 01:11:14 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/04/29 01:11:14 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-cde8d74c-cfce-47af-s347-08b8e4c09440
24/04/29 01:11:14 INFO MemoryStore: MemoryStore started with capacity 366.3 MiB
24/04/29 01:11:14 INFO OutputCommitCoordinator: Registering OutputCommitCoordinator
24/04/29 01:11:14 INFO JettyUtils: Start Jetty 0.0.0.0:44040 for SparkUI
24/04/29 01:11:14 INFO Utils: Successfully started service 'SparkUI' on port 4404.
24/04/29 01:11:14 INFO Executor: Starting executor ID driver on host e65c07e4c40f
24/04/29 01:11:14 INFO Executor: Starting executor ID driver on host e65c07e4c40f
24/04/29 01:11:14 INFO Executor: Java version: 18.0.42
24/04/29 01:11:14 INFO Executor: Starting executor with user classpath (user$classPathFirst = false): ''
24/04/29 01:11:14 INFO Executor: Created or updated repl class loader org.apache.spark.util.MutableURLClassLoader@7916442d for default.
24/04/29 01:11:14 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 39683.
24/04/29 01:11:14 INFO NettyBlockTransferService: Server created on e65c07e4c40f:39683
24/04/29 01:11:14 INFO BlockManagerMaster: Registering BlockManagerBlockManagerId(driver, e65c07e4c40f:39683, None)
24/04/29 01:11:14 INFO BlockManagerMasterEndpoint: Registering block manager e65c07e4c40f:39683 with 366.3 MiB RAM, BlockManagerId(driver, e65c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, e65c07e4c40f, 39683, None)
24/04/29 01:11:14 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, e65c07e4c40f, 39683, None)
Test Accuracy of wine prediction model = 0.96875
/opt/spark/python/lib/pyspark.zip/pyspark/sql/context.py:158: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Wrote File Path: /tmp/blockmgr-5ec1acf-a-b7-e4-89b0-65a002cef7bb/libsparkStaging-1714358124574_0002/.spark_libs_293978290278611398.zip
Exiting Spark Application

```