



Apache Spark MLlib vs Apache Mahout

Big Data Analytics Project Report

Syed Muzammil Ahmed

ERP ID: 25371

May 31st, 2022

TABLE OF CONTENT

Serial No	TOPICS	Page No.
1	ABSTRACT	03
2	INTRODUCTION	03
3	DATASET	03
4	SYSYEM CONFIGURATION	04
4	DOCKER COMPOSE FILE	04
5	APACHE SPARK METHODOLOGY	05
6	MACHINE LEARNING USING PYSPARK	07
7	APACHE MAHOUT METHODOLOGY	14
8	DIFFICULTIES	16
9	CONCLUSION	17

ABSTRACT: -

In this project, we are basically setting up Apache Spark MLlib and Apache Mahout pipeline using a docker container. The goal is to compare performance of both Apache MLlib and Apache Mahout by performing different machine learning algorithms on both MLlib and mahout and checking out their performance in terms of time, speed, accuracy, and efficiency.

INTRODUCTION: -

Machine learning is the new boom in the software industry which helps in training the computer to think, organize and process data by itself. The main intent of machine learning is that the machine learns to observe data, extract important information from it and grasp on its own to predict, recommend or alter any action without any human mediation. This requires various algorithms over varied systems. For the ease of these algorithms, Apache has come up with frameworks Mahout and Spark, which in its different ways helps in implementing machine learning in a better way.

Apache Spark is a data processing framework that can quickly perform processing tasks on very large data sets and can also distribute data processing tasks across multiple computers, either on its own or in tandem with other distributed computing tools. It is a lightning-fast unified analytics engine for big data and machine learning. To support Python with Spark, the Apache Spark community released a tool, PySpark. Using PySpark, one can work Python programming language and that's what we are using in this Project.

Whereas The Apache Mahout distributed linear algebra framework delivers new tools and methods for performing data analysis, building machine learning data pipelines, and implementing machine learning models in production. Mahout is not an in-memory framework like spark.

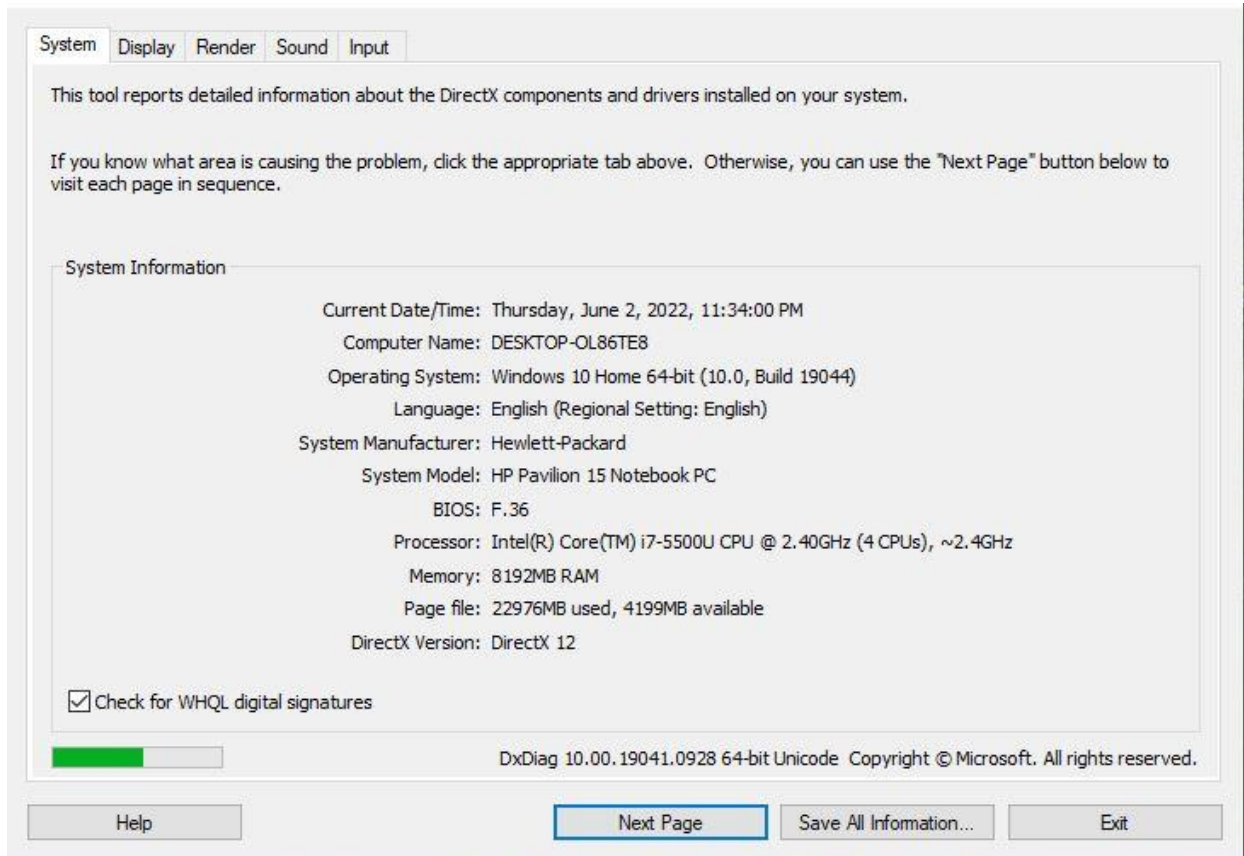
DATASET: -

Dataset we are using in this project is basically taken from famous dataset website Kaggle. Name of dataset is 1.6 million UK Accidents, this is a countrywide car accident dataset, which covers both urban and rural areas of the UK. Each accident record is described by a wide range of data attributes, including the accident location, weather, time, Date etc. This dataset is officially collected by UK Police department and contain mostly severe accidents in which police was also involved. Dataset is officially present on Government of UK website.

Size of our data is around 430mb, and it contain around **1.5million records** with **33 features**. Dataset consist of both categorical and continuous features.

SYSTEM CONFIGURATION: -

System configuration we are using in this project are as follows:



CREATING DOCKER COMPOSE FILE: -

In this project, we are creating YAML File which will help us to access our containers. Our YAML file contain 2 containers: one for Mahout and one for MLLIB. Docker Compose is a tool that was developed to help define and share multi-container applications. With Compose, we can create a YAML file to define the services and with a single command, can spin everything up or tear it all down. We can access our container using YAML file with command “**Docker-compose up -d**”.

```

1  version: "3.8"
2
3  services:
4
5    mahout:
6      image: michabirklbauer/mahout:latest
7      container_name: muzcont4
8
9      stdin_open: true
10
11     volumes:
12       - 'E:/Big_Data_Project/mydata:/data/'
13     restart: always
14     ports:
15       - 9000:9000
16     restart: always
17
18   mllib:
19     image: jupyter/pyspark-notebook
20     container_name: mycont3
21
22     stdin_open: true
23
24     volumes:
25       - 'E:/Big_Data_Project/mydata:/home/jovyan/data'
26     ports:
27       - 9088:9088
28     restart: always
29
30

```

Apache Spark Methodology: -

Apache Spark is a data processing framework that can quickly perform processing tasks on very large data sets and can also distribute data processing tasks across multiple computers, either on its own or in tandem with other distributed computing tools. It is a lightning-fast unified analytics engine for big data and machine learning

To support Python with Spark, the Apache Spark community released a tool, PySpark. Using PySpark, one can work with RDDs in Python programming language.

SETTING UP JUPYTER LAB ON DOCKER: -

- 1- First thing we need to do is to create a docker container on which we will pull docker image. We created container on docker named as “mycont2”.
- 2- After that we need to pull docker image from docker Hub into our container. Image size is around 2GB so it will take a quite time for downloading image.

```
C:\Users\Muzammil>docker run -it -d --name muzcont -p 8888:8888 jupyter/pyspark-notebook
Unable to find image 'jupyter/pyspark-notebook:latest' locally
latest: Pulling from jupyter/pyspark-notebook
d5fd17ec1767: Pull complete
3d97877476d9: Pull complete
7f2fba549eac: Pull complete
4f4fb700ef54: Pull complete
3c7df632b2eb: Pull complete
fe1eb68bb574: Pull complete
b08fbc176a13: Pull complete
9f7eb04384fc: Pull complete
86cab8aad649: Pull complete
7e6a3ce544a4: Pull complete
2ba4906ba99c: Pull complete
af17c5895b02: Pull complete
f60a067d4b86: Pull complete
dab8615aa3f2: Pull complete
bd83e344fa38: Pull complete
7081879b52c0: Pull complete
feb28b72853: Pull complete
2c4f66c5eb47: Pull complete
1e35175af6b1: Pull complete
e2702888a467: Pull complete
1b555678be34: Pull complete
b86aa2c799c8: Pull complete
3d3ea594328c: Pull complete
1002a301aa6e: Pull complete
01306fe31446: Pull complete
Digest: sha256:aaa476f803d711107999f195587ac7e27eb3ed9f596e57251435821f3726a16f
Status: Downloaded newer image for jupyter/pyspark-notebook:latest
30da10582d861b4bb3e3039474da47dc87650181dea410934a4336ed04feef99
```

- 3- Now, we need to setup jupyter lab on our container. For that reason, we execute jupyter lab command on container.

```
C:\Users\Muzammil>docker exec -it muzcont jupyter lab
[I 2022-05-28 20:50:50.224 ServerApp] jupyterlab | extension was successfully linked.
[W 2022-05-28 20:50:50.240 NotebookApp] 'ip' has moved from NotebookApp to ServerApp. This config will be passed to ServerApp. Be sure to update your config before our
next release.
[W 2022-05-28 20:50:50.240 NotebookApp] 'port' has moved from NotebookApp to ServerApp. This config will be passed to ServerApp. Be sure to update your config before ou
r next release.
[W 2022-05-28 20:50:50.241 NotebookApp] 'port' has moved from NotebookApp to ServerApp. This config will be passed to ServerApp. Be sure to update your config before ou
r next release.
[I 2022-05-28 20:50:50.261 ServerApp] nbclassic | extension was successfully linked.
[I 2022-05-28 20:50:50.745 ServerApp] notebook_shim | extension was successfully linked.
[I 2022-05-28 20:50:50.775 ServerApp] notebook_shim | extension was successfully loaded.
[I 2022-05-28 20:50:50.777 LabApp] JupyterLab extension loaded from /opt/conda/lib/python3.10/site-packages/jupyterlab
[I 2022-05-28 20:50:50.778 LabApp] JupyterLab application directory is /opt/conda/share/jupyter/lab
[I 2022-05-28 20:50:50.786 ServerApp] jupyterlab | extension was successfully loaded.
[I 2022-05-28 20:50:50.798 ServerApp] nbclassic | extension was successfully loaded.
[I 2022-05-28 20:50:50.799 ServerApp] The port 8888 is already in use, trying another port.
[I 2022-05-28 20:50:50.800 ServerApp] Serving notebooks from local directory: /home/jovyan
[I 2022-05-28 20:50:50.800 ServerApp] Jupyter Server 1.17.0 is running at:
[I 2022-05-28 20:50:50.801 ServerApp] http://30da10582d86:8889/lab?token=973802dd953cb782ac392ca720a1fb0b0b537985c6c47133
[I 2022-05-28 20:50:50.801 ServerApp] or http://127.0.0.1:8889/lab?token=973802dd953cb782ac392ca720a1fb0b0b537985c6c47133
[I 2022-05-28 20:50:50.801 ServerApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).

To access the server, open this file in a browser:
file:///home/jovyan/.local/share/jupyter/runtime/jpserver-169-open.html
Or copy and paste one of these URLs:
http://30da10582d86:8889/lab?token=973802dd953cb782ac392ca720a1fb0b0b537985c6c47133
or http://127.0.0.1:8889/lab?token=973802dd953cb782ac392ca720a1fb0b0b537985c6c47133
```

- 4- Now, we need to upload dataset on the container. So, that we can perform Machine learning on it so for that we will execute Docker cp command on another command prompt.

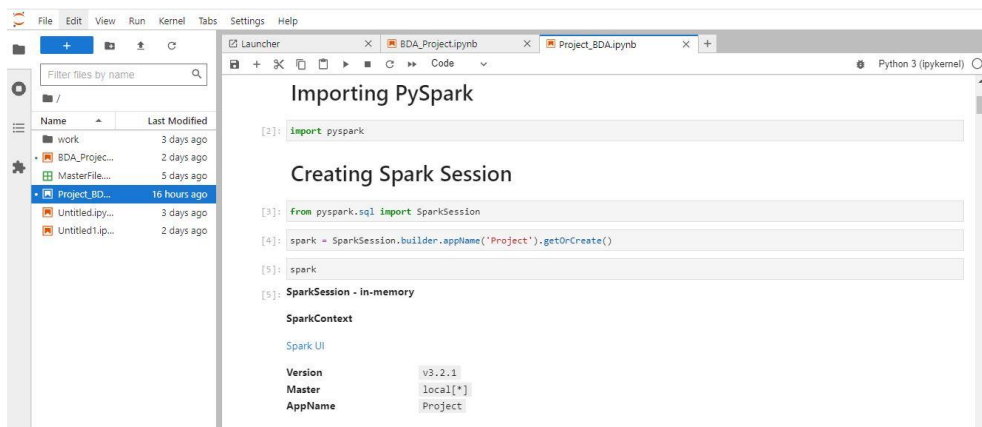
docker cp MasterFile.csv mycont2: \home\jovyan

MACHINE LEARNING USING PYSPARK: -

1) DATA LOADING AND SPARK SESSION: -

At first, we will import pyspark framework and create spark session. Spark session is a unified entry point of a spark application from Spark 2.0. It provides a way to interact with various spark's functionality with a lesser number of constructs. Instead of having a spark context, hive context, SQL context, now all of it is encapsulated in a Spark session.

After that we will import our dataset using pyspark commands and check the type of our data frame. Data frame created by pyspark is different than pandas they are `pyspark.sql. dataframe.DataFrame`.



2- DATA PRE – PROCESSING: -

Now, we are performing Data pre-processing using pyspark's code. We are performing different operations such as cleaning Null values, dropping null values, dropping irrelevant features, checking different data types of features etc.

1) Data Cleaning: -

Data cleaning is the process of eliminating or changing data that is inaccurate, incomplete, irrelevant, redundant, or incorrectly formatted to prepare it for analysis. When it comes to data analysis, this data is usually not necessary or beneficial because it can slow down the process or produce false results. Faulty data will badly impact our result, so we must remove it. So now after importing our dataset, we now check for null values using python pyspark's library and find that one of our column junction details is almost null. So, if we drop null values before dropping this column then it will make all our dataset empty. So, we are dropping this column junction detail. Now, we again check for null values and drop all the null values using pandas "dropna(how = "any")" function.

The screenshot shows a Jupyter Notebook with a file explorer on the left and a code editor on the right. The code editor has a tab for 'Project_BDA.ipynb'. The code in the cell is as follows:

```
[122]: from functools import reduce
df_agg_col = reduce(
    lambda a, b: a.union(b),
    (
        df_agg.select(F.lit(c).alias("Column_Name"), F.col(c).alias("NULL_Count"))
        for c in df_agg.columns
    )
)
df_agg_col.show()
```

The output of the code is a table with two columns: 'Column_Name' and 'NULL_Count'.

Column_Name	NULL_Count
Location_Easting_...	101
Location_Northing_...	101
Longitude	101
Latitude	101
Police_Force	0
Accident_Severity	0
Number_of_Vehicles	0
Number_of_Casualties	0
Day_of_Week	0
Local_Authority(...	0
1st_Road_Class	0

The screenshot shows the same Jupyter Notebook interface, but now the code editor has a tab for 'Project_BDA.ipynb'. The code in the cell is as follows:

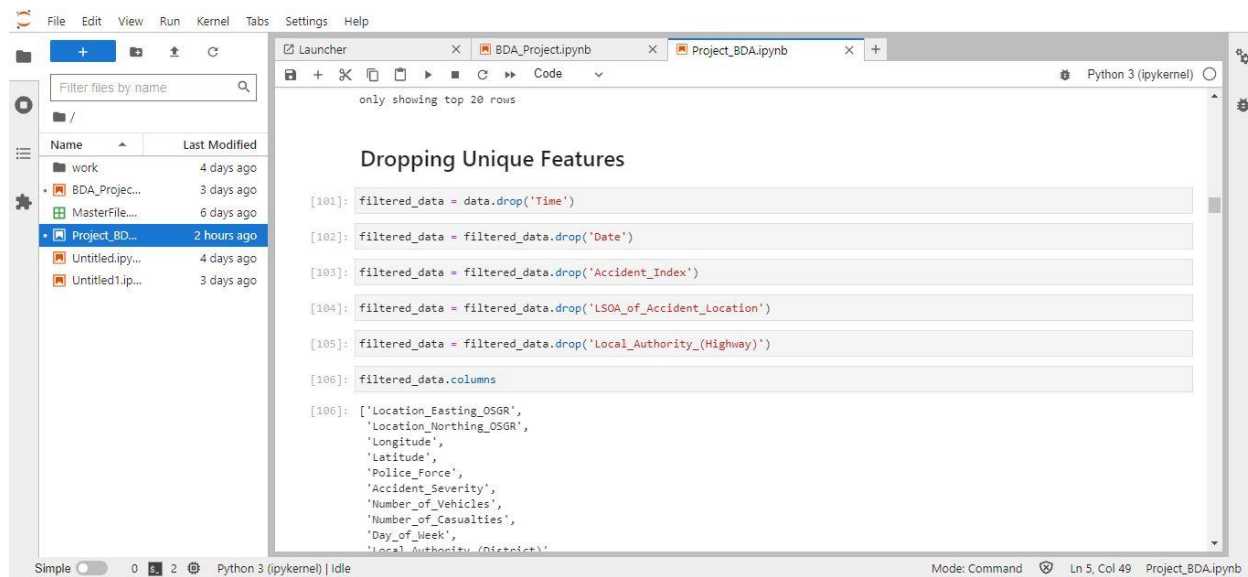
```
[132]: filter_data = filtered_data.dropna(how = 'any')
[133]: filter_data.count()
[133]: 1499613
[134]: filter_data.count()
[134]: 1504150
[135]: len(filter_data.columns)
[135]: 26
```

The output of the code is a table with two columns: 'Weather_Conditions' and '126'.

Weather_Conditions	126
only showing top 20 rows	

2) Data Transformation: -

Now we are checking our columns and dropping those columns which have all its values unique such as Date, time etc. Because in next step we must perform One-hot encoding. So, if also encode unique item columns, then it will generate too many columns which we cannot manage. So, we are dropping all the columns which have totally unique values. First, we are checking uniqueness of column using pandas. unique () function and check out different columns. The columns which we dropped here are Time, Date, Accident Index, LSOA Accident location, LSOA Highway.



The screenshot shows a Jupyter Notebook with a file explorer on the left and a code editor on the right. The file explorer shows a directory with files like 'work', 'BDA_Projec...', 'MasterFile...', 'Project_BD...', 'Untitled.ipyn...', and 'Untitled1.ip...'. The code editor has a tab titled 'Project_BDA.ipynb' and contains the following code:

```

Dropping Unique Features

[101]: filtered_data = data.drop('Time')

[102]: filtered_data = filtered_data.drop('Date')

[103]: filtered_data = filtered_data.drop('Accident_Index')

[104]: filtered_data = filtered_data.drop('LSOA_of_Accident_Location')

[105]: filtered_data = filtered_data.drop('Local_Authority_(Highway)')

[106]: filtered_data.columns

[106]: ['Location_Easting_OSGR',
'Location_Northing_OSGR',
'Longitude',
'Latitude',
'Police_Force',
'Accident_Severity',
'Number_of_Vehicles',
'Number_of_Casualties',
'Day_of_Week',
'Local_Authority_(District)']

```

The status bar at the bottom indicates 'Simple', 'Python 3 (ipykernel)', and 'Mode: Command | Ln 5, Col 49 | Project_BDA.ipynb'.

3) One-Hot Encoding: -

One hot encoding is a process of converting categorical data variables so they can be provided to machine learning algorithms to improve predictions. One hot encoding is a crucial part of feature engineering for machine learning. Our machine learning algorithms only understand numbers, so we must provide them numbers by converting categorical variables into continuous variables. It generates different columns and assigns them binary values. One hot encoding is essential before running a machine learning algorithm on a data set. Some algorithms can understand categorical data directly such as decision tree, but most of the supervised learning algorithms cannot operate on categorical data; they require all input variables to be numeric and also generate output in numeric value. This technique of transforming columns into binary variables in a data set is quite famous in supervised learning algorithms. These binary variables are also known as dummy variables in statistics.

Now, we must convert our categorical variable into continuous variables because Machine Learning only understand numbers. So, in pyspark, applying one hot encoding is not as simple as pandas because in this first you must implement string indexer. it is necessary to use the StringIndexer before OneHotEncoder, because OneHotEncoder needs a column of categorical indices as input. After string indexing, we apply One-Hot encoder.

The screenshot shows a Jupyter Notebook with the following code:

```
[141]: from pyspark.ml.feature import StringIndexer

[138]: catCols = [x for (x,dataType) in filter_data.dtypes if dataType == "string"]

[139]: print(catCols)

['Road_Type', 'Pedestrian_Crossing-Human_Control', 'Pedestrian_Crossing-Physical_Facilities', 'Light_Conditions', 'Weather_Conditions', 'Road_Surface_Conditions', 'Special_Conditions_at_Site', 'Carriageway_Hazards', 'Did_Police_Officer_Attend_Scene_of_Accident']

[142]: string_indexer = StringIndexer(inputCols = ['Road_Type', 'Pedestrian_Crossing-Human_Control', 'Pedestrian_Crossing-Physical_Facilities'],
                                     outputCols = ['Road_Type_StringIndexer', 'Pedestrian_Crossing-Human_Control_StringIndexer', 'Pedestrian_Crossing-Physical_Facilities_StringIndexer'])

[143]: indexed = string_indexer.fit(filter_data).transform(filter_data)

[144]: indexed.printSchema()

root
|-- Location_Easting_OSGR: double (nullable = true)
|-- Location_Northing_OSGR: double (nullable = true)
```

The screenshot shows a Jupyter Notebook with the following code:

```
[146]: from pyspark.ml.feature import OneHotEncoder

[148]: encoder = OneHotEncoder(inputCols = [f"{x}_StringIndexer" for x in catCols], outputCols = [f"{x}_OneHotEncoder" for x in catCols])

[149]: model = one_hot_encoder.fit(indexed).transform(indexed)

[150]: model.printSchema()

root
|-- Location_Easting_OSGR: double (nullable = true)
|-- Location_Northing_OSGR: double (nullable = true)
|-- Longitude: double (nullable = true)
|-- Latitude: double (nullable = true)
|-- Police_Force: integer (nullable = true)
|-- Accident_Severity: integer (nullable = true)
|-- Number_of_Vehicles: integer (nullable = true)
|-- Number_of_Casualties: integer (nullable = true)
|-- Day_of_Week: integer (nullable = true)
|-- Local_Authority_(District): integer (nullable = true)
|-- 1st_Road_Class: integer (nullable = true)
|-- 1st_Road_Number: integer (nullable = true)
|-- Road_Type: string (nullable = true)
|-- Speed_limit: integer (nullable = true)
|-- 2nd_Road_Class: integer (nullable = true)
|-- 2nd_Road_Number: integer (nullable = true)
|-- Pedestrian_Crossing-Human_Control: string (nullable = true)
|-- Pedestrian_Crossing-Physical_Facilities: string (nullable = true)
```

3- Vector Assembler: -

After One-Hot encoding, we use vector assembler to combine all features into one feature. The idea here is to assemble everything into a vector. This is reasonable since after one-hot encoding and stuff, you end up with a mishmash of integers, floats, sparse vectors, and may be dense vectors. And what we do next is bundle them altogether and call it features.

```
[151]: from pyspark.ml.feature import VectorAssembler

[152]: vector_assembler = VectorAssembler(
      inputCols = ["Location_Easting_OSGR", "Location_Northing_OSGR", "Longitude", "Latitude", "Police_Force", "Number_of_Vehicles", "Nu
    )

[153]: output = vector_assembler.transform(model)

[154]: model_filter_data = output.select("features", "Accident_Severity")
```

4- Train – Test Split: -

Next, we are doing is splitting our data into train and test dataset. Our train dataset contains 1051002 and test data contain 448611.

Train - Test Split

```
[166]: train , test = model_filter_data.randomSplit([0.7,0.3])
      print("Training Dataset Count: " + str(train.count()))
      print("Test Dataset Count: " + str(test.count()))

Training Dataset Count: 1051002
Test Dataset Count: 448611
```

Machine Learning Models: -

Now we are implementing our models and checking out their time and accuracy in MLlib. The models we are using are Decision Tree Classifier, Random Forest Classifier and Logistic Regression. For measurement of accuracy, we are using Multiclassification Evaluator module of Pyspark ML. So, below are the results we got after implementing models. All the models are implemented on basic parameters. Accuracy of each model is calculated using Multiclass evaluator module of a pyspark ml framework.

Models	Speed/Time	Accuracy
Decision Tree	176.33 sec	85.15%
Random Forest	347.45 sec	85.15%
Logistic Regression	107.33 sec	85.14%

```

[160]: %%time
decision_tree = DecisionTreeClassifier(featuresCol = 'features', labelCol = 'Accident_Severity')

CPU times: user 9.91 ms, sys: 0 ns, total: 9.91 ms
Wall time: 1.23 s

[165]: #import time
#start_time = time.time()
start = time.perf_counter()
decision_tree.fit(train)
end = time.perf_counter()
duration = format((end-start), '.4f')
print("Time in Seconds :", duration)

Time in Seconds : 176.3336

[210]: model = decision_tree.fit(train)

[211]: prediction_test = model.transform(test)

[212]: predictionAndLabels = prediction_test.select("Accident_Severity", "prediction").rdd

[213]: evaluator = MulticlassClassificationEvaluator(labelCol = "Accident_Severity", predictionCol = "prediction", metricName = "accuracy")
Accuracy_DT = evaluator.evaluate(prediction_test)
print("Accuracy Score of Decision Tree:" , Accuracy_DT)

```

Random Forest Classifier

```
[167]: from pyspark.ml.classification import RandomForestClassifier
       rf = RandomForestClassifier(featuresCol = 'features', labelCol = 'Accident_Severity')

[168]: start = time.perf_counter()
       rf.fit(train)
       end = time.perf_counter()
       duration = format((end-start), '.4f')
       print("Random Forest Classifier Time in Seconds :", duration)

Random Forest Classifier Time in Seconds : 347.4570

[205]: model = rf.fit(train)

[206]: prediction_test = model.transform(test)

[208]: predictionAndLabels = prediction_test.select("Accident_Severity", "prediction").rdd

[214]: evaluator = MulticlassClassificationEvaluator(labelCol = "Accident_Severity", predictionCol = "prediction", metricName = "accuracy")
       Accuracy_RF = evaluator.evaluate(prediction_test)
       print("Accuracy Score of Random Forest: ", Accuracy_RF)

Accuracy Score of Random Forest: 0.8515105514577217
```

Logistic Regression

```
[171]: from pyspark.ml.classification import LogisticRegression
       log_reg = LogisticRegression(featuresCol='features', labelCol='Accident_Severity')

[181]: start = time.perf_counter()
       log_reg.fit(train)
       end = time.perf_counter()
       duration = format((end-start), '.4f')
       print("Logistic Regression Time in Seconds :", duration)

Logistic Regression Time in Seconds : 107.3300

[189]: log_reg2 = LogisticRegression(featuresCol='features', labelCol='Accident_Severity', maxIter = 20)
       model = log_reg2.fit(train)

[191]: prediction_test = model.transform(test)

[192]: prediction_test.show()

+-----+-----+-----+-----+
|      features|Accident_Severity|      rawPrediction|      probability|prediction|
+-----+-----+-----+-----+
|(57,[0,1,2,3,4,5,...]|      3|[-8.7757203119566...|[3.38519414336258...|      3.0|
|(57,[0,1,2,3,4,5,...]|      3|[-8.7757203119566...|[3.38519414336258...|      3.0|
```

APACHE MAHOUT METHODOLOGY

Now, after getting results from apache spark, we must move forward to apply ML models on Apache Mahout, which is also framework of Hadoop, specially designed to perform ML models on Big Data sets. Mahout run over HDFS, so its speed is known to be quite slow as compared to apache spark which is in-memory framework. Mahout is rarely used, many of models are outdated and very less resource is available on internet on this framework.

SETTING UP MAHOUT CONTAINER: -

- 1- So, we need to create separate container for apache mahout, on which we will pull docker image. We created container on docker named as “mycont3”. Apache Mahout only support models in Java or Scala language so python is not supported in Mahout. So here we are also pulling image from Docker Hub.
- 2- After that we need to pull docker image from docker Hub into our container. Image size is around 3.9 GB so it will take a quite time for downloading image.

```

Microsoft Windows [Version 10.0.19044.1706]
(c) Microsoft Corporation. All rights reserved.

C:\Users\Muzamill>docker run -it --name muzcont3 michabirkibauer/mahout:latest
Unable to find image 'michabirkibauer/mahout:latest' locally
latest: Pulling from michabirkibauer/mahout
ecb575e29c: Pull complete
467d1831b69: Pull complete
eab2c490a3c: Pull complete
15a0f46f833: Pull complete
6cb1dfcbebb: Pull complete
b224ce6d4ea: Pull complete
932fe81bb40: Pull complete
c39e3902a25: Pull complete
8fadcbes53b: Pull complete
e3d54adf2d0: Pull complete
e3d56b510ec: Pull complete
ee173c37d0f: Pull complete
abd0bf9d9e1: Pull complete
6866a3d54fd: Pull complete
fb87913f780: Pull complete
bfe7931e2d: Pull complete
12fcec458b0: Downloading [=====] 70.88MB/185.3MB
ec3180c7883: Downloading [=====] 63.92MB/175.5MB
643bf1aed1b: Downloading [=====] 67.66MB/175.6MB
c6258235011: Waiting
36e64712f6c: Waiting
47d855c16f2: Pulling fs layer
aa6c9a8bbf4: Waiting
d6493e48cb8: Waiting
51c2f69bd7b: Waiting

```

```

spark context Web UI available at http://bef525381492:4040
spark context available as 'sc' (master = local[*], app id = local-1654122347868).
spark session available as 'spark'.

mahout version 0.14

warning: File 'spark-shell' does not exist.
welcome to

scala version 3.1.1

using Scala version 2.12.10 (OpenJDK 64-Bit Server VM, Java 1.8.0_282)
type in expressions to have them evaluated.
type :help for more information.

scala>

```



```
scala> df.counts
<console>:46: error: value counts is not a member of org.apache.spark.sql.DataFrame
df.counts
  ^

scala> df.counts()
<console>:46: error: value counts is not a member of org.apache.spark.sql.DataFrame
df.counts()
  ^

scala> df.count
res4: Long = 1504150
```

Importing some importing features

```
Command Prompt - docker exec -it muzcont4 /bin/bash

scala> import org.apache.mahout.math._
import org.apache.mahout.math._

scala> import org.apache.mahout.math._
import org.apache.mahout.math._

scala> import org.apache.mahout.math.scalabindings._
import org.apache.mahout.math.scalabindings._

scala> import org.apache.mahout.math.drm._
import org.apache.mahout.math.drm._

scala> import org.apache.mahout.math.scalabindings.RLikeOps._
import org.apache.mahout.math.scalabindings.RLikeOps._

scala> import org.apache.mahout.math.drm.RLikeDrmOps._
import org.apache.mahout.math.drm.RLikeDrmOps._

scala> import org.apache.mahout.sparkbindings._
import org.apache.mahout.sparkbindings._

scala>

scala> implicit val sdc: org.apache.mahout.sparkbindings.SparkDistributedContext = sc2sdc(sc)
sdc: org.apache.mahout.sparkbindings.SparkDistributedContext = org.apache.mahout.sparkbindings.SparkDistributedContext@7c8c5798
```

We can't process more and not able to run any model on Apache Mahout because of less resources available on Apache mahout on internet. I am not able to find anything useful enough to train my model on Apache mahout. Even Mahout official documentation is not providing any relevant information about mahout algorithms. So, it is very difficult for us at a moment to run any model for comparison. Only one algorithm of Least linear regression is available, but our dataset is classification, and no information is available on that, so that's why we are leaving it here.

DIFFICULTIES DURING PROJECT:

These are some difficulties I faced during this project.

- 1- Naïve Bayes algorithm on pyspark, not working on Jupyter lab and giving Java Error.
- 2- Gradient Boosting on Pyspark not supporting multi-class problem.
- 3- No proper resources of mahout available on internet.
- 4- Mahout algorithms and resources not even available on official sites.
- 5- Mahout does not support Python.

CONCLUSION: -

So, I am only able to complete pyspark modelling in this project and not able to perform modelling on mahout because of almost no resource available on modelling on mahout. I have tried my best to find some useful resource but not able to succeed. So that's why we are not able to compare performance of both Mllib and Mahout. I believe Apache Mahout is outdated now as people are using spark because of its high speed and efficiency. Apache mahout algorithms resources are also deleted from official page of apache mahout.