# Network Intrusion Detection

Big Data Analytics (W1)

Ali Munawar

F2019332007

Submitted To:

Mazhar Javed Awan

## Overview

The intrusion detection systems are an integral part of modern communication networks. The business environments require a high level of security to safeguard their private data from any unauthorized personnel. The current intrusion detection systems are a step upgrade from the conventional anti-virus software. Two main categories based on their working. These are: • Network Intrusion Detection Systems (NIDS): These systems continuously monitor the network traffic and analyze the packets for a possible rule infringement. • Host-based Intrusion Detection Systems (HIDS): These systems monitor the operating system files of an end-user system to detect malicious software that might temper with its normal functioning.

## About Dataset

Background
The dataset to be audited was provided which consists of a wide variety of intrusions simulated in a military network environment. It created an environment to acquire raw TCP/IP dump data for a network by simulating a typical US Air Force LAN. The LAN was focused like a real environment and blasted with multiple attacks. A connection is a sequence of TCP packets starting and ending at some time duration between which data flows to and from a source IP address to a target IP address under some well-defined protocol. Also, each connection is labeled as either normal or as an attack with exactly one specific attack type. Each connection record consists of about 100 bytes.
For each TCP/IP connection, 41 quantitative and qualitative features are obtained from normal and attack data (3 qualitative and 38 quantitative features) .The class variable has two categories:
• Normal
• Anomalous
The dataset is taken from Kaggle:
https://www.kaggle.com/datasets/sampadab17/network-intrusion-detection

# Operations on Resilient distributed datasets

Importing pyspark and SparkContext in databricks

```python
1   import pyspark
```
Command took 0.65 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 23:04:23 on My Cluster

```python
1   from pyspark.context import SparkContext
```
Command took 0.04 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:20 on My Cluster

## Action:

RDD actions are operations that return non-RDD values. Here we have the **count()** which will tell us about the total elements in the dataset.

```python
1   #action
2   network_data.count()
```
▸ (1) Spark Jobs

Out[6]: 25193

## Collect() is used to show the

```python
1   network_data.collect()
```
▸ (1) Spark Jobs

```
Out[7]: ['duration,protocol_type,service,flag,src_bytes,dst_bytes,land,wrong_fragment,urgent,hot,num_failed_logins,logged_in,num_compromised,root_shell,su_attempted,num_root,num_file_creations,num
_shells,num_access_files,num_outbound_cmds,is_host_login,is_guest_login,count,srv_count,serror_rate,srv_serror_rate,rerror_rate,srv_rerror_rate,same_srv_rate,diff_srv_rate,srv_diff_host_rate,dst_h
ost_count,dst_host_srv_count,dst_host_same_srv_rate,dst_host_diff_srv_rate,dst_host_same_src_port_rate,dst_host_srv_diff_host_rate,dst_host_serror_rate,dst_host_srv_serror_rate,dst_host_rerror_rat
e,dst_host_srv_rerror_rate,class',
 '0,tcp,ftp_data,SF,491,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,2,2,0,0,0,0,1,0,0,150,25,0.17,0.03,0.17,0,0,0,0.05,0,normal',
 '0,udp,other,SF,146,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,13,1,0,0,0,0.08,0.15,0,255,1,0,0.6,0.88,0,0,0,0,normal',
 '0,tcp,private,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,123,6,1,1,0,0,0.05,0.07,0,255,26,0.1,0.05,0,0,1,1,0,0,anomaly',
 '0,tcp,http,SF,232,8153,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,5,5,0.2,0.2,0,0,1,0,0,30,255,1,0,0.03,0.04,0.03,0.01,0,0.01,normal',
 '0,tcp,http,SF,199,420,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,30,32,0,0,0,0,1,0,0.09,255,255,1,0,0,0,0,0,0,0,normal',
 '0,tcp,private,REJ,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,121,19,0,0,1,1,0.16,0.06,0,255,19,0.07,0.07,0,0,0,0,1,1,anomaly',
 '0,tcp,private,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,166,9,1,1,0,0,0.05,0.06,0,255,9,0.04,0.05,0,0,1,1,0,0,anomaly',
 '0,tcp,private,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,117,16,1,1,0,0,0.14,0.06,0,255,15,0.06,0.07,0,0,1,1,0,0,anomaly',
 '0,tcp,remote_job,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,270,23,1,1,0,0,0.09,0.05,0,255,23,0.09,0.05,0,0,1,1,0,0,anomaly',
 '0,tcp,private,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,133,8,1,1,0,0,0.06,0.06,0,255,13,0.05,0.06,0,0,1,1,0,0,anomaly',
 '0,tcp,private,REJ,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,205,12,0,0,1,1,0.06,0.06,0,255,12,0.05,0.07,0,0,0,0,1,1,anomaly',
 '0,tcp,private,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,199,3,1,1,0,0,0.02,0.06,0,255,13,0.05,0.07,0,0,1,1,0,0,anomaly',
 '0,tcp,http,SF,287,2251,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,3,7,0,0,0,0,1,0,0.43,8,219,1,0,0.12,0.03,0,0,0,0,normal',
 '0,tcp,ftp_data,SF,334,0,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,2,2,0,0,0,0,1,0,0.2,20,1,0,1,0.2,0,0,0,0,anomaly',
 '0,tcp,name,S0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,233,1,1,1,0,0,0.06,0,255,1,0,0.07,0,0,1,1,0,0,anomaly',
```
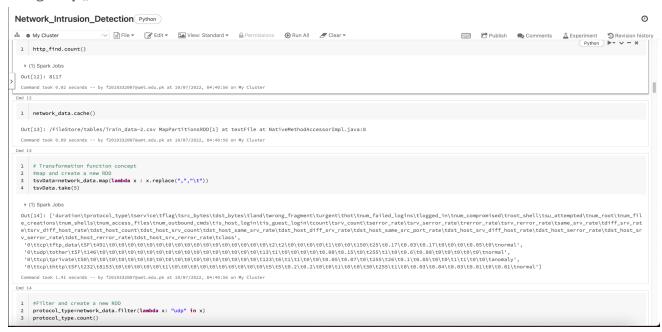
## Transformation:

Spark Transformation is a function that produces new RDD from the existing RDDs.

Here we have created a new RDD named http_find.

```python
1   #transformation
2   http_find=network_data.filter(lambda line:"http" in line)
```
Command took 0.05 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:20 on My Cluster

```python
1   http_find
```
Out[10]: PythonRDD[36] at RDD at PythonRDD.scala:58
Command took 0.12 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:20 on My Cluster

```python
1   http_find.collect()
```
▸ (1) Spark Jobs

```
Out[11]: ['0,tcp,http,SF,232,8153,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,8,0,5,5,0.2,0.2,0,0,1,0,0,30,255,1,0,0.03,0.04,0.03,0.01,0,0.01,normal',
 '0,tcp,http,SF,199,420,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,30,32,0,0,0,0,1,0,0.09,255,255,1,0,0,0,0,0,0,0,normal',
 '0,tcp,http,SF,287,2251,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,3,7,0,0,0,0,1,0,0.43,8,219,1,0,0.12,0.03,0,0,0,0,normal',
 '0,tcp,http,SF,300,13788,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,8,9,0,0.11,0,0,1,0,0.22,91,255,1,0.01,0.02,0,0,0,0,normal',
 '0,tcp,http,SF,233,616,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,3,3,0,0,0,0,1,0,0.66,255,1,0,0.02,0.03,0,0,0.02,0,normal',
 '0,tcp,http,SF,343,1178,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,9,10,0,0,0,0,1,0,0.2,157,255,1,0.01,0.04,0,0,0,0,normal',
 '0,tcp,http,SF,253,11905,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,8,10,0,0,0,0,1,0,0.2,87,255,1,0.01,0.02,0,0,0,0,normal',
 '0,tcp,http,SF,227,6588,0,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,5,22,0,0,0,0,1,0,0.18,43,255,1,0.02,0.14,0,0,0.56,0.57,normal',
```
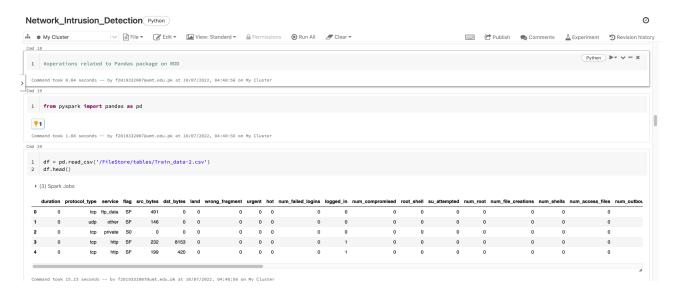
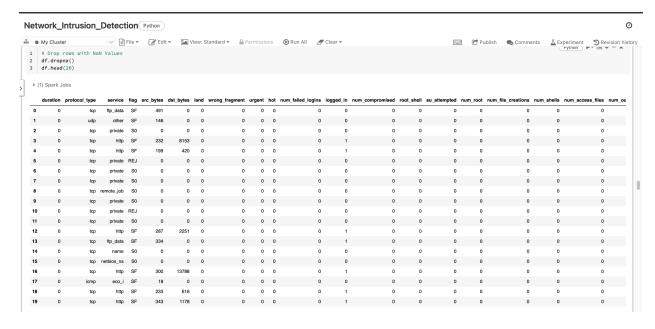## Using **map()** to create a new RDD



## Filtering using **flatmap()**

## Operations related to pandas in RDD

Pandas is a software library written for the Python programming language for data manipulation and analysis, it offers data structures and operations for manipulating numerical tables and time series.
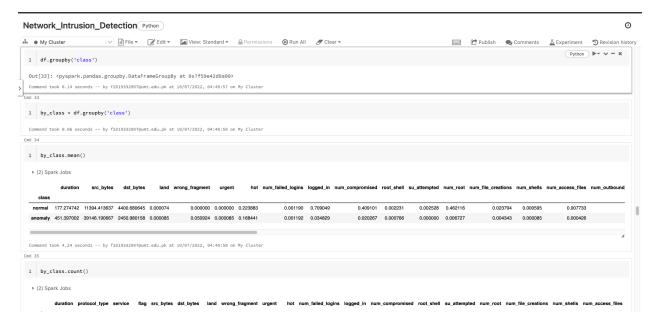

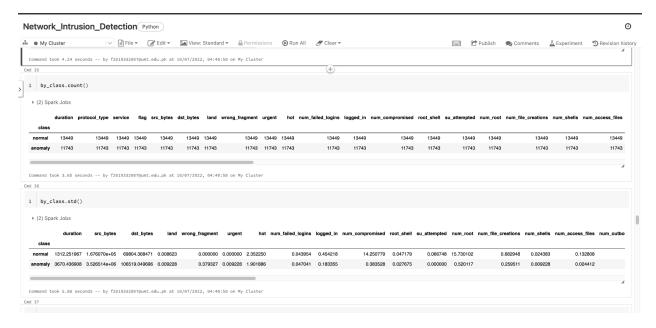
Dropping the null values in our dataset.

## Applying pandas **GroupBy()**

It allows you to split your data into separate groups to perform computations for better analysis. So, we have split the class and taken the mean of it, to see the total normals and anomalies.



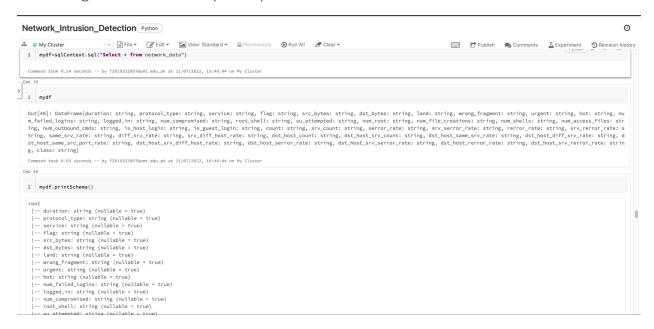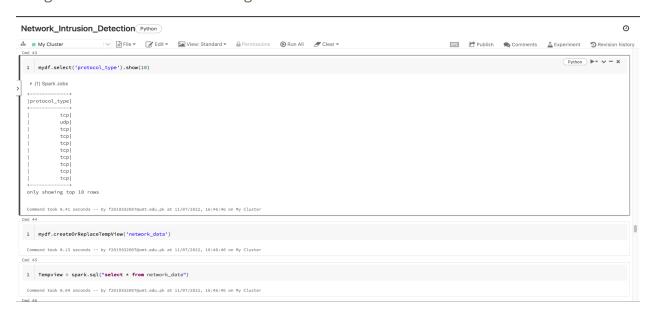Taking the count and std of our class variable.

# Operations on Spark SQL

**Spark SQL:**

Spark SQL is a Spark module for structured data processing. It provides a programming abstraction called DataFrames and can also act as a distributed SQL query engine.

Converting our dataset into spark sql.



Using select command and making new table

The **printSchema()** method is used to display the schema of the dataframe.

```
1   mydf.printSchema()
```

```
root
 |-- duration: string (nullable = true)
 |-- protocol_type: string (nullable = true)
 |-- service: string (nullable = true)
 |-- flag: string (nullable = true)
 |-- src_bytes: string (nullable = true)
 |-- dst_bytes: string (nullable = true)
 |-- land: string (nullable = true)
 |-- wrong_fragment: string (nullable = true)
 |-- urgent: string (nullable = true)
 |-- hot: string (nullable = true)
 |-- num_failed_logins: string (nullable = true)
 |-- logged_in: string (nullable = true)
 |-- num_compromised: string (nullable = true)
 |-- root_shell: string (nullable = true)
 |-- su_attempted: string (nullable = true)
 |-- num_root: string (nullable = true)
 |-- num_file_creations: string (nullable = true)
 |-- num_shells: string (nullable = true)
 |-- num_access_files: string (nullable = true)
 |-- num_outbound_cmds: string (nullable = true)
```
Command took 0.04 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:23 on My Cluster

**Cache()** DataFrame **(nework_data)** into memory

```
1   network_data.cache()
```

Out[13]: /FileStore/tables/Train_data.csv MapPartitionsRDD[33] at textFile at NativeMethodAccessorImpl.java:0

Command took 0.13 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:21 on My Cluster

## Visualization using matplotlib and seaborn:

Matplotlib is a python library used extensively for the visualization of data. While Seaborn is a python library based on matplotlib. Seaborn provides a high-level interface for drawing attractive and informative statistical graphics.

## Plotting the total value counts of service

```python
1  df['service'].value_counts().plot(kind='bar')
```



## Displot using seaborn



## Boxplot using seaborn

# Shared Variables

Spark supports two types of shared variables: broadcast variables, which can be used to cache a value in memory on all nodes, and accumulators, which are variables that are only "added" to, such as counters and sums.

### Accumulators:

Accumulators are variables that are only "added" to through an associative operation and can therefore, be efficiently supported in parallel. They can be used to implement counters (as in MapReduce) or sums.

### Broadcast:

Broadcast variables allow the programmer to keep a read-only variable cached on each machine rather than shipping a copy of it with tasks. They can be used, for example, to give every node a copy of a large input dataset in an efficient manner.

```python
1  ac=sc.accumulator(0)
2  ac1=sc.accumulator(0)
3  b=sc.broadcast("normal")
4  c=sc.broadcast("anomaly")
```

Command took 0.14 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:21 on My Cluster

Cmd 21

```python
1  def count(data):
2      global ac
3      global ac1
4      if b.value in data:
5          ac=ac+1
6      if c.value in data:
7          ac1=ac1+1
8      return count
9
10  print(ac)
11  print(ac1)
12
13  shared=network_data.map(count)

0
0
```
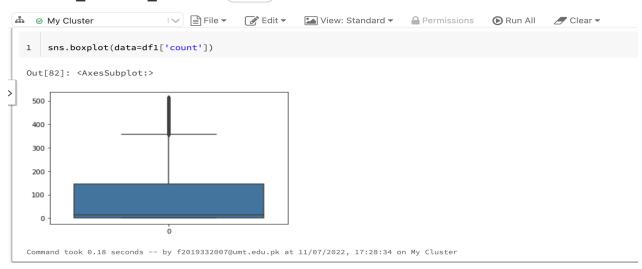
Command took 0.13 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:21 on My Cluster

# Spark Machine Learning

MLlib is Sparks machine learning (ML) library. Its goal is to make practical machine learning scalable and easy.

**Use cases of machine learning with Spark:**

- ML Algorithms: common learning algorithms such as classification, regression, clustering, and collaborative filtering
- Featurization: feature extraction, transformation, dimensionality reduction, and selection

- Pipelines: tools for constructing, evaluating, and tuning ML Pipelines
- Persistence: saving and load algorithms, models, and Pipelines
- Utilities: linear algebra, statistics, data handling, etc.

**Spark MLlib (w/ RDDs) versus Spark ML (w/ DataFrames)**

Spark.mllib carries the original API built on top of RDDs whereas Spark.ml contains higher-level API built on top of DataFrames for constructing ML pipelines.

Importing **StringIndexer** from **pyspark.ml.feature**

<div>Markdown</div>

## SPARK ML

Cmd 72

```python
1  from pyspark.ml.feature import StringIndexer
2
```

Command took 0.07 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:26 on My Cluster

Cmd 73

```python
1  from pyspark.ml.feature import StringIndexer
2
```

Command took 0.05 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:26 on My Cluster

**Extracting, Transforming** and **selecting features** in our model through **StringIndexer**

<div>Python</div>

```python
1  s_ml = StringIndexer(
2      inputCol='class',
3      outputCol='label',
4      handleInvalid='keep').fit(s_ml).transform(s_ml)
5  s_ml= StringIndexer(
6      inputCol='protocol_type',
7      outputCol='Type_Of_Protocol',
8      handleInvalid='keep').fit(s_ml).transform(s_ml)
9
```

▶ (4) Spark Jobs

Command took 2.63 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:26 on My Cluster

Cmd 75

```python
1  s_ml.show(5)
```

▶ (1) Spark Jobs

```
+-------------+---------+------+----------+-------+-----+----------------+
|protocol_type|src_bytes|urgent|rerror_rate|  class|label|Type_Of_Protocol|
+-------------+---------+------+----------+-------+-----+----------------+
|          tcp|      491|     0|         0| normal|  0.0|             0.0|
|          udp|      146|     0|         0| normal|  0.0|             1.0|
|          tcp|        0|     0|         0|anomaly|  1.0|             0.0|
|          tcp|      232|     0|         0| normal|  0.0|             0.0|
|          tcp|      199|     0|         0| normal|  0.0|             0.0|
+-------------+---------+------+----------+-------+-----+----------------+
only showing top 5 rows
```

Required Features

Cmd 79

```python
1  required_features = ['src_bytes',
2                       'urgent',
3                       'rerror_rate',
4                       'label',
5                       'Type_Of_Protocol'
6                       ]
```

Command took 0.03 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:27 on My Cluster

## Vector Assembler:

VectorAssembler is a transformer that combines a given list of columns into a single vector column. It is useful for combining raw features and features generated by different feature transformers into a single feature vector, in order to train ML models.

Using VectorAssembler to convert the features into vector columns.

```
1   from pyspark.ml.feature import VectorAssembler
```

Command took 0.03 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:27 on My Cluster

Cmd 81

```
1   from pyspark.sql.functions import col
```

Command took 0.04 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:27 on My Cluster

Cmd 82

```
1   s_ml = s_ml.select(col('src_bytes').cast('int'),
2                      col('urgent').cast('int'),
3                      col('rerror_rate').cast('int'),
4                      col('label').cast('int'),
5                      col('Type_Of_Protocol').cast('int')
6
7                     )
8   s_ml.show()
```

▶ (1) Spark Jobs

```
+---------+------+-----------+-----+----------------+
|src_bytes|urgent|rerror_rate|label|Type_Of_Protocol|
+---------+------+-----------+-----+----------------+
|      491|     0|          0|    0|               0|
|      146|     0|          0|    0|               1|
|        0|     0|          0|    1|               0|
|      232|     0|          0|    0|               0|
|      199|     0|          0|    0|               0|
|        0|     0|          1|    1|               0|
|        0|     0|          0|    1|               0|
|        0|     0|          0|    1|               0|
|        0|     0|          0|    1|               0|
|        0|     0|          0|    1|               0|
|        0|     0|          1|    1|               0|
|        0|     0|          0|    1|               0|
```

Cmd 83

```
1   assembler = VectorAssembler(inputCols=required_features, outputCol='features')
2
```

Command took 0.07 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:27 on My Cluster

Cmd 84

```
1   transformed_data = assembler.transform(s_ml)
```

Command took 0.15 seconds -- by f2019332007@umt.edu.pk at 17/07/2022, 15:25:27 on My Cluster

## Machine Learning Algorithms

Following are the algorithms which I have implemented in this project:

- Random Forest Classifier
- Decision Tree Classification
- Logistic Regression
- Naive Bayes Classifier

### RANDOM FOREST CLASSIFIER

Cmd 87

```python
from pyspark.ml.classification import RandomForestClassifier
rf = RandomForestClassifier(labelCol='label',
                            featuresCol='features',
                            maxDepth=5)
```

Command took 0.15 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:40 on My Cluster

## Test Accuracy of Random Forest Classifier

Cmd 92

```python
accuracy = evaluator.evaluate(predictions)
print('Test Accuracy = ', accuracy)
```

▶ (1) Spark Jobs

Test Accuracy =  1.0

Command took 1.47 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:40 on My Cluster

### DECISION TREE CLASSIFICATION

Cmd 94

```python
from pyspark.ml.classification import DecisionTreeClassifier
```

Command took 0.02 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:41 on My Cluster

Cmd 95

```python
DT=DecisionTreeClassifier(labelCol='label',
                          featuresCol='features',
                          maxDepth=5)
```

Command took 0.08 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:41 on My Cluster

## Test Accuracy of Decision Tree

Cmd 100

```python
accuracy = evaluator.evaluate(pre)
print('Test Accuracy = ', accuracy)
```

▶ (1) Spark Jobs

Test Accuracy =  1.0

Command took 1.12 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:42 on My Cluster

## LOGISTIC REGRESSION

Cmd 102

```
1  from pyspark.ml.classification import LogisticRegression
```

Command took 0.03 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:42 on My Cluster

Cmd 103

```
1  Lr=LogisticRegression(labelCol='label',
2                        featuresCol='features',
3                        maxIter=100)
```

## Test Accuracy of Logistic Regression

Cmd 107

```
1  accuracy = evaluator.evaluate(pr)
2  print('Test Accuracy = ', accuracy)
```

▸ (1) Spark Jobs

Test Accuracy =  1.0

Command took 1.17 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:44 on My Cluster

## NAIVE BAYES

Cmd 106

```
1  from pyspark.ml.classification import NaiveBayes
```

Cmd 107

```
1  N=NaiveBayes(featuresCol='features',labelCol='label')
```

Cmd 108

## Test Accuracy of Naive Bayes

```
1  accuracy = evaluator.evaluate(pre1)
2  print('Test Accuracy = ', accuracy)
```

▸ (1) Spark Jobs

Test Accuracy =  0.9991982361194628

Command took 0.99 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 20:12:45 on My Cluster

## Confusion Matrix:

Confusion matrix is a table that is used to define the performance of a classification algorithm. It visualizes and summarizes the performance of a classification algorithm.

## CONFUSION MATRIX

Cmd 116

```
1  from pyspark.mllib.evaluation import MulticlassMetrics
2  from pyspark.sql.types import FloatType
3  import pyspark.sql.functions as F
```

Command took 0.03 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 21:34:43 on My Cluster

Cmd 117

```
1  preds_and_labels = predictions.select(['prediction','label']).withColumn('label', F.col('label').cast(FloatType())).orderBy('prediction')
2  preds_and_labels = preds_and_labels.select(['prediction','label'])
3  metrics = MulticlassMetrics(preds_and_labels.rdd.map(tuple))
4  print('CONFUSION MATRIX',metrics.confusionMatrix().toArray())
```

▸ (4) Spark Jobs

/databricks/spark/python/pyspark/sql/context.py:134: FutureWarning:

Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.

CONFUSION MATRIX [[2729.    0.]
 [   0. 2260.]]

Command took 4.20 seconds -- by f2019332007@umt.edu.pk at 18/07/2022, 21:34:59 on My Cluster

**Deployment on Flask and Herokuapp**