Big Data Programming CSEE5590/490

Module 2 Lab 2

Report

Submitted

By:

Sneha

Mishra:

Class ID: 11

Aditya

Soman

Class ID: 19

Team: 12

Professor: Yugyung Lee

Name: Sneha Mishra

Class ID: 11

Email: smccr@mail.umkc.edu

<u>MyGitHub</u>

Technical Partner: Name: Aditya Soman

Class ID: 19

Email: as9f3@mail.umkc.edu

<u>GitHub</u>

YouTube Link explaining the Lab work can be found here
The report for the Lab work is here
The source code for this lab work can be found here
The available datasets formats can be found here

Objective

Understanding Spark Classification, Spark Streaming and Spark Graphx Task.

Features

- 1. Use of Classification Algorithms such as Naïve Bayes, Decision Tree, Random Forest for attribute classification.
- 2. Report the Confusion matrix, Accuracy based on FMeasure, Precision & Recall for all the algorithms.
- 3. Reason why one of algorithms out performs the rest.
- 4. Perform Word-Count on Twitter Streaming Data using Spark.
- 5. Perform Page Rank on given Dataset.
- 6. State importance of using graphx on the chosen dataset.

Steps:

Part 1: Spark Classification Task

This task contains working on 3 algorithms namely:

1. Naïve Bayes:

It is a classification technique based on Bayes' theorem. Naïve Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Code for the Algorithm:

```
### Section of the first property of the prope
```

```
rb4 ) 🎏 naive_bayes.py :
ot + 🔘 🛊 | 🕸 - 🔄 👸 DecisionTree py 🗴 👸 naive_bayes.ay 🖈 👸 RandomForest.ay
d2_Lab4 [Lab 4] ~/Dow 3
                                 data = assem.transform(data)
 netastore_db
                                 # Split the data into train and test
splits = data.randomSplit([0.6, 0.4], 1234)
spark-warehouse
DecisionTree.py
                                train = splits[0]
test = splits[1]
derby.log
getTweets.pv
                               # create the trainer and set its parameters
nb = NaiveBayes(smoothing=1.8, modelType="gultinomial")
naive_bayes.py
RendomForest.py
twitterstream.py
                                model = nb.fit(train)
atches and Consoles
                                # select example rows to display.
predictions = model.transform(test)
                                 y_true = data.select["DTI"].rdd.flatMap(lastds z: x).collect()
y_pred = data.select["RDW"].rdd.flatMap[lastds x: x).collect()
                                 accuracy = evaluator.evaluate(predictions)
                                 confusionmatrix = confusion_matrix(y_true, y_pred)
                                 precision = precision_score(y_true, y_pred, average='nicro')
                                 print("Naive Bayes - Test set accuracy = " + striaccuracyi)
                                 print("The precision score for Naive Bayes Model is: " + str(precision))
```

Output after running the Algorithm:

```
But Section of the se
```

2. Decision Tree:

Code for the Algorithm:

```
### Comparison of the Comparis
```

```
The problem of the pr
```

Output after running the Algorithm:

```
Substitution of the property o
```

3. Random Tree:

Code for the Algorithm:

```
| Comparison | Com
```

```
b4 🛜 RandomForest.py
文 🔻 🔘 中中位 🏸 🥻 DecisionTree.py 👋 🔓 naive_bayes.py 🔞 🎁 RandomForest.py ...
(2_Lab4 [Lab4] ~/Dow
                                @ Train a Randomforest model.
rf = RandomforestClassifier[labelCol="indexedLabel", featuresCol="indexedFeatures", numTrees=18)
netastore_db
park-warehouse 43
DecisionTree.py 44
45
                                lerby.log
etTweets.py
                                y_true = data.select("DRI").rdd.flatMap(lambda x: x).collect()
y_pred = data.select("BOA").rdd.flatMap(lambda x: x).collect()
ceive_bayes_py
landomForest.py
                                @ Chain indexers and forest in a Pipeline
pipeline = Pipeline[stages=[labelIndexer, featureIndexer, rf, labelConverter]]
rnal Libraries
tches and Consoles
                                 # Train model. This also runs the indexers.
model = pipeline.fit(trainingData)
                                 # Make predictions.
                                 predictions = model.transform(testData)
                                 predictions.select("predictedLabel", "label", "features").show[5]
                                 # Select iprediction, true label| and compute test error
evaluator = MulticlassClassificationEvaluator(
                                              ol="indexedtabel", predictionCol="prediction", metricHarse="accuracy")
                                 accuracy = evaluator.evaluate/predictions)
                                 confusionmatrix = confusion_matrix(y_true, y_pred)
                                 precision = precision_scorely_true, y_pred, average='nicro'l
                                 rfModel = model.stages[2]
print(rfModel) # summary only
                                 print("Random Forest - Test Error = %g" % (1.0 - accuracy!)
                                 print("The precision score for Random Forest Model is: " + str[precision])
                                 print("The recall score for Random Forest Model is: " + str(recall))
```

Output after running the Algorithm:

```
Claim in the confirmation

Claim in the confirmation of the confirmation of
```

State the reasons on why one of algorithms out performs the rest:

The results indicate that the classification accuracy comparison between Naïve Bayes, Random Forest and Decision Trees that Decision Tree has got the highest average accuracy value than the Naïve Bayes and Random Forest but the difference is not statistically significant.

Part 2: Spark Streaming Task

In this task we perform Word-Count on Twitter Streaming Data using Spark. First we get the Twitter data and then we stream it.

Collecting Tweets code:

```
Mcd2_Lab4 [~/Downloads/Mod2_Lab4] - .../get IW
  getTweets.py
    \ominus 😤 | 🏇 - 🎦 🌠 DecisionTree.py 🗵 🌈 naive_bayes.py 🗵 🌈 RandomForest.py 🗵 🌈 getTweets.py 💥
Lab4 [Lab 4] ~/Dow 1
                             from tweepy import OAuthHandler
from tweepy import Stream
:astore_db
                             from tweepy.streaming import StreamListener
rk-warehouse
                             import socket
import json
import time
isionTree.py
by.log
Tweets.py
re_bayes.py
                             consumer_key = 'LTWzE5RIZrW1BL2NWVpmnq9Rd'
                             consumer_secret = '020tURnBZHMdpEfyVpk0YVrsau,B3tJrF8Ww58PYWgxmCAEJJzg'
domForest.py
                             access_taken = '474841965-z6HE681ecahhBcNngLwf4Bb25vJW60Bw9vi2w0N7'
terstream.py
                             access_secret = '4riBdHwe8DLD4Ygsp56aZLT7QWiNIhiBqPXpQfq@PJHkL'
al Libraries
hes and Consoles
                             auth = OAuthHandler(consumer_key, consumer_secret)
                             auth.set_access_token(access_token, access_secret)
                             class TweetsListener(StreamListener):
                                 def __init__(self, csocket):
                                      self.client_socket = csocket
                                 def on_data(self, data):
                                          msg = json.loads(data)
                                          print(msg['text'l.encode('utf-8'))
                                          self.client_socket.send(msg['text'].encode('utf-8'))
                                      return True except BaseException as e:
                                          print("Error on_data: %s" % str(e))
                                  def on_error(self, status):
                                      print(status)
```

```
👸 getTweets.py
                                                                                                                                      RandomForest py × SetTweets.py × Set
               Lab4 [Lab 4] -/Dow 71
                                                                                                                           self.client_socket = csocket
tastore_db
                                                                                                            def on_data(self, data):
irk-warehouse
 dsionTree.py
                                                                                                                                      msg = json.loads(data)
by.log
                                                                                                                                        print(msg['text'].encode('wtf-8'))
                                                                                                                                        self.client_socket.send(msg['text'].encode('utf-8'))
Tweets.py
                                                                                                                          return True
except BaseException as er
ve_bayes.py
idomForest.py
tterstream.py
al Libraries
                                                                                                           def on_error(self, status):
hes and Consoles
                                                                                                                          print(status)
                                                                                               def sendData(c_socket):
                                                                                                            auth = QAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_secret)
                                                                                                            twitter_stream = Stream(auth, TweetsListener(c_socket))
twitter_stream.filter(track=['Football'])
                                                                                                          __nane__ == "__main__":
s = socket.socket() # Create a socket object
host = "localhost" # Get local machine nane
port = 5555 # Reserve a port for your service.
s.bind((host, port)) # Bind to the port
                                                                                               11 __nane_
                                                                                                           print("Listening on port: %s" % str(port))
                                                                                                            s.listen(5) # Now wait for client connection.
                                                                                                            c, addr = s.accept() # Establish connection with client.
                                                                                                            print("Received request from: " + str(addr))
                                                                                                            time.sleep(5)
                                                                                                             sendData(c)
```

Output of Collecting Tweets code:

```
| Part | September | Company | Compa
```

Stream Twitter data code:

```
MODE_EDD-1 [*/DOWNROADS/MODE_EDD-1] * .../EMITOLSTONIE.Dy (EDD-4)
 🖰 twitterstream.py
    \ominus 💠 | 🌣 - 🇺 | 👸 DecisionTree.py × - 👸 naive_bayes.py × - 👸 RandomForest.py × - 👸 getTweets.py × - 📸 twitterstream.py -
                                from pysperk import SperkContext
from pysperk.streaming import StreamingContext
astore_cb
rk-warehouse
isionTree.py
                                from collections import namedtuple
y.log
(weets.py
                                import os
e_bayes.py
                                os.environ["SPARK_HOME"] = "/usr/local/spark/spark=2.3.1-bin=hadoop2.7/"
domForest.py
terstream.py
al Libraries
                                def main():
                                   sc = SparkContext[appName="PysparkStreaming"]
nes and Consoles
                                    wardcount = {}
ssc = StreamingContext(sc, 5)
                                     lines = ssc.socketTextStream("localhost", 5555)
                                     fields = ("word", "count")
Tweet = namedtuple('Text', fields)
                                    # lines = socket_stream.window(20)
counts = lines.flatMap(lambda text: text.split(" ")]\
    .nap(lambda x: (x, 1))\
    .reduceByKey(lambda a, b: a + b).map(lambda rec: Tweet(rec[0], rec[1]))
                                    counts.pprint()
                                     ssc.start()
                                     ssc.awaitTermination()
                                11 __name__ -- "__main__":
                                     main()
```

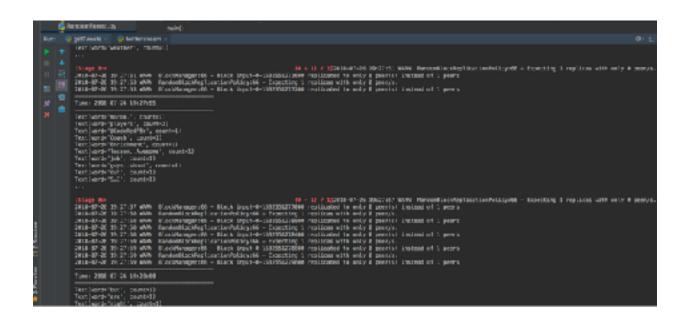
Output of the Twitter Streaming data:

```
### Processor | Processor | Processor |

### Processor | Processor | Processor |

### Processor | Processor | Processor | Processor |

#### Processor | Processor
```



References:

- 1. https://www.linkedin.com/pulse/apache-spark-streaming-twitter-python-laurent-weichberger/
- 2. https://github.com/stefanobaghino/spark-twitter-stream-example

Data-sets provided:

- 1. Absenteeism at work: https://archive.ics.uci.edu/ml/datasets/Absenteeism+at+work
- 2. Immunotherapy Dataset: https://archive.ics.uci.edu/ml/datasets/Immunotherapy+Dataset
- 3. Nashville-meetup Dataset: https://www.kaggle.com/stkbailey/nashville-meetup
- 4. Word Game Dataset: https://www.kaggle.com/anneloes/wordgame
- 5. Cyber Crime Motive: https://www.kaggle.com/sunilkumarsv/indiacybercrimestats2013