CSE – 587       Data Intensive Computing         Spring 2018

# Report

## Lab3: Data Analysis Pipeline Using Apache Spark

Submitted By:
Rahul Mallu (50245594)
Likhith Kumar Miryala (50249280)

# 1. Understand Apache Spark with Titanic Data Analysis

a. We have loaded the titanic data and parsed the RDD to DataFrame.

b. We removed the header from the data and separated each row by "," and converted it to a tuple.

c. Gave names as each column by the header. Figure 1 is the screenshot of the data frame

```
18/05/11 16:43:21 INFO TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool
18/05/11 16:43:21 INFO CodeGenerator: Code generated in 63.610977 ms
+-----------+--------+------+---------+--------------------+------+---+-----+-----+-------------+-------+-----+--------+
|PassengerId|Survived|Pclass|FirstName|                Name|   Sex|Age|SibSp|Parch|       Ticket|   Fare|Cabin|Embarked|
+-----------+--------+------+---------+--------------------+------+---+-----+-----+-------------+-------+-----+--------+
|          1|       0|     3|  "Braund|    Mr. Owen Harris"|  male| 22|    1|    0|    A/5 21171|   7.25|     |       S|
|          2|       1|     1| "Cumings| Mrs. John Bradle...|female| 38|    1|    0|     PC 17599|71.2833|  C85|       C|
|          3|       1|     3|"Heikkinen|        Miss. Laina"|female| 26|    0|    0|STON/O2. 3101282|  7.925|     |       S|
|          4|       1|     1|"Futrelle| Mrs. Jacques Hea...|female| 35|    1|    0|       113803|   53.1| C123|       S|
|          5|       0|     3|  "Allen|  Mr. William Henry"|  male| 35|    0|    0|       373450|   8.05|     |       S|
+-----------+--------+------+---------+--------------------+------+---+-----+-----+-------------+-------+-----+--------+
only showing top 5 rows

root
 |-- PassengerId: string (nullable = true)
 |-- Survived: string (nullable = true)
 |-- Pclass: string (nullable = true)
 |-- FirstName: string (nullable = true)
 |-- Name: string (nullable = true)
 |-- Sex: string (nullable = true)
 |-- Age: string (nullable = true)
 |-- SibSp: string (nullable = true)
 |-- Parch: string (nullable = true)
 |-- Ticket: string (nullable = true)
 |-- Fare: string (nullable = true)
 |-- Cabin: string (nullable = true)
 |-- Embarked: string (nullable = true)
```

Figure: 1

d. After converting to DataFrame, we have manipulated the data by converting some of the features such as age, fare to numeric and Figure 2 is the screenshot of the DataFrame after this manipulation.

```
18/05/11 16:39:59 INFO DAGScheduler: ResultStage 13 (showString at NativeMethodAccessorImpl.java:0) finished in 0.043 s
18/05/11 16:39:59 INFO DAGScheduler: Job 8 finished: showString at NativeMethodAccessorImpl.java:0, took 0.053147 s
18/05/11 16:39:59 INFO CodeGenerator: Code generated in 18.777631 ms
+-------+------+----+
|Age_Sex|female|male|
+-------+------+----+
|   null|    53| 124|
|   24.0|    16|  14|
|   18.0|    13|  13|
|   22.0|    12|  15|
|   30.0|    11|  14|
|   35.0|     8|  10|
|   31.0|     7|  10|
|   36.0|     7|  15|
|   19.0|     7|  18|
|   29.0|     7|  13|
|   28.0|     7|  18|
|   21.0|     7|  17|
|   33.0|     6|   9|
|    2.0|     6|   4|
|   27.0|     6|  12|
|   17.0|     6|   7|
|   45.0|     6|   6|
|   39.0|     6|   8|
|   16.0|     6|  11|
|   40.0|     6|   7|
|   26.0|     5|  13|
|   23.0|     5|  10|
|   25.0|     5|  18|
|   50.0|     5|   5|
|    4.0|     5|   5|
|   38.0|     5|   6|
|   14.0|     4|   2|
|    9.0|     4|   4|
|   48.0|     4|   5|
|   41.0|     4|   2|
|   15.0|     4|   1|
|    5.0|     4|   0|
|   34.0|     4|  11|
|   32.0|     3|  15|
|   44.0|     3|   6|
|   42.0|     3|  10|
|   54.0|     3|   5|
|   58.0|     3|   2|
|   52.0|     2|   4|
|   20.0|     2|  13|
+-------+------+----+
only showing top 40 rows
```
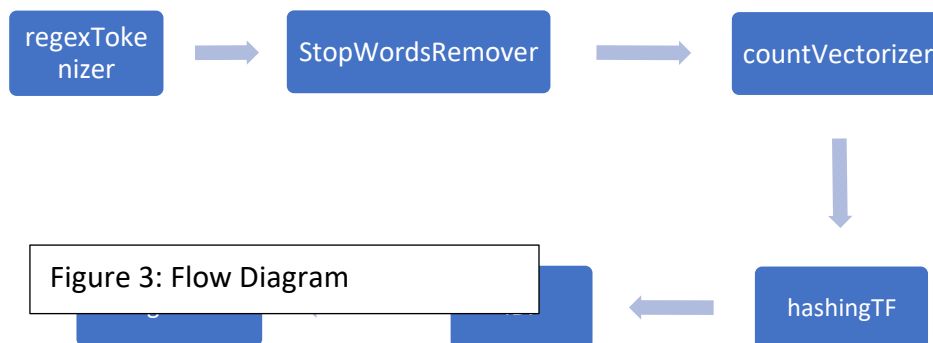
Figure: 2

# 2. Collect and Clean data

a. We have collected the news articles for four categories: Politics, Sports, Business and Technology.
b. We formatted the data to store it in a .csv file so that we can use this for further processing down the line.
c. The further explanations and details have been mentioned in the Jupyter Notebook

# 3. Feature Engineering

a. We have extracted the words characterizing the category and also cleaned the data by removing stop words, punctuations, numbers, etc.
b. This was done by building a model pipeline in Spark which has six stages:
   i. We have used regexTokenizer which takes articles text as input and gives the words as output
   ii. We have created a custom list of stop words and used the StopWordsRemover function to remove those list of stop words
   iii. We have used the countVectorizer function to compute the term frequency
   iv. We have used HashingTF to map words to their frequencies
   v. We have used the IDF function to remove the sparse terms and it gives the TF-IDF
   vi. String Indexer function has been used to encode a string column of labels to a column of label indices. Here we convert the category to label.
c. By the end of this step, we have characterized all the categories

regexTokenizer → StopWordsRemover → countVectorizer

hashingTF

Figure 3: Flow Diagram

# 4. Multi-class Classification

a. We have done the classification using three different classifying algorithms: Logistic Regression, Naive Bayes and Random Forest.
b. We have made predictions and score on the test and training set and calculated the accuracy using multi-class classification evaluator
c. The accuracies obtained for train, test and random data are as follows:

```
Predictions on Trainingset Results:
Logistic Regression Acc: 0.8240890919303049
Naive Bayes Acc: 0.7377569256931171
Random Forest Acc: 0.7577477098758036




Predictions on Testingset Results:
Logistic Regression Acc: 0.4616541353383459
Naive Bayes Acc: 0.42656641604010026
Random Forest Acc: 0.47180451127819545




Final Test Results (New Data):
Logistic Regression Acc: 0.4603174603174603
Naive Bayes Acc: 0.43206349206349204
Random Forest Acc: 0.37333333333333335
```

Figure: 4

d. From the above results, we understand that random forest is robust and versatile but for high dimensional sparse data, logistic regression model is preferable
e. Since there is discrepancy in the accuracy between the training, test and random data and as we have used cross validation we can say that the statistical properties of test, train and random data are different

# 5. Screenshots of Working Program

# 6. References

i. https://datascienceplus.com/multi-class-text-classification-with-pyspark

ii. https://6chaoran.wordpress.com/2016/08/13/__trashed/