**UNIVERSITY OF MISSOURI-KANSAS CITY**

**Bigdata Hadoop Programming**

**# Lab 4 Assignment**

**Team Members:**

Pragathi Thammaneni

Sridevi Mallipudi

**Introduction:**

The main core concept for executing the Lab 4 Assignment is to implement Spark classification task and Spark Streaming Task**.**

**Objectives:**

To code for the 2 questions the below concepts are implemented.

- Implementing the Classification algorithms – Naive Bayes, Decision Tree and Random Forest

- Twitter Streaming Data using Spark

**Approaches /Methods:**

  PySpark , PyCharm

**Workflow &Datasets/Parameters and Evaluation:**

The below each question will follow different approaches to solve. Coding is done to perform the evaluation of each individual snippet to execute the datasets which are provided as the input parameters.

**Question 1:**

1. **Spark Classification Task**

   **Use one of the following datasets**
   - **Absenteeism at work:**
     https://archive.ics.uci.edu/ml/datasets/Absenteeism+at+work

   - **Immunotherapy Dataset:**
     https://archive.ics.uci.edu/ml/datasets/Immunotherapy+Dataset#

   **Perform the following tasks**
   a. **Use the following Classification Algorithms: Naïve Bayes, Decision Tree and Random Forest for the same attribute classification**

   b. **Report the Confusion matrix, Accuracy based on FMeasure, Presion & Recall for all the algorithms**

   c. **State the reasons on why one of algorithms out performs the rest.**

**Solution:**

**Implemented the three classification algorithms -Naïve Bayes, Decision Tree and Random Forest**

**Data Set: Absenteeism**

**Training and testing split ratio :70 and 30**

**Naïve Bayes Classification**

Naive Bayes classifiers are highly scalable, requiring a number of parameters linear in the number of variables (features/predictors) in a learning problem. Maximum-likelihood training can be done by evaluating a closed-form expression which takes linear time, rather than by expensive iterative approximation as used for many other types of classifiers.

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a feature in a class is unrelated to the presence of any other feature.

**Accuracy, precision and recall:**

Accuracy refers to the closeness of a measured value to a standard or known value. Precision refers to the closeness of two or more measurements to each other. Recall referred to as the true positive rate or sensitivity.

Accuracy – 0.1545

Precision – 0.0297

Recall –0.0297

**Output Screens:**

**Decision Tree:**

Decision Tree - Classification. Decision tree builds classification or regression models in the form of a tree structure. The result is a tree with decision nodes and leaf nodes. A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy).
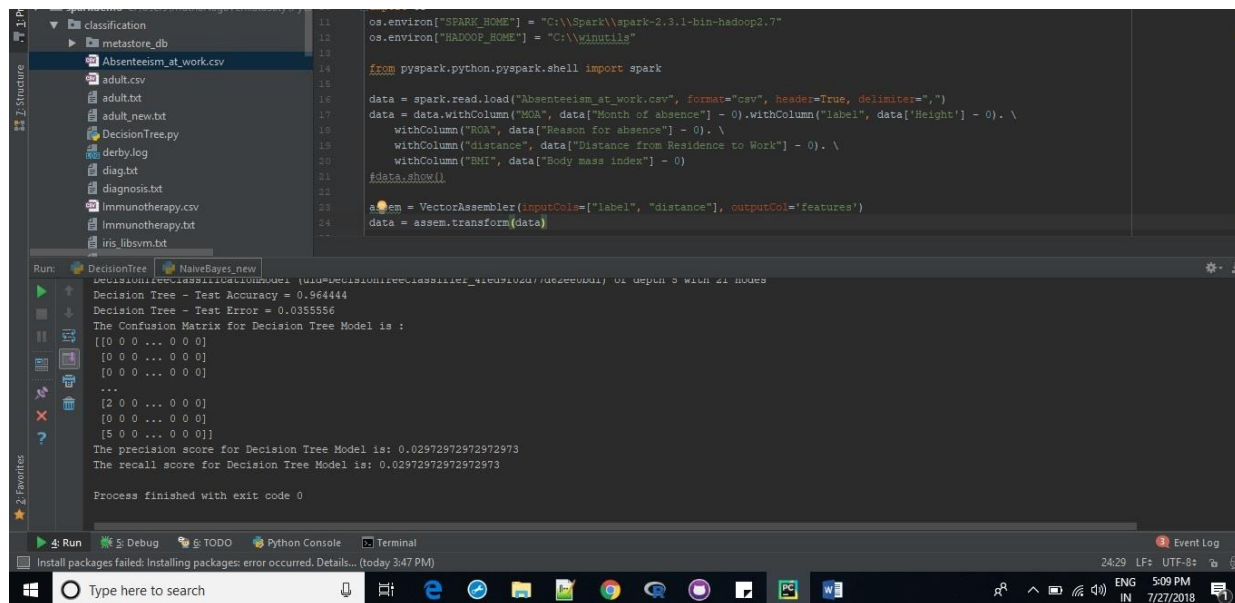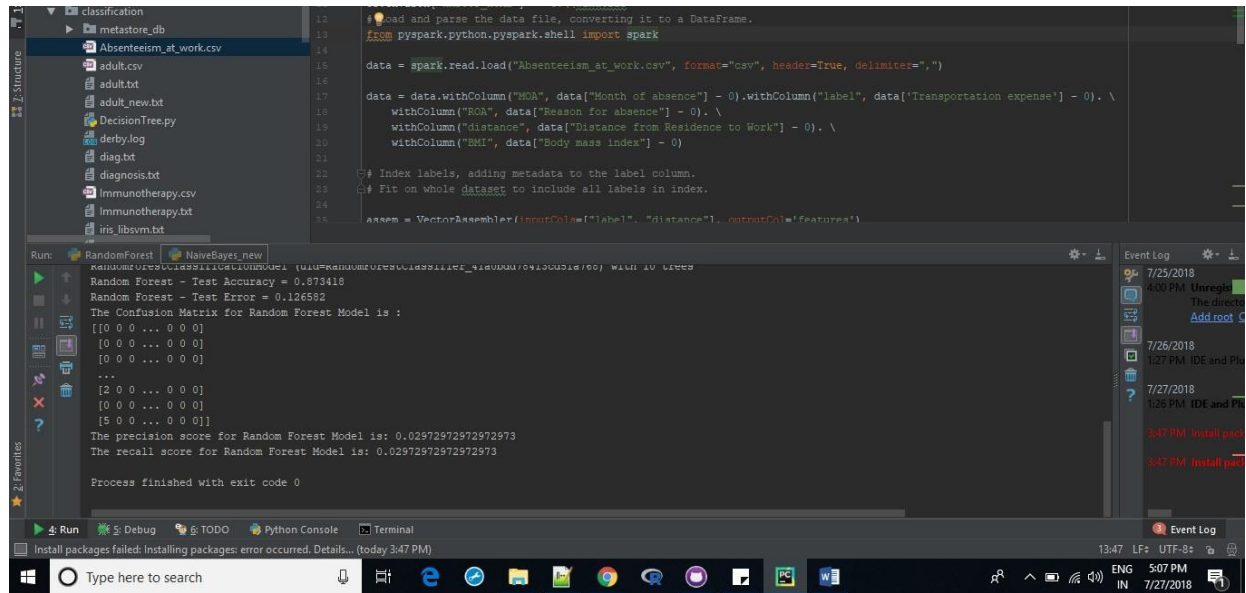
**Accuracy, precision and recall:**

Accuracy refers to the closeness of a measured value to a standard or known value. Precision refers to the closeness of two or more measurements to each other. Recall referred to as the true positive rate or sensitivity.

Accuracy – 0.96444

Precision – 0.0297

Recall –0.0297



**Random Forest:**

Random forest (or random forests) is a trademark term for an ensemble classifier that consists of many decision trees and outputs the class that is the mode of the classes output by individual trees. Random forests are collections of trees, all slightly different. It randomizes the algorithm, not the training data.

**Accuracy, precision and recall:**

Accuracy refers to the closeness of a measured value to a standard or known value. Precision refers to the closeness of two or more measurements to each other. Recall referred to as the true positive rate or sensitivity.

Accuracy – 0.873

Precision – 0.0297

Recall –0.0297

**Output Screen Shots:**



**Final Observations:**

As Random forest (or random forests) is a trademark term for an ensemble classifier that consists of many decision trees and outputs the class that is the mode of the classes output by individual trees. But in this scenario Decision tree out performs with 96% accuracy may be due to "**A single decision tree will train much more quickly, and computing a prediction is also much quicker**" as random forest randomize the algorithm.

**Question 2:**

2. **Spark Streaming Task**

   **Perform Word-Count on Twitter Streaming Data using Spark.**

   https://www.linkedin.com/pulse/apache-spark-streaming-twitter-python-laurent-weichberger/

   https://github.com/stefanobaghino/spark-twitter-stream-example

**Solution:**

Twitter Stream is a simple plugin designed to simply show a user's Twitter timeline. It includes file caching to stop overuse of Twitter's API. You can also choose how many updates to return (maximum of 200). It also includes auto linking for URL's found within the timeline.

- Firstly, we have created a developer account in Twitter using below link. https://apps.twitter.com/
- Below are the variables that contains the user credentials to access Twitter API

  - API ACCESS_TOKEN = " 1329248834-R3X8SrLzchMQyBOSRtV1QSJYfis2lidRex4RXeq"
  - ACCESS_SECRET    = " g6pwJFjRJTQK2cRi01rGV3kpr6C9v7sQNTjrtPLtiETc4"
  - CONSUMER_KEY     = " l82JLsLffFzGMrytXhF4DAnf2"
  - CONSUMER_SECRET = " Fh4985bvDNeF2RurZsseuDYBnbFv3cLmpdyyizFapBkr1uS954"
- We have written python program that is used to fetch tweets in JSON format.
- The tweet data is collected and performed word count on the streaming data
  Spark streaming ,the data can be -ingested – from different sources like twitter and can perform high level complex algorithms like queries   and the processed data can be pushed out to the file systems.To stream the twitter data twitter.utils contains all the built in functionality -to- stream data from twitter



**Code Snippet :**

**Twitter Streaming**

## Wordcount



## Output Screen shot :

**Conclusion:**

As stated, the above workflow describes Spark classification task and Spark Streaming Task.

**Source code link:**

**https://github.com/PragathiThammaneni/Bigdata-Programming--Hadoop-Spark/tree /master/Lab%204/Source%20Code**

**Video Link: Provided in wiki link** https://youtu.be/UNqLOIICUL0

**Wiki Link:**
**https://github.com/PragathiThammaneni/Bigdata-Programming--Hadoop-Spark/wiki/Lab-4-A ssignment**