# Final Term Project – CS634 Data Mining

Instructed by Prof. Jason Wang

| | |
|---|---|
| **First-name:** | Rahul Gautham |
| **Last-name:** | Putcha |
| **NJIT-ID:** | 31524074 |
| **UCID** | RP39 |
| **Email-address:** | rp39@njit.edu |
| **Option:** | 5 |

**<u>Topic on focus</u>**

Data Mining Using Hadoop/Spark in the Cloud

# TABLE OF CONTENT

## List of Tables and Figures

# Chapter 1. General Introduction

With the amount of data increasing at an astounding rate, it has become difficult for people to get useful insights from the large dataset in real-time. Data mining has provided us with a vast set of algorithms that are being used in a variety of applications and have proved to be reliable from time to time.

After years of research, the field of Mathematics have evolved with thirst of automation in computer science. Automation in computer science led us with to new ways to resolve a problem using sophisticated mathematical concepts, one of which is by means of Machine Learning. Machine Learning is a process of making a computer (or a machine) learn a task by using repeated experimentation, with an associated performance for improvement included at every experiment conducted on that task.

Data mining, is a crucial part of machine learning which help a computer find (or mine) useful patterns over a large dataset conquered from various sources. Therefore, Machine learning leverages data mining and computational intelligence algorithms to improve decision making models [3].

## 1.1 The Classification problem

In this project manual we will be focusing one of many and also an important decision-making process, namely the Classification problem.

*Classification is the problem of identifying which of a set of categories (sub-populations) an observation, (or observations) belongs to* [4]*.*

In other word, if you are given few pictures of animals and you are asked to identify what kind of animal it is from a predefined set of known animals, for example: cat, dog, cow, rabbit, giraffe or lion, we are conducting a classification and the problem is popularly known as image classification problem.

In above example, the image is called as a feature and the predefined set (for eg: cow, dog or rabbit) are called labels or targets, of classification. The features need not be an image, it can be collection of attributes with values, such as whether it is an invertebrate or vertebrate, warm or cold blooded, joint or has no legs, etc. Similarly, the label we classify need not be of one type (aka, multi-class classification), but it can be more than one (aka, multi-label classification).

For the remainder of this project, we will be focusing on solving the problem of predicting possibility of Breast-Cancer over a record of patient taken from a hospital in Wisconsin, or the *Breast Cancer Wisconsin Dataset* [1][2]. This dataset is publicly available and the details about the dataset are given in *Chapter 5: Implementation*. In the next section, we will be briefing on the problem for predicting possibility of Breast-Cancer for patients in a hospital at Wisconsin.

## 1.2 Problem Statement

Breast cancer is most common among women in Wisconsin regardless of race. It accounts for nearly one-third of all cancers diagnosed among women [5]. This dataset was first published and is currently publicly available at https://archive.ics.uci.edu/ml/datasets repository. Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image [1].

Following are the list and details of attributes found within the dataset [1]:

1) ID number
2) Diagnosis (M = malignant, B = benign)
3-32) Ten real-valued features are computed for each cell nucleus:

      a) radius (mean of distances from center to points on the perimeter)

      b) texture (standard deviation of gray-scale values)

      c) perimeter

      d) area

      e) smoothness (local variation in radius lengths)

      f) compactness (perimeter^2 / area - 1.0)

      g) concavity (severity of concave portions of the contour)

      h) concave points (number of concave portions of the contour)

      i) symmetry

      j) fractal dimension ("coastline approximation" - 1)

With the details mentioned in above list, we get to know that the labels we are predicting are the outcome of diagnosis results, given in Diagnosis attribute. And the rest of the attributes are all collectively used within our features for predicting the target / label.

## 1.3 What will we be doing and achieving by the end of this manual?

Now that we are familiar with our dataset, lets discuss on what we will be doing for the remainder of the project, i.e., how will we carry out the job of automating prediction of a task by learning a dataset.

Firstly, in the *Chapter 2: Focusing on Algorithm*, we will be discussing on one of the widely used classification algorithm, the Random Forest (RF) Classification algorithm. We will understand the working of an RF and then shift to what we require for carrying the classification task, in *Chapter 3: Requirement Specifications*. In Chapter 3, we will discuss briefly about Apache Spark, a tool used for parallel processing, with its special built-in library for Machine Learning, the Spark MLlib library.

We will also gain hands-on for using EMR (Elastic Map Reduce), a popular platform as a service (PaaS) present of the world-renowned cloud provider, the Amazon Web Service (AWS). This saves us the burden of setting up and installing/configuring software and unknown dependencies, especially Apache Spark and Hadoop HDFS.

In *Chapter 4: Setting up the Environment*, we will be presenting step by step details on setting up the application from scratch with the application source code presented in *Chapter 5: Implementation*.

And lastly, in *Chapter 6: Running the application on Apache Spark inside the EMR (Elastic MapReduce)*, we will guide you through the process of running the application on the cloud both for training and testing over the Breast-Cancer-Wisconsin's dataset.

Therefore, by the end of the project manual we will be building a trained model for predicting the possibility of breast cancer in patients at Wisconsin and present an evaluation metric that estimate the correctness of the model.

***Perquisite:*** *It is assumed that the user has slight familiarity with AWS services and Management Console. Explanation related to AWS services and tool are presented in scope to the Final Term Project.*

# Chapter 2. Focusing on Algorithm

In this chapter we will be focusing on one of the widely used classification algorithm, the *Random Forest (RF) Classification* algorithm and give a detail on how & why Random Forest algorithm efficiently solves classification problem.

## 2.1 Using Random Forest for Classification

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time [6]. In ensemble learning, we combine prediction of two or more models.

In Random Forest, we do this by constructing multiple decision tree over random subsets chosen from the training set. The model then classifies the label, via voting mechanism taken over the results of all random decision tree. Hence, RFs chooses the label that is produced by majority number of decision trees. Although the making of a number of decision tree is a rather overwhelming task, the correctness of model evaluated is more accurate compared to a model trained using other classification algorithms. In the next section we will details on working of Random Forest with an example.

## 2.2 How does it solve the problem?

The key behind random forest problem is a decision tree. Decision trees are basic building block of Random forests. It's a classification algorithm that give the solution by traversing levels involving series of questions. Each question determines conditions facing the observation to determine the kind of answer label. By this we narrow the possible values till we are confident enough to make a single prediction similar to our daily lives.



**Fig 2.1 Decision Tree**

As given in the figure, a tree is constructed over the entire dataset by using the entropy and information gain measure. While testing, we take the observation/sample data and run it over series of question to classify the label. Although a decision tree enforces the consideration of all possible outcomes of a decision and traces each path to a conclusion, they are unstable. This means, a small change in the data can lead to a large change in the structure of the optimal decision tree.

Hence, to resolve this problem we are using Random Forest for classification. In the below figure, we see a sample on how a random forest is generated. Random Forest trains itself by taking random subsets of training set and constructing a number of decision trees over each individual subset.



**Fig 2.2 Random Forest Training**

Similarly, the testing is conducted by grouping the individual resulting labels and aggregating the counts. Finally, the result is picked as whichever label that is picked by the majority number of decision trees, as shown in the figure below. Here, we see that the final result is negative or Malignant label.



**Fig 2.3 Random Forest Testing**

By this we conclude Chapter 2. In the next chapter we will focus on the requirements for the application to classification problem.

# Chapter 3. Requirement Specifications

This chapter presents the basic requirement to run the application for classification using Apache Spark on the cloud. This chapter also highlights all necessary task/requirements that we are going to accomplish in the later chapters.

## 3.1 System Requirements

In order to run the application for a general classification problem, such as Breast Cancer at Wisconsin, the following mentions the hardware and software requirements.

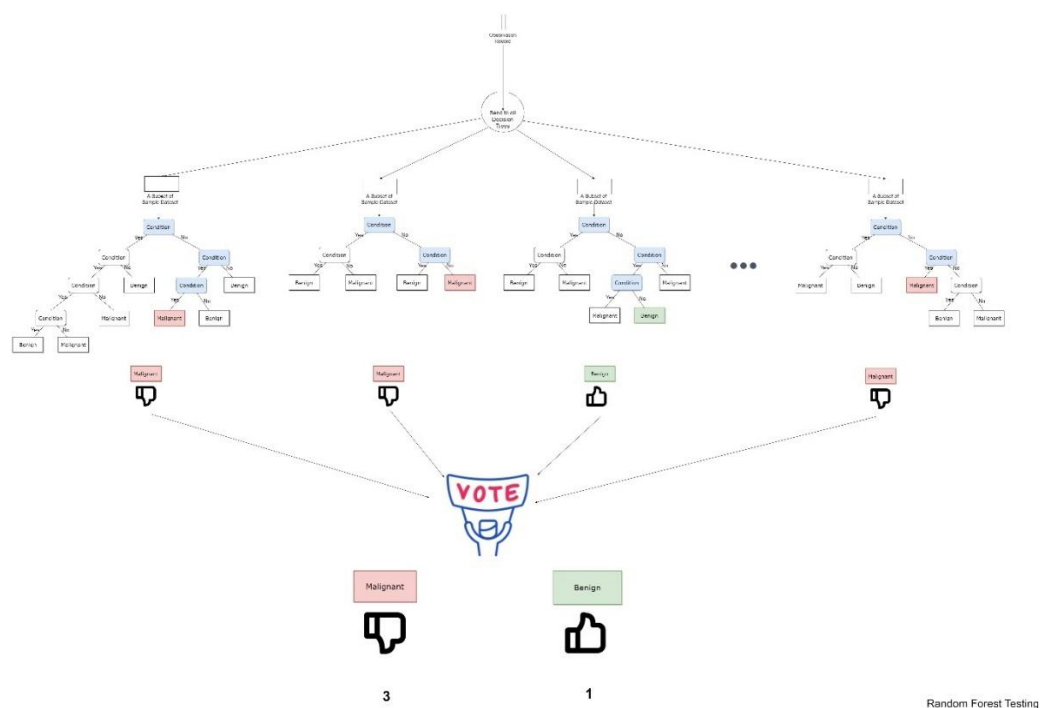| ReqID | Requirements | Details |
|---|---|---|
| Sys-1 | AWS Elastic Map Reduce (EMR) <br> • Apache Spark 2.4.7 <br> • Hadoop 2.10 <br> • Cluster mode launch <br> • 4 EC2 instances used <br> • m5.xlarge instance type | Requires AWS account. Warning: Additional cost involved!! General Linux cluster. 1 Master and 3 Slaves nodes. Spark 2.4.7 on Hadoop 2.10. <br><br> For more details visit, AWS EMR page or please review Chapter 4: Setting up the Environment. |
| Sys-2 | Programming Language: Java | Java v 1.8 (Java JDK8) (Programming Language used) <br> *// Other Language Required: Scala v 2.12 (For Apache Spark)* |
| Sys-3 | S3 (Simple Storage Service) | For storing dataset and final trained model for future predictions. |
| Sys-4 | AWS CLI | AWS CLI for working with Spark on AWS resources in an efficient manner. Requires AWS programmer credential key setup. |
| Sys-5 | SSH-Terminal App | Mobaxterm or Command Prompt with Open SSH installed. |
| Sys-6 | Apache Maven | Apache Maven version 3.8.1 |
| Sys-7 | Apache Spark Core + Spark ML library + Spark SQL Library | spark-core_2.12 v 3.1.2, spark-mllib_2.12 v 3.1.2, spark-sql_2.12 v 3.1.2 |
| Sys-8 | Kaggle Dataset: Breast Cancer Wisconsin's | To download please visit the following link, Breast Cancer Wisconsin. <br> Requires Kaggle account. |

**Tab 3.1 System Requirements**

## 3.2 Functional Requirements

This section presents the requirements that we want our final-build of our application to fulfill. This is specific to the Final Term Project:

| ReqID | Details |
|---|---|
| FR-1 | Setup EMR on AWS to access any AWS resources needed for the project |
| FR-2 | Use Apache Spark for designing ML model, trained parallel over a cluster of 4 nodes (instances). |
| FR-3 | Connect to file storage service for storing and accessing data (for trained model and datasets) |
| FR-4 | Take input (Command Line Argument) from user for taking dataset of choice |
| FR-5 | Automate process for making dataset more efficient by necessary Data Preprocessing. |
| FR-6 | Split the dataset into 2 parts. One used as a training set and another as a testing set. |
| FR-7 | Build a model using a Data Mining algorithm (Random Forest) on any dataset (for example: Breast Cancer Wisconsin dataset) |
| FR-8 | Evaluate model accuracy using Classification Matrix. |

**Tab 3.2 Functional Requirements**

# Chapter 4. Setting up the Environment

This chapter will give a tutorial on setting up infrastructure resources on AWS cloud. The infrastructure we will be working on is Elastic Map Reduce (EMR), which is a popular PaaS for building distributed application on cloud. We will then look into setting of S3, a cloud storage service on AWS, which we will be using for storing our dataset and final trained model for future testing and predictions.

And finally wind up with a general introduction to Apache Spark and Apache Maven which we will be using for setting our application in the next chapter.

## 4.1 Setting up the AWS Account

In order to use AWS resources we are required to setup an AWS account. AWS give us variety ways to setup an account.

## Personal Account

This involves:

- Filling up the form with user details, password credentials and unique account name.
- AWS also requires contact details for billing or promotional contacting in future.
- A personal account with a free-tier is sufficient for this project needs.
- You may also require a credit card, incase you have exceeded your free-tier limit.



**Fig 4.1 Sign-in Page (Step-1)**



**Fig 4.2 Sign-in Page (Step-2)**

**Fig 4.3 Credit-card Page**

(Source: https://docs.sendwithses.com/how/aws-setups/how-to-create-a-free-aws-account)

- After verifying your contact and credit card details, you will land in to the AWS service home page, aka, the Management Console page



**Fig 4.4 AWS Service Home Page (a.k.a Management Console page)**

*Note that the UI changes often depending on AWS satisfaction to the customer in the future, but the access to the resources may remain same to a certain extent as show in this project manual.*

## Student Account

AWS also offers educational account for student studying in the university. Students can access AWS cloud infrastructure without a need for any credit card details. This may require their university email id. For more details please visit, AWS Educate.

## 4.2 Setting up EMR (Elastic MapReduce) on AWS

Let's now focus on setting up the AWS service, Elastic MapReduce (EMR), which will be acting as the base infrastructure on which our application runs. EMR provides a simple way to setup master-slave architecture, another term for a cluster of nodes which commonly referred in distributed systems.

### General Idea about a Cluster

A computer cluster is a set of computers that work together so that they can be viewed as a single system [7]. Theses cluster are connected over a local network with each having their own operating system and follow a common set of protocols with a setup that make it part of the localized network. A cluster consist of a managing entity, or a master, that makes a computer system or worker/slave node part of the network, and it accepts large workload from user. This large workload is split into chunks of sequential and parallel tasks that are later assigned to each of the worker/slave nodes.

This way we can see that any task can be executed efficiently and parallelly among a number of computers in the distributed system. For the project's sake, we are focusing on using this ability of cluster wherein we feed in the dataset to the master to create random forest. This will be done by making individual slaves in the cluster handle the task of creating the individual decision trees from a chosen sample subset of the dataset (Refer to Chapter 2 for more details).

Making the setting up and designing a program for exhibiting such behavior in network of computer systems from scratch by means of a cluster is daunting task. Although we are capable of doing the same using only Hadoop, Apache Spark provides a simple built-in library that resolves the task of setting up a cluster to assigning random decision trees training over worker nodes, via Spark ML library (We will be seeing this in later section of this chapter).
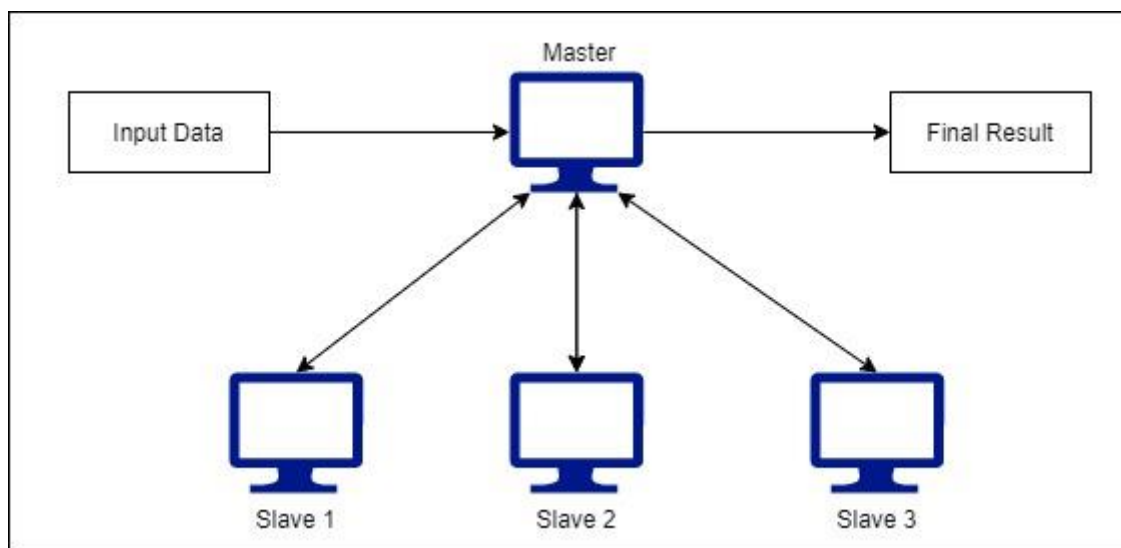


**Fig 4.5 Cluster Network Architecture**

### About EMR and EC2(s)

Amazon EMR is the cloud big data platform for processing vast amounts of data using open-source tools such as Apache Spark, Apache Hive, Apache HBase, Apache Flink, Apache Hudi, and Presto [8]. It automatically makes the setup, operate and scaling of a big data environments on a chosen platform and

tools. This setting may incur cost, but the hourly billing is nominal and we only have to pay for what we use per hour, i.e., pay for each of the EC2 instances we use per hour.

## Creating a Cluster on EMR

Before we create our first EMR cluster, we would require a EC2 access key-pair for securely accessing our individual nodes using our remote terminal of choice. To do this:

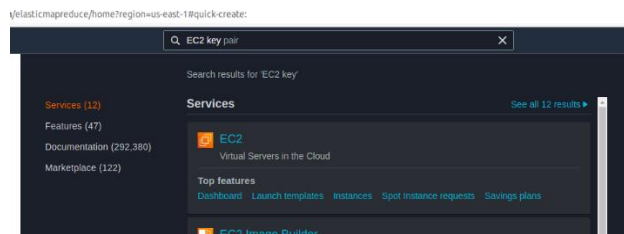- Search for EC2 and click on EC2 service
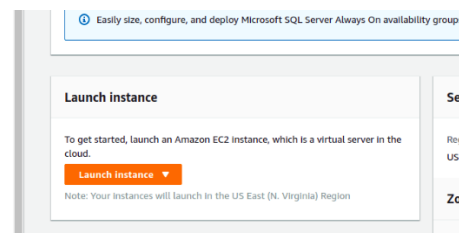


**Fig 4.6 Searching for EC2 service**



**Fig 4.7 Key-Pair Tab option**

- Click on Key-pair option on the left tab menu.
- Fill up the form. If you have a Linux/Mac terminal, please select the .pem and for Windows, please select .ppk option. Finally click on Create key pair button.



**Fig 4.8 Create key pair button**



**Fig 4.9 Create your first key pair**

Your key pair is created and downloaded. Please make sure to save this file, we will be using this in an upcoming section. Now, you are now ready to setup the EMR cluster. Following are the steps to setup an EMR on AWS cloud:

- On management console, use the search bar to search for EMR service and select *EMR.*
- After landing to the EMR page, click on Create Cluster.
- Fill up the form as show in, *Fig 4.8,* and then click on Create cluster.



**Fig 4.10 Searching for EMR service**

**Fig 4.11 EMR service: Home page (For creating a cluster)**



**Fig 4.12 Creating a cluster using EMR service**

After this your cluster is now created and organized using EMR. We will now access the master node using the cluster we have just setup. In order to access our cluster, we first need to define a SSH rule in our EMR-master security groups. Follow steps below to setup the Security group for Master node to access it via SSH.

- After EMR cluster indicates that its running, click on Security group for Master (highlighted in the Fig below)



**Fig 4.13 EMR Post-Cluster creation and Security group for Master Highlighted**

- Select *ElasticMapReduce-master ('s) security group* and proceed to *Inbound rules tab.*
- Click on *Edit Inbound rules button* to see all rule available for a EMR by default.



**Fig 4.14 Elastic MapReduce-master security group**

- Add an SSH rule as shown below in *Fig 4.12* and click save rule.



**Fig 4.15 Adding a SSH rule to EMR-master**

This will allow us to securely access the EMR's EC2 master node using any remote terminal of choice. The next step is to use SSH within our terminal of choice, which we will be discussing in a later section. Before that as we are now still on AWS, lets setup a file storage for holding our dataset and model files.

## 4.3 Setting up a S3 (Simple Storage Service) on AWS

Every Machine Learning (ML) application requires a persistent storage container for dataset file, especially a storage service that provides easy to use API for accessing containing files. AWS S3 provides a simple way to access files on EMR using the resource URI (Uniform Resource Identifier), which is associated with the file/resource we are referring.

- Using the search bar, search for S3 and select *S3 storage service*.
- In the S3 page, click on *Create bucket*.



**Fig 4.16 Searching for S3 service**



**Fig 4.17 Create Bucket button (Top Right-side)**

- This will show us a form. Just fill in the *Bucket name,* must be a globally unique name, if needed adjust the region to your closest or most compliant region of choice.
- Scroll down and click *Create Bucket button*.



**Fig 4.18 Filling Bucket name (must be unique)**
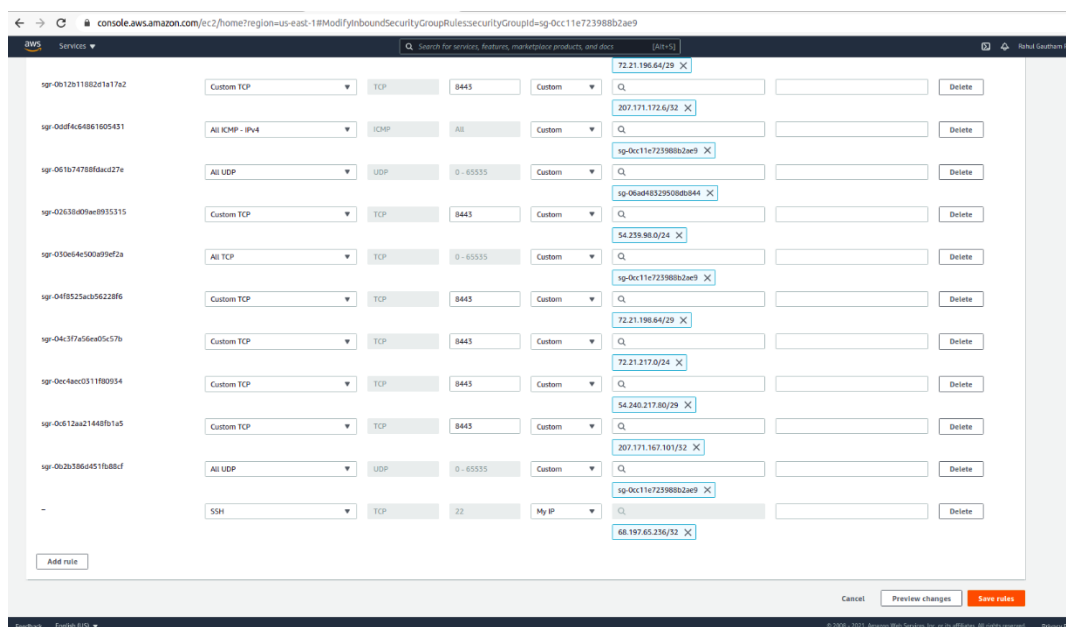


**Fig 4.19 Create Bucket button (in the form)**

This will make a new bucket for us which we can use to store our datasets as we will look at it in detail, *Chapter 5: Implementation.*

## 4.4 Setting up the Application Environment (inside Master node of EMR)

We are now done with setting up services inside management console. We will be focusing on using SFTP and SSH to get the application files loaded into EMR and accessing EMR using a remote terminal or Command Prompt. We will also take a brief look Apache Spark, a popular framework for working on large-scale data using clusters that exhibits parallelism in running workloads. This section also presents a tutorial for setting Apache Maven to use Apache Spark API and then go with the setup of AWS SDK/CLI credentials inside of EMR's Master EC2 node.

**Using SSH to access EMR master:**

We will see how one can use SSH to access EMR master EC2 node system remotely from our computer.

- Since we have already downloaded our key-pair for accessing a EC2, below image show how one can access EMR using the key pair stored in the folder from where the terminal is currently working at. For sake of convenience let name this terminal as, *EMR-master's SSH accessed terminal.*



**Fig 4.20 Accessing EMR master using a Linux/Mac Terminal**

We can also do the same inside Windows, although we may require Putty and PuttyGen tools for privately accessing our EMR master node instance. Follow the tutorial at this link for this purpose.

**Using SFTP for transferring application files to EMR master:**

Now let's focus on transferring the application files to EMR master node. We use SFTP command to do this. We will delve into the folder structure and files in detail in the *Chapter 5: Implementation*.

- Suppose we have the application files in a folder as shown below in *Fig 4.17* and this (project) folder is residing inside of another folder also containing our *AWS access credentials*. Having this pictured, let's proceed with using SSH within the folder containing the AWS access credential and project folder.



**Fig 4.21 Folder containing application files**

**Fig 4.22 Folder containing AWS access credentials**



**Fig 4.23 Transferring file to EMR master node**

- Opening the terminal in the folder similar to the one presented in *Fig 4.18*, use the commands shown in the figure above to access the SFTP. Note that we may have to change the IP-address as allotted by AWS, shown in the AWS EMR page.

If you are not able to find the IP-address, one quick way to find is to visit the EMR cluster page and click on *Connect to the Master Node Using SSH* link. This opens a dialog box, this also shows the IP-address along with the user-name of the EMR-master EC2 node, which are collectively used to access via SFTP or SSH.



**Fig 4.24 Connect to the Master Node Using SSH**

Now, in order for our instances to use any AWS resources using CLI, there are 2 ways to setup either by using Roles or by setting access credentials. Setting by Roles is out of scope for the project, therefore, let's focus on much a simpler way, i.e., by using AWS SDK/CLI access credentials. Please follow the steps below:

- Using the search bar, go to *IAM > Users tab (left tab menu) > Create a user* with at least access to EMR. (For the sake of ease for the project, you can a lot *AdministratorAccess* or role with less privilege that give access to EMR and S3 resources in AWS).

- After user is created, within the user select *Security Credentials tab > Create Access key*
- Download the csv file and proceed to the *EMR-master's SSH accessed terminal. Input the following commands.* For more details in the process, please visit this link.

```
$ mkdir .aws
$ vi ./.aws/credentials
```

Paste your access key credentials presented inside the csv to this file as shown below, and save it.

```
[default]
aws_access_key_id=AKIAIOSFODNN7EXAMPLE
aws_secret_access_key=wJalrXUtnFEMI/K7MDENG/bPxRfiCYEXAMPLEKEY
```

Make sure to substitute you own credentials details for the above. We are done with the configuration.

## About Apache Spark

Apache Spark is an open-source unified analytics engine for large-scale data processing. Spark provides an interface for programming entire clusters with implicit data parallelism and fault tolerance [9]. If we are to build a program using Apache Spark, the easiest way is to use Apache Maven. To install Apache Maven please follow the steps below:

- Within our *EMR-master's SSH accessed terminal,* execute following steps to install maven and setup project files:

```
# Linux-step by step setup for Apache Maven
$ cd ~
$ wget https://mirrors.ocf.berkeley.edu/apache/maven/maven-3/3.6.3/binaries/apache-maven-
3.6.3-bin.tar.gz
$ tar xvf apache-maven-3.6.3-bin.tar.gz
$ sudo mv apache-maven-3.6.3  /usr/local/apache-maven
$ export M2_HOME=/usr/local/apache-maven
$ export M2=$M2_HOME/bin
$ export PATH=$M2:$PATH
$ mvn -version
```

This will install Apache Maven inside of EMR master. The reason we are doing this is to create and organize our application files and finally let it build into a *.jar* file, getting us to producing the full-fledge application. This build is later submitted to Apache Spark, which is already preinstalled due to our choice made previously during creation of EMR cluster.

We are also going to use two popular Apache Spark's built-in libraries, namely Apache Spark ML library and Apache Spark SQL. All of the details of installation of these are already defined in Maven dependency management's pom.xml file supplied in section named *Miscellaneous* of *Chapter 6: Implementation*.

## *General Warning:* Things to NOTE while working on AWS

While working in AWS we must always beware of cost billing associated with resources we use and also those we keep alive after requirement is fulfilled. It is always recommended to terminate, delete or clean up all resource that we have issued on AWS after the requirement and necessities are fulfilled.

Now we proceed to the next chapter which presents the details on setting up Dataset in S3 and demonstrating the code-wise implementations.

# Chapter 5. Implementation

Now that we have done the setup our environment on AWS cloud, we can go into the details on how our application is implemented on *Apache Spark in Java*. First, we will be looking at how we can setup our dataset on S3 and then move on to implementation for parallel machine learning for solving Breast cancer at Wisconsin's problem.

## 5.1 Details on the setting up Dataset in S3

In this section we will setup our files on S3 in a systematic manner.

- Download the dataset for Breast Cancer Wisconsin from Kaggle.
- Go to S3 page in your Management Console
- Open the S3 bucket we have created in, *4.3 Setting up a S3 (Simple Storage Service) on AWS*



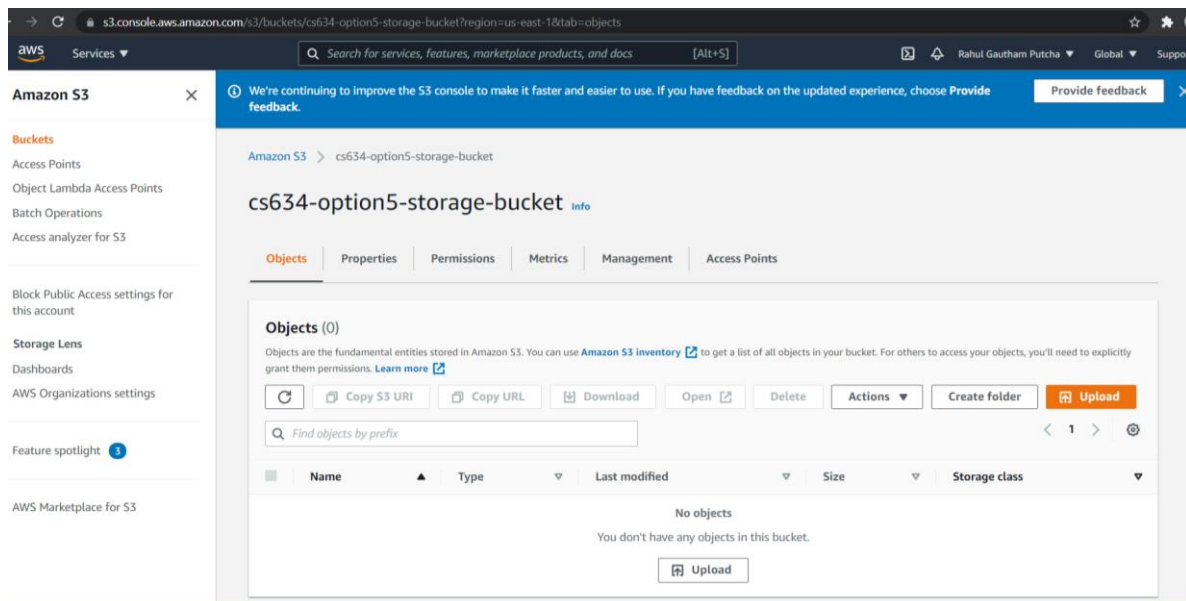**Fig 5.1 Initial Object Listing in S3 Bucket**

- Click on Upload button. An Upload form appears, Click on *Add files* and choose the dataset file.



**Fig 5.2 Uploading a file on S3 (1)**

- Before we submit the form after choosing the dataset, the destination section shows the URI link that will be allotted to the resource we are currently submitting, as shown in *Fig 5.2.*
- Do copy this link as we are going to use it as our argument to our Random Forest Training program that's using Apache Spark.

  Note: A sample link is shown below,
  <div align="center">s3://cs634-option5-storage-bucket</div>
  If we want to access a specific file for this bucket, eg: data.csv file which we have recently uploaded, set the prefix and as above link and suffix as the file name as shown below.
  <div align="center">s3://cs634-option5-storage-bucket/data.csv</div>

- Finally, hit the *Upload button* at the bottom of the form.



**Fig 5.3 Uploading a file on S3 (2)**

## 5.2 Folder Structure

Below figure show the project folder structure details. All the implementation file details are presented in the later sections. Note that *src/main/java/com/njit/rp39/cs634/FinalTermProject* are all individual folders, i.e, src contains main, main contains java, java contains com folder and so on. The original reason for the structure was for maintaining the project as a package, where *com.njit.rp39.cs634.FinalTermProject* is the package our application belongs to. The screenshot below is taken from Visual Studio code.

Since we are done installing with Apache Maven on to EMR from the previous chapter, we will now move on to setting up our project file by build the application structure using files given in the next few sections, all from scratch by referring to the folder structure given below in next page.

**Fig 5.4 Application Folder structure**

After integrating all the file in this chapter, please re-review the section 4.4 Setting up the Application Environment (inside Master node of EMR) from Chapter 4: Setting up the Environment. Below is the first piece of the code, belonging to *App.java,* which is going to be the entry point for our spark program.

**App .java:**

The piece of code below is the entry point for our program, whenever it runs. It gives user an option to decide whether to train/test, run on cluster/thread mode, where the save point for our model and test.csv file be, what dataset we are going to use and what is the label part of our dataset which we want to classify. All of the details above are taken as program argument (i.e., as a command line argument).

```java
package com.njit.rp39.cs634.FinalTermProject;
/**
* Final Term Project
* Author: Rahul Gautham Putcha
* Details:
* - Implement the Data Mining Algorithm on Apache Spark on AWS Elastic MapReduce EMR cluster
**/

import org.apache.log4j.Level;
import org.apache.spark.sql.SparkSession;
import org.apache.log4j.Logger;

/**
* Class: App
* Accepts: Command Line argument for
* Details: Main Application. Conducts Training/Testing
**/
public class App {

  public static void main( String[] args ) {
    SparkSession session;
    Logger.getLogger("org.apache").setLevel(Level.WARN);

    if(args.length==5) {
```

```java
        if(args[1].compareTo("cluster") == 0) {
            session = SparkSession.builder()
                    .appName("Final Term Project: ")
                    .getOrCreate();
        } else {
            session = SparkSession.builder()
                    .appName("Final Term Project: ")
                    .master("local[*]")
                    .getOrCreate();
        }


        if(args[0].compareTo("train") == 0) {
            new ModelTraining(session, args[3], args[2], args[4]);
        } else if(args[0].compareTo("test") == 0) {
            new ModelTesting(session, args[2], args[3], args[4]);
        } else {
            System.out.println(
                "[ERROR]: CLI SYNTAX ERROR\n" +
                "[INFO]: CLI SYNTAX arg: " +
                "{0:train / test} {1:cluster / thread} {2:save/load point location} " +
                "{3:dataset filepath} {4:label}"
            );
        }

    } else{
        System.out.println(
            "[ERROR]: CLI SYNTAX ERROR\n" +
            "[INFO]: CLI SYNTAX arg: " +
            "{0:train / test} {1:cluster / thread} {2:save/load point location} " +
            "{3:dataset filepath} {4:label}"
        );
    }
  }
}
```

**Code 5.1 Code for App.java**

## Data Preprocessing .java

Below code is shared between both ModelTesting.java and ModelTraining.java. This piece of code does all Data preprocessing: Data Cleaning, Data Imputation, Normalization etc., before we send it over to Training or Testing phase.

```java
package com.njit.rp39.cs634.FinalTermProject;

// Import Files (Apache Spark)
import org.apache.spark.ml.feature.Imputer;
import org.apache.spark.ml.feature.Normalizer;
import org.apache.spark.ml.feature.VectorAssembler;

import org.apache.spark.sql.Column;
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import static org.apache.spark.sql.functions.col;


/**
 * Class: Data Preprocessing
 * Details: Contains Common utilities for Data Preprocessing (used in Training and Testing phase)
 **/
public class DataPreprocessing {
    protected Dataset<Row> feature, target;
    protected String[] getFeatureNames(Dataset<Row> data, String label){
```

```java
/**
 * Function: getFeatureNames
 * Argument:
 * Dataset<Row> data: Spark SQL dataset with a Header row
 * String label: Label name present in dataset's header row
 * Details:
 * Parse Header row from a given dataset and disregards label from features.
 * Returns: String[] list_of_features in the dataset `data`
 **/

    String[] columnName = data.schema().fieldNames();
    int j = 0, columnLength = 0;
    for (String s : columnName) {
        if (s.compareTo(label) != 0 && s.charAt(0) != '_' && s.compareTo("id")!=0) {
            columnLength++;
        }
    }
    String[] featureCol = new String[columnLength];
    for (String s : columnName) {
        if (s.compareTo(label) != 0 && s.charAt(0) != '_' && s.compareTo("id")!=0) {
            featureCol[j++] = s;
        }
    }
    return featureCol;
}


protected Dataset<Row> getFeatures(Dataset<Row> data, String label){
    /**
     * Function: getFeatures
     * Argument:
     * Dataset<Row> data: a Dataset with features and label
     * String label: label name present as column header in the dataset `data`
     * Details:
     * - Removes label column inside of dataset `data`
     * Returns: Dataset<Row> feature_columns (...excludes label_column)
     **/
    String[] featureCols = getFeatureNames(data, label);
    VectorAssembler vectorAssembler = new VectorAssembler();
    vectorAssembler.setInputCols(featureCols);
    vectorAssembler.setOutputCol("features");

    Column[] featureColumns = { col("id"), col("features") };

    return vectorAssembler.transform(data).select(featureColumns);
}

protected Dataset<Row> cleanAndImputeData(Dataset<Row> data, String[] columnNames, String label){
    /**
     * Function: cleanAndImputeData
     * Argument:
     * Dataset<Row> data: a Dataset with features and label
     * String[] columnNames: List of columns names
     * String label: label name present as column header in the dataset `data`
     * Details:
     * - Clean dataset by removing columns with over 30% of NULL values and Impute dataset (... on features only)
     * - Hence, excludes `label` column in dataset `data`
     * Returns: Dataset<Row> cleaned_dataset
     **/

    System.out.println("[INFO]: Cleaning data.");

    // Cleaning Data
    for (String s: columnNames) {
        if( data.where(data.col(s).isNull()).count() > data.count()*(30.0f/100.0f) ){
```

```java
            System.out.println("[INFO]: Dropped column "+s+" for >30% null values.");
            data = data.drop(s);
        }
    }

    //Impute Data
    columnNames = getFeatureNames(data,label);
    data = (new Imputer()
            .setInputCols(columnNames)
            .setOutputCols(columnNames)
            .setStrategy("mean")).fit(data).transform(data);

    return data;
}
protected void preprocessData(Dataset<Row> dataset, String label){
    /**
     * Function: preprocessData
     * Argument:
     * Dataset<Row> data: a Dataset with features and label
     * String label: label name present as column header in the dataset `data`
     * Details: (Basic)
     * * 1. Clean Null Population      : Done
     * * 2. Imputation              : Done
     * * 3. Normalization           : Done
     * After Execution: (Result)
     * - Stores `this.feature` variable
     * - Stores `this.target` variable
     **/

    Dataset<Row> cleanedData= cleanAndImputeData(dataset, dataset.schema().fieldNames() ,label);

    // If necessary, drop other
    //for (String s:manualDroppedColumns) { cleanedData = cleanedData.drop(s); }
    //cleanedData.show(5);

    // Getting features from Dataset
    Dataset<Row> features = getFeatures(cleanedData, label);

    // Normalization
    Normalizer normalizer = new Normalizer()
            .setInputCol("features")
            .setOutputCol("normFeatures")
            .setP(1.0);
    Column[] featureColumns = { col("id"), col("normFeatures") };
    Dataset<Row> normData = normalizer.transform(features).select(featureColumns);

    // Storing individual feature and target to class member vaiables
    Column[] targetColumns = { col("id"), col(label) };
    feature = normData;
    target = dataset.select(targetColumns);

}
}
```
**Code 5.2 Code for DataPreprocessing.java**

## 5.3 Code for Random Forest (RF) Classification training on Apache Spark

Below is the code used for the purpose of training our model. This code loads the dataset and splits the dataset into 2 subsets, Training and testing. The testing dataset is saved at the save point entered as argument into our App and the training dataset is used for building our model. Before we train the model using the training dataset, we preprocess the data using function inherited from DataPreprocessing class.

After the training is complete, we save the model to the same save point location where we are storing the test file.

**ModelTraining.java**

```java
package com.njit.rp39.cs634.FinalTermProject;

// Import Files (Apache Spark)
import org.apache.spark.ml.Pipeline;
import org.apache.spark.ml.PipelineStage;
import org.apache.spark.ml.classification.RandomForestClassifier;
import org.apache.spark.ml.evaluation.BinaryClassificationEvaluator;
import org.apache.spark.ml.feature.StringIndexer;
import org.apache.spark.ml.param.ParamMap;
import org.apache.spark.ml.tuning.CrossValidator;
import org.apache.spark.ml.tuning.CrossValidatorModel;
import org.apache.spark.ml.tuning.ParamGridBuilder;

import org.apache.spark.sql.SparkSession;
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import static org.apache.spark.sql.functions.col;
import org.apache.spark.sql.Column;

import java.io.IOException;


/**
 * Class: Model Training
 * Details: Conducts Training on Dataset using Random Forest Model using Apache Spark
 **/
public class ModelTraining extends DataPreprocessing{
    String[] manualDroppedColumns;

    public ModelTraining(SparkSession session, String datasetLoc, String saveLoc, String label){
        /**
         * Function: ModelTraining (Constructor)
         * Argument:
         * SparkSession session: A spark session for working with Spark utilities
         * String datasetLoc: Location of the initial dataset (containing Training+Testing dataset)
         * String saveLoc: A Save point location for saving the test dataset (test.csv) and final trained model (rf_model.model)
         * String label: Label name belonging to the dataset present at location `datasetLoc`
         * Details:
         * - Loads dataset into spark,
         * - Splits dataset into random train & test,
         * - Conducts Data-Preprocessing, Normalization
         * - Trains model
         * * All tasks are conducted over cluster with parallelism using Apache Spark
         * Dataset requirement (Assumptions):
         * - Works on any dataset with a mandatory `id` column
         * - must have at least one feature columns (other than `id`)
         * - must have at least one label column
         **/

        // Manual Editing columns
        manualDroppedColumns = new String[] {};

        // Fetching Entire Dataset/CSV file from
        Dataset<Row> data = session.read()
            .option("header", true)
            .option("inferSchema", true)
            .csv(datasetLoc);
        data.show(5);
```

```java
    // Splitting Dataset to Train-Test set Split here ...
    Dataset<Row>[] splits = data.randomSplit(new double[] {0.8, 0.2}, 12345);
    Dataset<Row> train = splits[0];
    Dataset<Row> test = splits[1];

    // Saving Testing set to use it later while Running Program with TestingModel.java
    try {
        test.write().format("csv")
            .option("inferSchema", "true")
            .option("header", "true")
            .csv(saveLoc+"/test.csv");
    }catch(Exception e){
        System.out.println("[Error]: Couldn't save the test split file to location - "+saveLoc+"/test.csv");
        System.out.println(e.getMessage());
    }

    // Data Preprocessing (Basic)
    preprocessData(train, label);

    System.out.println("Training set:\nFeatures:");
    feature.show(3);

    System.out.println("Target(Labels):");
    target.show(3);

    // Splitting Dataset to into Features and Labels: Generating Training Dataset
    Column[] colNames = {col("normFeatures"), col("diagnosis")};
    Dataset<Row> train_set = feature
        .join(target, feature.col("id").equalTo( target.col("id") ))
        .select(colNames)
        .withColumnRenamed("normFeatures","features")
        .withColumnRenamed(label,"label");
    StringIndexer stringIndexer = new StringIndexer().setInputCol("label").setOutputCol("labelIndex");

    Column[] trainCol = {col("features"), col("labelIndex")};
    train_set = stringIndexer.fit(train_set)
        .transform(train_set).select(trainCol)
        .withColumnRenamed("labelIndex", "label");

    System.out.println("Final training set:");
    train_set.show(3);

    // Start Training a model on the Training Dataset
    train(train_set, saveLoc);
}

private void train(Dataset<Row> train, String saveLocation){
    /**
     * Function: train
     * Argument:
     * Dataset<Row> data: a Dataset with features and label
     * String saveLocation: save point for saving the model
     * Details:
     * 1. Model Selection: Random Forest (RF) Classifier
     * 2. Model Training and Cross Validation
     * 4. Best Model chosen using areaUnderROC metric
     * 3. Save Model for future testing and predictions
     * After Execution: (Result)
     * - Best Model will be saved at save point `saveLocation`
     **/

    // Training
    RandomForestClassifier rf = new RandomForestClassifier();

    //Iterable<String> impurity = Arrays.asList(new String[]{"gini", "entropy"});
```

```java
    Pipeline pipeline = new Pipeline().setStages(new PipelineStage[] {rf});
    ParamMap[] paramMap = new ParamGridBuilder()
        .addGrid(rf.maxDepth(), new int[]{1,2,5, 10, 15})
        .addGrid(rf.minInstancesPerNode(), new int[]{1, 2, 4,5,10})
        .build();

    // We now treat the Pipeline as an Estimator, wrapping it in a CrossValidator instance.
    // A CrossValidator requires an Estimator, a set of Estimator ParamMaps, and an Evaluator.
    // Note that the evaluator here is a BinaryClassificationEvaluator and its default metric
    // is areaUnderROC.
    CrossValidator cv = new CrossValidator()
        .setEstimator(pipeline)
        .setEstimatorParamMaps(paramMap)
        .setEvaluator(new BinaryClassificationEvaluator())
        .setNumFolds(10) // 10-Fold validation
        .setParallelism(4); // 4 setting in parallel

    // Run cross-validation, and choose the best set of parameters.
    CrossValidatorModel cvModel = cv.fit(train);

    try {
       cvModel.save(saveLocation+"/rf_model.model");
       System.out.println("[INFO]: Model Trained Successfully");
    }catch(IOException e){
       System.out.println("[Error]: Couldn't Write model to location - "+saveLocation+"/rf_model.model");
    }

  }
}
```

<div align="center">Code 5.3 Code for TrainingModel.java</div>

## 5.4 Code for Testing the RF model for solving Breast-Cancer Problem

Below is the code use for testing purpose. Here we take in location of the model we have saved at end of training phase, aka, load point, hence we load the model in to our application and then make Spark use it to predict and test over the testing dataset. We also accept the testing dataset location as an individual argument separately. By end of the test predication conducted, we showcase the model performance by conducting evaluation using *F1-score*, *Accuracy score* (internally performed by means of using *Confusion Matrix*) and lastly, by using *Area under ROC* metric.

**ModelTesting.java**

```java
package com.njit.rp39.cs634.FinalTermProject;

// Import Files (Apache Spark)
import org.apache.spark.ml.evaluation.BinaryClassificationEvaluator;
import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator;
import org.apache.spark.ml.feature.StringIndexer;
import org.apache.spark.ml.tuning.CrossValidatorModel;

import org.apache.spark.sql.SparkSession;
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import static org.apache.spark.sql.functions.col;
import org.apache.spark.sql.Column;




/**
* Class: Model Testing
```

```java
 * Details: Conducts Testing on Dataset using Random Forest Model using Apache Spark
 **/
public class ModelTesting extends DataPreprocessing{
  String[] manualDroppedColumns;

  public ModelTesting(SparkSession session, String loadLoc, String datasetLoc , String label){
    /**
     * Function: ModelTesting (Constructor)
     * Argument:
     * SparkSession session: A spark session for working with Spark utilities
     * String datasetLoc: Location of the test dataset (Testing dataset only)
     * String loadLoc: A Load point location holding the final trained model (rf_model.model)
     * String label: Label name belonging to the dataset present at location `datasetLoc`
     * Details:
     * - Loads test dataset into spark,
     * - Conducts Data-Preprocessing, Normalization
     * - Test and Evaluates the model
     * * All tasks are conducted over cluster with parallelism using Apache Spark
     * Dataset requirement (Assumptions):
     * - Works on any dataset with a mandatory `id` column
     * - must have at least one feature columns (other than `id`)
     * - must have at least one label column
     **/

    System.out.println("loadLoc:"+loadLoc);
    System.out.println("datasetLoc:"+datasetLoc);
    // Manual Editing columns
    manualDroppedColumns = new String[] {};

    // Fetching Entire Dataset/CSV file from
    Dataset<Row> test = session.read()
        .option("header", true)
        .option("inferSchema", true)
        .csv(datasetLoc);
    test.show(5);


    // Data Preprocessing (Basic)
    preprocessData(test, label);

    System.out.println("Testing set:\nFeatures:");
    feature.show(3);

    System.out.println("Target(Labels):");
    target.show(3);

    // Splitting Dataset to into Features and Labels: Generating Training Dataset
    Column[] colNames = {col("normFeatures"), col("diagnosis")};
    Dataset<Row> test_set = feature
        .join(target, feature.col("id").equalTo( target.col("id") ))
        .select(colNames)
        .withColumnRenamed("normFeatures","features")
        .withColumnRenamed(label,"label");
    StringIndexer stringIndexer = new StringIndexer().setInputCol("label").setOutputCol("labelIndex");

    Column[] testCol = {col("features"), col("labelIndex")};
    test_set = stringIndexer.fit(test_set)
        .transform(test_set).select(testCol)
        .withColumnRenamed("labelIndex", "label");

    System.out.println("Final Test set (After preprocessing):");
    test_set.show(3);

    // Start Testing a model on the Testing Dataset
    test(test_set, loadLoc);
```

```java
    }

    private void test(Dataset<Row> test, String loadLocation){
        /**
         * Function: test
         * Argument:
         * Dataset<Row> test: a Dataset with features and label
         * String loadLocation: load point for fetching the pre-trained RF model
         * Details:
         * 1. Loads the Best Pre-Trained model
         * 2. Does Prediction on the Test dataset
         * 4. Evaluates using different metrics:
         * * - Accuracy
         * * - F1-score, and
         * * - Area under ROC metric
         **/

        // Testing
        CrossValidatorModel cvModel = CrossValidatorModel.load(loadLocation+"/rf_model.model");

        // Make predictions on test documents. cvModel uses the best model found.
        Dataset<Row> predictions = cvModel.transform(test);

        MulticlassClassificationEvaluator evaluator = new MulticlassClassificationEvaluator()
                .setLabelCol("label")
                .setPredictionCol("prediction")
                .setMetricName("accuracy");
        System.out.println("Accuracy: "+evaluator.evaluate(predictions));
        evaluator = new MulticlassClassificationEvaluator()
                .setLabelCol("label")
                .setPredictionCol("prediction")
                .setMetricName("f1");
        System.out.println("F1-score: "+evaluator.evaluate(predictions));

        System.out.println("Area under ROC metric: "+new BinaryClassificationEvaluator().evaluate(predictions));

    }
}
```

**Code 5.4 Code for TestingModel.java**

## 5.5 Miscellaneous

Below code is also a crucial part of our application as it is used to perform the build using Maven and ultimately used to create the functional & independent .jar file, which we use to submit jobs to Apache Spark. The pom.xml file showcases the dependencies that our application need, to properly handle the execution. Apache Maven use this file to install all dependencies from the popular Maven repository.

**pom.xml**

```xml
<?xml version="1.0" encoding="UTF-8"?>

<project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd/maven-4.0.0.xsd">
    <modelVersion>4.0.0</modelVersion>

    <groupId>com.njit.rp39.cs634.FinalTermProject</groupId>
    <artifactId>cs634-ftp-opt5</artifactId>
    <version>1.0-SNAPSHOT</version>

    <name>cs634-ftp-opt5</name>
```

```xml
<properties>
  <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
  <maven.compiler.source>1.7</maven.compiler.source>
  <maven.compiler.target>1.7</maven.compiler.target>
</properties>

<dependencies>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_2.12</artifactId>
    <version>3.1.0</version>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-sql_2.12</artifactId>
    <version>3.1.0</version>
  </dependency>
  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-hdfs</artifactId>
    <version>3.2.1</version>
  </dependency>
  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-core</artifactId>
    <version>1.2.1</version>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-mllib_2.12</artifactId>
    <version>3.1.2</version>
  </dependency>
  <dependency>
    <groupId>junit</groupId>
    <artifactId>junit</artifactId>
    <version>4.11</version>
    <scope>test</scope>
  </dependency>
</dependencies>

<build>
  <pluginManagement><!-- lock down plugins versions to avoid using Maven defaults (may be moved to parent pom) -->
    <plugins>
      <!-- clean lifecycle, see https://maven.apache.org/ref/current/maven-core/lifecycles.html#clean_Lifecycle -->
      <plugin>
        <artifactId>maven-clean-plugin</artifactId>
        <version>3.1.0</version>
      </plugin>
      <!-- default lifecycle, jar packaging: see https://maven.apache.org/ref/current/maven-core/default-
bindings.html#Plugin_bindings_for_jar_packaging -->
      <plugin>
        <artifactId>maven-resources-plugin</artifactId>
        <version>3.0.2</version>
      </plugin>
      <plugin>
        <artifactId>maven-compiler-plugin</artifactId>
        <version>3.8.0</version>
      </plugin>
      <plugin>
        <artifactId>maven-surefire-plugin</artifactId>
        <version>2.22.1</version>
      </plugin>
      <plugin>
        <artifactId>maven-jar-plugin</artifactId>
```

```xml
      <version>3.0.2</version>
    </plugin>
    <plugin>
     <artifactId>maven-install-plugin</artifactId>
     <version>2.5.2</version>
    </plugin>
    <plugin>
     <artifactId>maven-deploy-plugin</artifactId>
     <version>2.8.2</version>
    </plugin>
    <!-- site lifecycle, see https://maven.apache.org/ref/current/maven-core/lifecycles.html#site_Lifecycle -->
    <plugin>
     <artifactId>maven-site-plugin</artifactId>
     <version>3.7.1</version>
    </plugin>
    <plugin>
     <artifactId>maven-project-info-reports-plugin</artifactId>
     <version>3.0.0</version>
    </plugin>
   </plugins>
  </pluginManagement>
  <plugins>
   <plugin>
    <groupId>org.apache.maven.plugins</groupId>
    <artifactId>maven-compiler-plugin</artifactId>
    <configuration>
     <source>8</source>
     <target>8</target>
    </configuration>
   </plugin>
  </plugins>
 </build>
</project>
```

**Code 5.5 Code for pom.xml**

We now also look at the code for Random Forest Classifier that Spark uses to build the Random Forest Classifier model. As described in the *Chapter 2: Focusing on the algorithms,*

```
/*
 * Licensed to the Apache Software Foundation (ASF) under one or more
 * contributor license agreements.  See the NOTICE file distributed with
 * this work for additional information regarding copyright ownership.
 * The ASF licenses this file to You under the Apache License, Version 2.0
 * (the "License"); you may not use this file except in compliance with
 * the License.  You may obtain a copy of the License at
 *
 *    http://www.apache.org/licenses/LICENSE-2.0
 *
 * Unless required by applicable law or agreed to in writing, software
 * distributed under the License is distributed on an "AS IS" BASIS,
 * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
 * See the License for the specific language governing permissions and
 * limitations under the License.
 */

class RandomForestClassifier @Since("1.4.0") (
  @Since("1.4.0") override val uid: String)
  extends ProbabilisticClassifier[Vector, RandomForestClassifier, RandomForestClassificationModel]
  with RandomForestClassifierParams with DefaultParamsWritable {
  @Since("1.4.0")
  def this() = this(Identifiable.randomUID("rfc"))
```

```
// Override parameter setters from parent trait for Java API compatibility.
// Parameters from TreeClassifierParams:
/** @group setParam */
@Since("1.4.0")
def setMaxDepth(value: Int): this.type = set(maxDepth, value)
/** @group setParam */
@Since("1.4.0")
def setMaxBins(value: Int): this.type = set(maxBins, value)
/** @group setParam */
@Since("1.4.0")
def setMinInstancesPerNode(value: Int): this.type = set(minInstancesPerNode, value)
/** @group setParam */
@Since("3.0.0")
def setMinWeightFractionPerNode(value: Double): this.type = set(minWeightFractionPerNode, value)
/** @group setParam */
@Since("1.4.0")
def setMinInfoGain(value: Double): this.type = set(minInfoGain, value)
/** @group expertSetParam */
@Since("1.4.0")
def setMaxMemoryInMB(value: Int): this.type = set(maxMemoryInMB, value)
/** @group expertSetParam */
@Since("1.4.0")
def setCacheNodeIds(value: Boolean): this.type = set(cacheNodeIds, value)
/**
 * Specifies how often to checkpoint the cached node IDs.
 * E.g. 10 means that the cache will get checkpointed every 10 iterations.
 * This is only used if cacheNodeIds is true and if the checkpoint directory is set in
 * [[org.apache.spark.SparkContext]].
 * Must be at least 1.
 * (default = 10)
 * @group setParam
 */
@Since("1.4.0")
def setCheckpointInterval(value: Int): this.type = set(checkpointInterval, value)
/** @group setParam */
@Since("1.4.0")
def setImpurity(value: String): this.type = set(impurity, value)
// Parameters from TreeEnsembleParams:
/** @group setParam */
@Since("1.4.0")
def setSubsamplingRate(value: Double): this.type = set(subsamplingRate, value)
/** @group setParam */
@Since("1.4.0")
def setSeed(value: Long): this.type = set(seed, value)
// Parameters from RandomForestParams:
/** @group setParam */
@Since("1.4.0")
def setNumTrees(value: Int): this.type = set(numTrees, value)
/** @group setParam */
@Since("3.0.0")
def setBootstrap(value: Boolean): this.type = set(bootstrap, value)
/** @group setParam */
@Since("1.4.0")
def setFeatureSubsetStrategy(value: String): this.type =
  set(featureSubsetStrategy, value)

/**
 * Sets the value of param [[weightCol]].
 * If this is not set or empty, we treat all instance weights as 1.0.
```

```
 * By default the weightCol is not set, so all instances have weight 1.0.
 *
 * @group setParam
 */
@Since("3.0.0")
def setWeightCol(value: String): this.type = set(weightCol, value)

override protected def train(
    dataset: Dataset[_]): RandomForestClassificationModel = instrumented { instr =>
  instr.logPipelineStage(this)
  instr.logDataset(dataset)
  val categoricalFeatures: Map[Int, Int] =
    MetadataUtils.getCategoricalFeatures(dataset.schema($(featuresCol)))
  val numClasses: Int = getNumClasses(dataset)

  if (isDefined(thresholds)) {
    require($(thresholds).length == numClasses, this.getClass.getSimpleName +
      ".train() called with non-matching numClasses and thresholds.length." +
      s" numClasses=$numClasses, but thresholds has length ${$(thresholds).length}")
  }

  val instances = extractInstances(dataset, numClasses)
  val strategy =
    super.getOldStrategy(categoricalFeatures, numClasses, OldAlgo.Classification, getOldImpurity)
  strategy.bootstrap = $(bootstrap)

  instr.logParams(this, labelCol, featuresCol, weightCol, predictionCol, probabilityCol,
    rawPredictionCol, leafCol, impurity, numTrees, featureSubsetStrategy, maxDepth, maxBins,
    maxMemoryInMB, minInfoGain, minInstancesPerNode, minWeightFractionPerNode, seed,
    subsamplingRate, thresholds, cacheNodeIds, checkpointInterval, bootstrap)

  val trees = RandomForest
    .run(instances, strategy, getNumTrees, getFeatureSubsetStrategy, getSeed, Some(instr))
    .map(_.asInstanceOf[DecisionTreeClassificationModel])
  trees.foreach(copyValues(_))

  val numFeatures = trees.head.numFeatures
  instr.logNumClasses(numClasses)
  instr.logNumFeatures(numFeatures)
  createModel(dataset, trees, numFeatures, numClasses)
}

private def createModel(
    dataset: Dataset[_],
    trees: Array[DecisionTreeClassificationModel],
    numFeatures: Int,
    numClasses: Int): RandomForestClassificationModel = {
  val model = copyValues(new RandomForestClassificationModel(uid, trees, numFeatures, numClasses))
  val weightColName = if (!isDefined(weightCol)) "weightCol" else $(weightCol)

  val (summaryModel, probabilityColName, predictionColName) = model.findSummaryModel()
  val rfSummary = if (numClasses <= 2) {
    new BinaryRandomForestClassificationTrainingSummaryImpl(
      summaryModel.transform(dataset),
      probabilityColName,
      predictionColName
$(labelCol),
      weightColName,
      Array(0.0))
  } else {
```

```scala
    new RandomForestClassificationTrainingSummaryImpl(
      summaryModel.transform(dataset),
      predictionColName,
      $(labelCol),
      weightColName,
      Array(0.0))
  }
  model.setSummary(Some(rfSummary))
}

@Since("1.4.1")
override def copy(extra: ParamMap): RandomForestClassifier = defaultCopy(extra)
}

@Since("1.4.0")
object RandomForestClassifier extends DefaultParamsReadable[RandomForestClassifier] {
 /** Accessor for supported impurity settings: entropy, gini */
 @Since("1.4.0")
 final val supportedImpurities: Array[String] = TreeClassifierParams.supportedImpurities

 /** Accessor for supported featureSubsetStrategy settings: auto, all, onethird, sqrt, log2 */
 @Since("1.4.0")
 final val supportedFeatureSubsetStrategies: Array[String] =
   TreeEnsembleParams.supportedFeatureSubsetStrategies

 @Since("2.0.0")
 override def load(path: String): RandomForestClassifier = super.load(path)
}

/**
 * <a href="http://en.wikipedia.org/wiki/Random_forest">Random Forest</a> model for classification.
 * It supports both binary and multiclass labels, as well as both continuous and categorical
 * features.
 *
 * @param _trees  Decision trees in the ensemble.
 *                Warning: These have null parents.
 */
@Since("1.4.0")
class RandomForestClassificationModel private[ml] (
    @Since("1.5.0") override val uid: String,
    private val _trees: Array[DecisionTreeClassificationModel],
    @Since("1.6.0") override val numFeatures: Int,
    @Since("1.5.0") override val numClasses: Int)
  extends ProbabilisticClassificationModel[Vector, RandomForestClassificationModel]
  with RandomForestClassifierParams with TreeEnsembleModel[DecisionTreeClassificationModel]
  with MLWritable with Serializable
  with HasTrainingSummary[RandomForestClassificationTrainingSummary] {

  require(_trees.nonEmpty, "RandomForestClassificationModel requires at least 1 tree.")

  /**
   * Construct a random forest classification model, with all trees weighted equally.
   *
   * @param trees  Component trees
   */
  private[ml] def this(
    trees: Array[DecisionTreeClassificationModel],
    numFeatures: Int,
    numClasses: Int) =
   this(Identifiable.randomUID("rfc"), trees, numFeatures, numClasses)
```

```scala
@Since("1.4.0")
override def trees: Array[DecisionTreeClassificationModel] = _trees

// Note: We may add support for weights (based on tree performance) later on.
private lazy val _treeWeights: Array[Double] = Array.fill[Double](_trees.length)(1.0)

@Since("1.4.0")
override def treeWeights: Array[Double] = _treeWeights

/**
 * Gets summary of model on training set. An exception is thrown
 * if `hasSummary` is false.
 */
@Since("3.1.0")
override def summary: RandomForestClassificationTrainingSummary = super.summary

/**
 * Gets summary of model on training set. An exception is thrown
 * if `hasSummary` is false or it is a multiclass model.
 */
@Since("3.1.0")
def binarySummary: BinaryRandomForestClassificationTrainingSummary = summary match {
  case b: BinaryRandomForestClassificationTrainingSummary => b
  case _ =>
    throw new RuntimeException("Cannot create a binary summary for a non-binary model" +
      s"(numClasses=${numClasses}), use summary instead.")
}

/**
 * Evaluates the model on a test dataset.
 *
 * @param dataset Test dataset to evaluate model on.
 */
@Since("3.1.0")
def evaluate(dataset: Dataset[_]): RandomForestClassificationSummary = {
  val weightColName = if (!isDefined(weightCol)) "weightCol" else $(weightCol)
  // Handle possible missing or invalid prediction columns
  val (summaryModel, probabilityColName, predictionColName) = findSummaryModel()
  if (numClasses > 2) {
    new RandomForestClassificationSummaryImpl(summaryModel.transform(dataset),
      predictionColName, $(labelCol), weightColName)
  } else {
    new BinaryRandomForestClassificationSummaryImpl(summaryModel.transform(dataset),
      probabilityColName, predictionColName, $(labelCol), weightColName)
  }
}

@Since("1.4.0")
override def transformSchema(schema: StructType): StructType = {
  var outputSchema = super.transformSchema(schema)
  if ($(leafCol).nonEmpty) {
    outputSchema = SchemaUtils.updateField(outputSchema, getLeafField($(leafCol)))
  }
  outputSchema
}
override def transform(dataset: Dataset[_]): DataFrame = {
  val outputSchema = transformSchema(dataset.schema, logging = true)

  val outputData = super.transform(dataset)
```

```scala
  if ($(leafCol).nonEmpty) {
    val leafUDF = udf { features: Vector => predictLeaf(features) }
    outputData.withColumn($(leafCol), leafUDF(col($(featuresCol))),
      outputSchema($(leafCol)).metadata)
  } else {
    outputData
  }
}

@Since("3.0.0")
override def predictRaw(features: Vector): Vector = {
  // TODO: When we add a generic Bagging class, handle transform there: SPARK-7128
  // Classifies using majority votes.
  // Ignore the tree weights since all are 1.0 for now.
  val votes = Array.ofDim[Double](numClasses)
  _trees.foreach { tree =>
    val classCounts = tree.rootNode.predictImpl(features).impurityStats.stats
    val total = classCounts.sum
    if (total != 0) {
      var i = 0
      while (i < numClasses) {
        votes(i) += classCounts(i) / total
        i += 1
      }
    }
  }
  Vectors.dense(votes)
}

override protected def raw2probabilityInPlace(rawPrediction: Vector): Vector = {
  rawPrediction match {
    case dv: DenseVector =>
      ProbabilisticClassificationModel.normalizeToProbabilitiesInPlace(dv)
      dv
    case sv: SparseVector =>
      throw new RuntimeException("Unexpected error in RandomForestClassificationModel:" +
        " raw2probabilityInPlace encountered SparseVector")
  }
}

@Since("1.4.0")
override def copy(extra: ParamMap): RandomForestClassificationModel = {
  copyValues(new RandomForestClassificationModel(uid, _trees, numFeatures, numClasses), extra)
    .setParent(parent)
}

@Since("1.4.0")
override def toString: String = {
  s"RandomForestClassificationModel: uid=$uid, numTrees=$getNumTrees, numClasses=$numClasses, " +
    s"numFeatures=$numFeatures"
}

/**
 * Estimate of the importance of each feature.
 *
 * Each feature's importance is the average of its importance across all trees in the ensemble
 * The importance vector is normalized to sum to 1. This method is suggested by Hastie et al.
 * (Hastie, Tibshirani, Friedman. "The Elements of Statistical Learning, 2nd Edition." 2001.)
 * and follows the implementation from scikit-learn.
 *
```

```scala
 * @see `DecisionTreeClassificationModel.featureImportances`
 */
@Since("1.5.0")
lazy val featureImportances: Vector = TreeEnsembleModel.featureImportances(trees, numFeatures)

/** (private[ml]) Convert to a model in the old API */
private[ml] def toOld: OldRandomForestModel = {
  new OldRandomForestModel(OldAlgo.Classification, _trees.map(_.toOld))
}

@Since("2.0.0")
override def write: MLWriter =
  new RandomForestClassificationModel.RandomForestClassificationModelWriter(this)
}

@Since("2.0.0")
object RandomForestClassificationModel extends MLReadable[RandomForestClassificationModel] {

 @Since("2.0.0")
 override def read: MLReader[RandomForestClassificationModel] =
   new RandomForestClassificationModelReader

 @Since("2.0.0")
 override def load(path: String): RandomForestClassificationModel = super.load(path)

 private[RandomForestClassificationModel]
 class RandomForestClassificationModelWriter(instance: RandomForestClassificationModel)
   extends MLWriter {

   override protected def saveImpl(path: String): Unit = {
     // Note: numTrees is not currently used, but could be nice to store for fast querying.
     val extraMetadata: JObject = Map(
       "numFeatures" -> instance.numFeatures,
       "numClasses" -> instance.numClasses,
       "numTrees" -> instance.getNumTrees)
     EnsembleModelReadWrite.saveImpl(instance, path, sparkSession, extraMetadata)
   }
 }

 private class RandomForestClassificationModelReader
   extends MLReader[RandomForestClassificationModel] {

 /** Checked against metadata when loading model */
 private val className = classOf[RandomForestClassificationModel].getName
 private val treeClassName = classOf[DecisionTreeClassificationModel].getName

   override def load(path: String): RandomForestClassificationModel = {
     implicit val format = DefaultFormats
     val (metadata: Metadata, treesData: Array[(Metadata, Node)], _) =
       EnsembleModelReadWrite.loadImpl(path, sparkSession, className, treeClassName)
     val numFeatures = (metadata.metadata \ "numFeatures").extract[Int]
     val numClasses = (metadata.metadata \ "numClasses").extract[Int]
     val numTrees = (metadata.metadata \ "numTrees").extract[Int]

     val trees: Array[DecisionTreeClassificationModel] = treesData.map {
       case (treeMetadata, root) =>
       val tree =
         new DecisionTreeClassificationModel(treeMetadata.uid, root, numFeatures, numClasses)
       treeMetadata.getAndSetParams(tree)
       tree
```

```
  }
  require(numTrees == trees.length, s"RandomForestClassificationModel.load expected $numTrees" +
    s" trees based on metadata but found ${trees.length} trees.")
  val model = new RandomForestClassificationModel(metadata.uid, trees, numFeatures, numClasses)
  metadata.getAndSetParams(model)
  model
  }
 }
```

<div align="center">Code 5.5 Random Forest Classifier (open-source: under Apache Licence)</div>

By this we end the chapter and proceed with running the application.

# Chapter 6. Running the application on Apache Spark inside the EMR (Elastic MapReduce)

In the previous chapters we have done setup of necessary infrastructure on AWS cloud and together with it we made through with setting up the application implementation. Now that we are set with the implementation details, let execute the program to train the model and then test the model it on EMR using Apache Spark.

After setting up the Apache Maven and setting up the files, use the following command to build the .jar package.

```
$ # [After installing Apache Maven, execute below line in directory where pom.xml resides]
Important!!
$ mvn package
```



**Fig 6.1 Performing a Maven build**

## 6.1 Running application for Training the model

In order to run the program using .jar package, execute the following command.

```
$ # [Make sure to put in your S3 bucket URI's in place of mine.]

$ spark-submit --class com.njit.rp39.cs634.FinalTermProject.App target/cs634-ftp-opt5-1.0-
SNAPSHOT.jar train cluster s3n://cs634-option5-storage-bucket s3n://cs634-option5-storage-
bucket/data.csv diagnosis
```

This produces the following result shown in the screenshots below.

**Fig 6.2 Executing the Model Training process (1)**



**Fig 6.3 Executing the Model Training process (2)**



**Fig 6.4 Executing the Model Training process (3)**

**Fig 6.5 Complete look at the Training Process**

## 6.2 Running application for Testing the model

For running a test using testing dataset execute the following:

```
$ spark-submit --class com.njit.rp39.cs634.FinalTermProject.App target/cs634-ftp-opt5-1.0-
SNAPSHOT.jar test cluster s3n://cs634-option5-storage-bucket s3n://cs634-option5-storage-
bucket/test.csv diagnosis
```


**Fig 6.6 Executing the Model Testing process (1)**


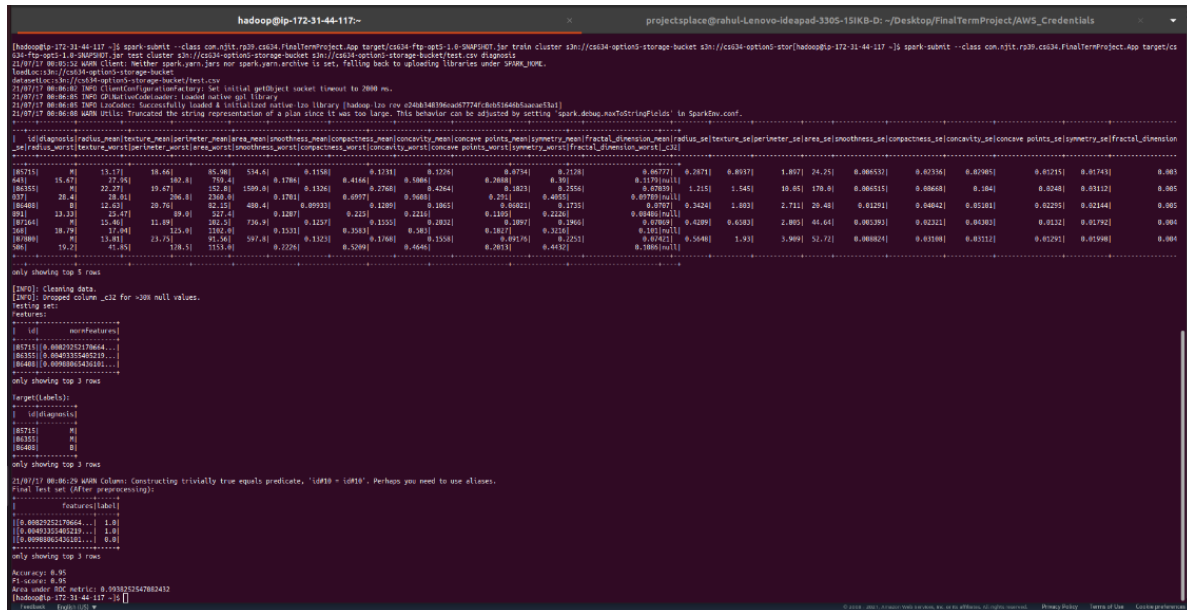**Fig 6.7 Executing the Model Testing process (2)**

**Fig 6.9 Complete look at the Testing Process**

By this we complete the project execution and our model is trained and tested.

*Warning: Clean-up!! Make sure to <u>delete</u> the resources and turn them in to the AWS. This may save us the trouble if getting any unexpected cost incurred by running the unused resources. Keep whatever is needed. Thank you.*

# Chapter 7. Conclusion

We have finally reached the conclusion for this project document. We are now familiar with using Apache Spark and working with it. Also, we have leveraged the power of Cloud resources, specifically on AWS. We now know how to work on AWS Cloud environment and have worked on one of its most widely used infrastructure and service for working with Data mining and Big-data analytics problems, namely the Elastic MapReduce (EMR).

We also worked on another AWS service named the S3, or the Simple Storage Service, and used it as a base for storing our big-data files in the backend, that would later work with our application program for solving a general classification problem, such as predicting Breast-Cancer within patient records.

We have also seen the advantages of using this approach for solving big-data problems, namely working on large-dataset in an efficient manner, performance of training model on large application, using flexibility of the cloud resources and ease in deploying model for testing.

Apache Spark is an open-source unified analytics engine for large-scale data processing [9]. Spark is the outcome of years of endeavored researches carried on large-scale data processing with the aim of making the data mining process and algorithms efficient with the evolving computing infrastructures.

Along with it came the rise of Cloud Computing field, AWS being the pioneer at cloud computing [10], introducing a wide variety of resources for developers and data scientist to leverage the full potential and creativity of designing a needful application. It is now ideal to say that the field of data mining is just at its tip of the ice and will be blooming to become a large one in the near future.

By this we may now conclude the Final Term Project.

# References

[1] Original Dataset: Raw form
        https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29

[2] Dataset source: https://www.kaggle.com/uciml/breast-cancer-wisconsin-data

[3] https://www.tdktech.com/tech-talks/data-mining-and-machine-learning/

[4] https://en.wikipedia.org/wiki/Statistical_classification

[5] Details on Breast-Cancer-Wisconsin: https://www.dhs.wisconsin.gov/publications/p01573a.pdf

[6] https://en.wikipedia.org/wiki/Random_forest

[7] https://en.wikipedia.org/wiki/Computer_cluster

[8] https://aws.amazon.com/emr/

[9] https://en.wikipedia.org/wiki/Apache_Spark

[10] Amazon Web Services: The Pioneer in Cloud Computing, casecent.re/p/133541

[11] Spark MLlib: https://spark.apache.org/docs/latest/ml-guide.html

[12] Spark SQL: https://spark.apache.org/docs/latest/sql-programming-guide.html

[11] Data Mining: Concepts and Techniques, Han et al., Elsevier, 2011, ISBN 978-0-12-381479-1

[12] Introduction to Data Mining, Tan et al., Pearson, 2019, ISBN-13: 978-0-13-312890-1