

PROJECT REPORT

TITLE : CREDIT CARD FRAUD PREDICTION USING JAVA AND MACHINE LEARNING

INTRODUCTION:

Main challenges involved in credit card fraud detection are:

1. Enormous Data is processed every day and the model build must be fast enough to respond to the scam in time.
2. Imbalanced Data i.e most of the transactions (99.8%) are not fraudulent which makes it really hard for detecting the fraudulent ones
3. Data availability as the data is mostly private.
4. Misclassified Data can be another major issue, as not every fraudulent transaction is caught and reported.
5. Adaptive techniques used against the model by the scammers.

How to tackle these challenges?

1. The model used must be simple and fast enough to detect the anomaly and classify it as a fraudulent transaction as quickly as possible.
2. Imbalance can be dealt with by properly using some methods which we will talk about in the next paragraph
3. For protecting the privacy of the user the dimensionality of the data can be reduced.
4. A more trustworthy source must be taken which double-check the data, at least for training the model.
5. We can make the model simple and interpretable so that when the scammer adapts to it with just some tweaks we can have a new model up and running to deploy.

Hardware/Software REQUIREMENTS:

HARDWARE:

Processor : Intel i5 CPU at 1.60 Ghz

RAM : 4GB

System type : 64 bit OS

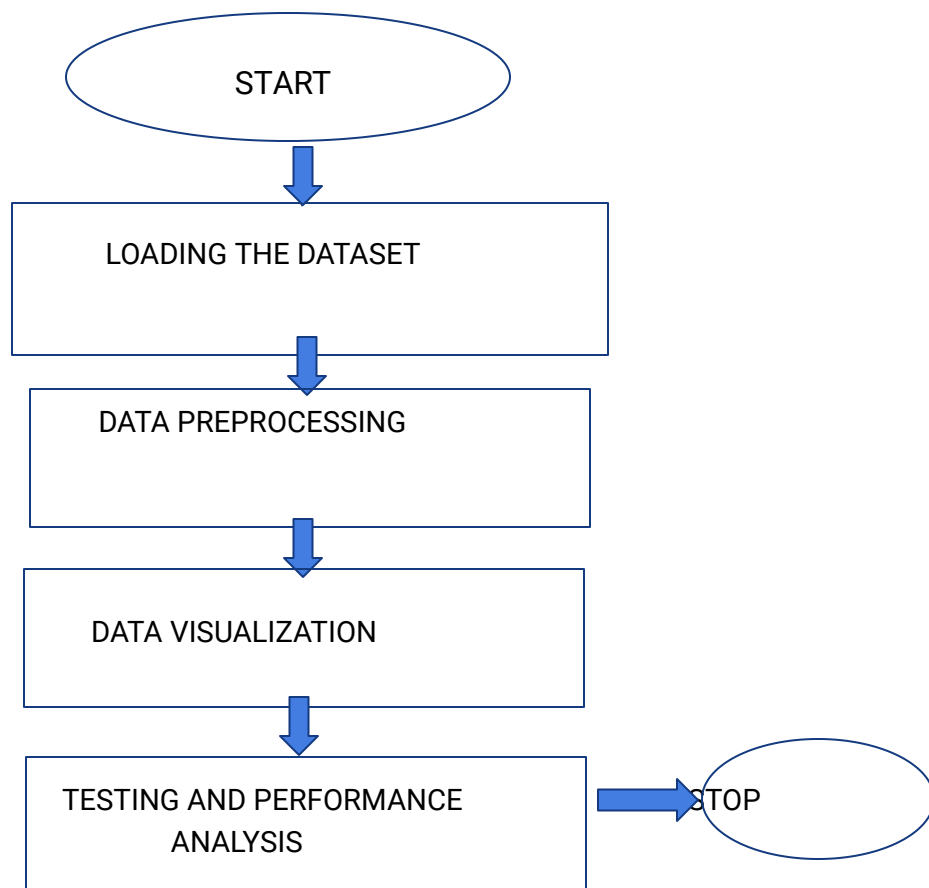
SOFTWARE:

OS : WINDOWS 7 OR ABOVE /ARCH LINUX

Programming language : Java

IDE : Eclipse

FLOW DIAGRAM:



SOURCE CODE:

```

1 package org.ml;
2
3 import tech.tablesaw.api.Table;
4 import weka.classifiers.Evaluation;
5 import weka.core.Instances;
6 import weka.core.converters.CSVLoader;
7 import java.io.File;
8 import java.io.IOException;
9 import weka.classifiers.trees.RandomForest;
10
11 public class RandomForestDemo
12 {
13
14     public static Instances getDataSet(String fileName) throws IOException
15     {
16
17         /**
18          * we can set the file i.e., loader.setFile("filename") to load the data
19          */
20         int classIdx = 1;
21         /** the CSVLoader to load the CSV file */
22         CSVLoader loader = new CSVLoader();
23         /** load the training data */
24         //loader.setSource(RandomForestDemo.class.getResourceAsStream("/") + fileName));
25         /**
26          * we can also set the file like loader3.setFile(new
27          * File("test-confused.arff"));
28          */
29         loader.setFile(new File(fileName));
30         Instances dataSet = loader.getDataSet();
31         /** set the index based on the data given in the CSV files */
32         dataSet.setClassIndex(classIdx);
33         return dataSet;
34     }
35
36     /**
37      * This method is used to process the input and return the statistics.
38      *
39      * @throws Exception
40      */
41     public static void main(String args[]) throws Exception
42     {
43
44         try
45         {
46             Table fraud_data=Table.read().csv("C:\\Users\\A. SRINIDHI\\eclipse-workspace\\org.ml\\src\\main\\java\\org\\ml\\fraud_data
47             System.out.println(fraud_data.shape());
48         }
49         catch(IOException e)
50         {
51             e.printStackTrace();
52         }
53
54         Instances trainingDataSet = getDataSet("C:\\Users\\A. SRINIDHI\\eclipse-workspace\\org.ml\\src\\main\\java\\org\\ml\\credit_fra
55         Instances testingDataSet = getDataSet("C:\\Users\\A. SRINIDHI\\eclipse-workspace\\org.ml\\src\\main\\java\\org\\ml\\credit_fra
56
57         RandomForest forest=new RandomForest();
58         forest.setNumFeatures(10);
59         forest.buildClassifier(trainingDataSet);
60         /**
61          * train the algorithm with the training data and evaluate the
62          * algorithm with testing data
63          */
64         Evaluation eval = new Evaluation(trainingDataSet);
65         eval.evaluateModel(forest, testingDataSet);
66
67         /** Print the algorithm summary */
68         System.out.println("*** Random Forest Evaluation with Datasets ***");
69         System.out.println(eval.toSummaryString());
70         System.out.print("The expression for the input data as per alogorithm is : ");
71         System.out.println(forest);
72     }
73 }

```

OUTPUT OF THE CODE:

```

827 rows X 13 cols
** Random Forest Evaluation with Datasets **

Correlation coefficient      0.9916
Mean absolute error         0.0056
Root mean squared error     0.0386
Relative absolute error     1.334 %
Root relative squared error  9.1558 %
Total Number of Instances   248

The expression for the input data as per algorithm is : RandomForest

Bagging with 100 iterations and base learner

```

REASONS FOR USING RANDOM FOREST CLASSIFIER:

1) INCREASE THE PREDICTIVE POWER:

Firstly, there is the `n_estimators` hyperparameter, which is just the number of trees the algorithm builds before taking the maximum voting or taking the averages of predictions. In general, a higher number of trees increases the performance and makes the predictions more stable, but it also slows down the computation.

Another important hyperparameter is `max_features`, which is the maximum number of features random forest considers to split a node. Sklearn provides several options, all described in the documentation.

The last important hyperparameter is `min_sample_leaf`. This determines the minimum number of leafs required to split an internal node.

2) INCREASE IN MODEL SPEED

The **`n_jobs`** hyperparameter tells the engine how many processors it is allowed to use. If it has a value of one, it can only use one processor. A value of “-1” means that there is no limit. The **`random_state`** hyperparameter makes the model’s output replicable. The model will always produce the same results when it has a definite value of `random_state` and if it has been given the same hyperparameters and the same training data. Lastly, there is the **`oob_score`** (also called oob sampling), which is a random forest cross-validation method. In this

sampling, about one-third of the data is not used to train the model and can be used to evaluate its performance. These samples are called the out-of-bag samples. It's very similar to the leave-one-out-cross-validation method, but almost no additional computational burden goes along with it.

CONCLUSION:

The result obtained by using random forest classifier is around 0.991 i.e 99.1%. Hence the proposed method works well when compared to existing models based on the decision tree algorithms.

SCOPE:

A practical mobile application based on fraud alert can be a suitable prototype as for the future scope.

References:

<https://www.romexsoft.com/blog/implement-credit-card-fraud-detection/>

Github link:

<https://github.com/smartinternz02/SPS-10809-Creditcard-Fraud-Prediction->