# Detecting Financial Fraud Using Machine Learning: Winning the War Against Imbalanced Data

N.Ashok[*1], D.Rajanipriya[2], M.Akhila[3], P.Yamini[4]

[*1]Assistant Professor(phd), Department of Information Technology, Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India

[2, 3, 4] Department of Information Technology, Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India

\------------------------------------------------------------------------**\*\*\*\***\------------------------------------------------------------------------

**ABSTRACT-Fraud detection is a process of monitoring the transaction behavior of a cardholder in order to detect whether an incoming transaction is done by the cardholder or others.To** obtain normal/fraud behavior features based on machine learning techniques, and then utilize these features to check if a transaction is fraud or not. The data used in our experiments come from an e-commerce company in Europe. In this project, we divide the data-set into train and test set by making use of the different machine learning algorithms like smote, under sampling, over sampling, both under and over sampling, rose. We have compared the accuracy and performance of these models using confusion matrix and roc curve

*Keywords: DFF(Detecting Financial Fraud),SMOTE(Synthetic Minority Over Sampling Technique),AUC(Area under the curve),Roc(Receiver Operating Curve)*

## I. INTRODUCTION

Credit cards are most familiar to every one.Since the physical card is not needed in the online transaction environment and the card's information is enough to complete a payment, it is easier to conduct a fraud than before. Transaction fraud has become a top barrier to the development of e-commerce and has a dramatic influence on the economy. Hence, fraud detection is essential and necessary. Due to these frauds many of the common people and sometimes countries revenue may lose. Fraud detection is a process of monitoring the transaction behavior of a cardholder in order to detect whether an in-coming transaction is done by the cardholder or others. Generally, there are two kinds of methods for fraud detection. Misuse detection and anomaly detection. Misuse detection uses classification methods to determine whether an incoming transaction is fraud or not. Usually, such an approach has to know about the existing types of fraud to make models by learning the various fraud patterns. Anomaly detection

is to build the profile of normal transaction behavior of a cardholder based on his/her historical transaction data, and decide a newly transaction as a potential fraud if it deviates from the normal transaction behavior. However, an anomaly detection method needs enough successive sample data to characterize thenormal transaction behavior of a cardholder.This paper is about Anomaly method.
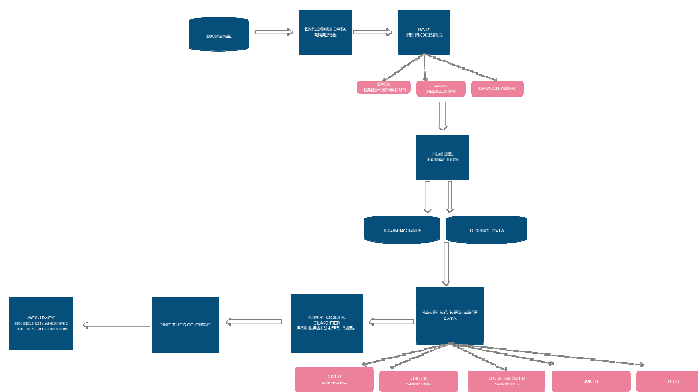
## II.SYSTEM ARCHITECTURE



**Fig 1 Block Diagram**

Framework design is the decided model that depicts the structure, direct, and more perspectives on a framework. A planning portrayal is a standard outline and delineation of a framework, filtered through with the ultimate objective that supports thinking about the structures and practices of the framework. A structure arrangement can contain framework parts and the sub-frameworks created that will team up to finish the general framework. There have been attempts to formalize vernaculars to portray structure working; as a rule these are called plan depiction tongues.

## III. DATASET

The datasets contains transactions made by credit cardsin September 2013 by european cardholders.This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions.It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, we cannot provide the original features and more background information about the data. Features V1, V2, … V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependant cost-senstive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.



**Fig 2 Screenshot of Dataset for Credit Card**

## IV. IMPLEMENTATION

We are using Rstudio to implement algorithm and with R language we are writing the code.**RStudio** is an integrated development environment (IDE) for R,programming language for statistical computing and graphics. It is available in two formats: R Studio.Desktop is a regular desktop application while R Studio Server runs on a remote server and allows accessing R Studio using a web browser.An IDE that was built just for R. Syntax highlighting, code completion, and smart indentation. Execute R code directly from the sourceeditor. R made its first appearance in 1993.A large group of individuals has contributed to R by sending code and bug reports. Since mid-1997 there has been a core group (the "R Core Team") who can modify the R source code archive. As stated earlier,R is a programming language and software environment for statistical analysis, graphics representation and reporting. The following are the important features of R is a well-developed,simple and effective programming language which includes conditionals, loops, user defined recursivefunctions and input and output facilities. R has an effective data handling and storage facility, R provides a suite of operators for calculations on arrays, lists, vectors and matrices. R provides a large, coherent and integrated collection of tools for data analysis. R provides graphical facilities for data analysis and display either directly at the computer or printing at the papers. As a conclusion, R is world's most widely used statistics programming language. It's the1st choice of data scientists andsupported by a vibrant and talented community of contributors.You can download the Windows installer version of R from R-3.2.x for Windows (32/64 bit) and save it in a local directory. As it is a Windows installer (.exe) with a name "R-version-win.exe".You can just double click and run the installer accepting the default settings. If your Windows is32-bit version, it installs the 32-bit version. But if your windows is 64-bit, then it installs both the 32-bit and 64-bit versions. After installation you can locate the icon to run the Program in a directory structure "R\R3.2.x\bin\i386\Rgui.exe" under the Windows Program Files.Clicking this icon brings up the R-GUI which is the R console to do R Programming.There are five algorithms which we are implementing undersampling,over sampling,both under and over sampling,rose,smote for sampling the imbalance data.

## V. RESULTS AND DISCUSSION

**Data Extraction:**The data sets contains transactions made by credit cards in September 2013 by European cardholders. This data-set presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The data-set is highly unbalanced, the positive class(frauds) account for 0.172% of all transactions. Dueto confidentiality reasons, the data was anonymized variable names wererenamed to V1, V2, V3 until V28. Moreover, most of it was scaled, except for the Amount and Class variables.

**Explotary Data Analysis:**In the exploratory data analysis we will do five tasks

1)Calculate the Class Count and Class Percentage

2) Visualize theData

3)Split the transactions into days

4) Find the time of transaction of a particular day

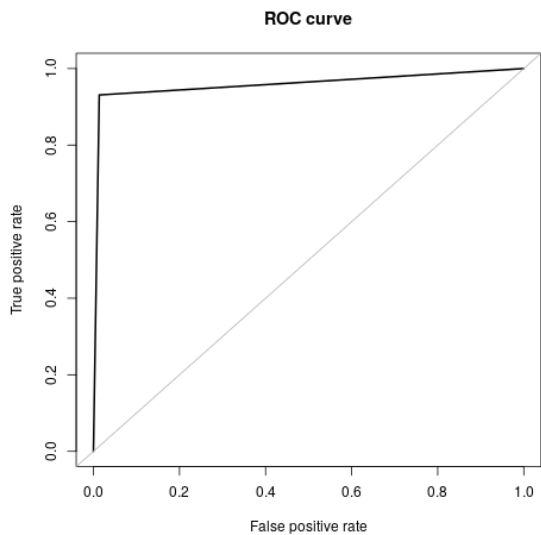5)Find the time quarters for transaction

We will create a data frame which consists of a class group, count of records of each class, and calculates the percentage of each class. The input data is CSV format and the input filename "input.csv" and the output file is "output.csv". After this we will visualize the data-set to check how many transactions are fraudulent and non-fraudulent plotting the graph(Using ggplot2). The input data is CSV format and the input filename "input.csv", the output is an image file. split the transactions into day 1 and day 2 based on recorded the time. Add the column named day into the data-set and then find the time of transaction made for that particular day, based on the time and the day details the time quarter of each transaction which falls on 4 different time quarters like0 - 6 hr - 1st quarter ,7 - 12 hr - 2nd quarter ,13 - 18 hr - 3rd quarter ,19 - 24 hr - 4th quarter based on the time_day details and add the column named "Time_qtr".Plot the graph of time quarters for the fraudulent and non-fraudulent transactions.

**Data Preprocessing**:Split the dataset into train and test datasets, and scale all the fields except the Class field.The dataset should divide 70% of records as training data_set and rest records as test dataset(using catools library to split the data).Then find the logistic classifier on the train set andpredict the response on the test set(exclude the class field), and output the model's accuracy using the ROC curve, Use Generalized Linear Models(glm) function to predict with binomial. If the predicted value is less than 0.5 make it 0 else make it 1 with variable name "y_pred_normal"..Omit the na values.

**Sampling The Imbalance Data**:Sample the imbalanced training data-set, using over method ."Oversample" means to artificially create observations in our data set belonging to the class that is under represented in our data.use the ovun.sample function for sampling on training dataset ,N = Double of the sample data,Method = over. Sample the imbalanced training data-set, using both method. use the ovun.sample function for sampling on the training dataset. seed = 222, Method = both ,N = Number of the sample data, Probability = 0.5. Sample the imbalanced training data-set, using the rose method. use the ovun.sample function for sampling on the training dataset, seed = 111.Sample the imbalanced training data-set, using smote method. se the smote function , and the size of Minority class – 100, the size of the Majority class – 200.
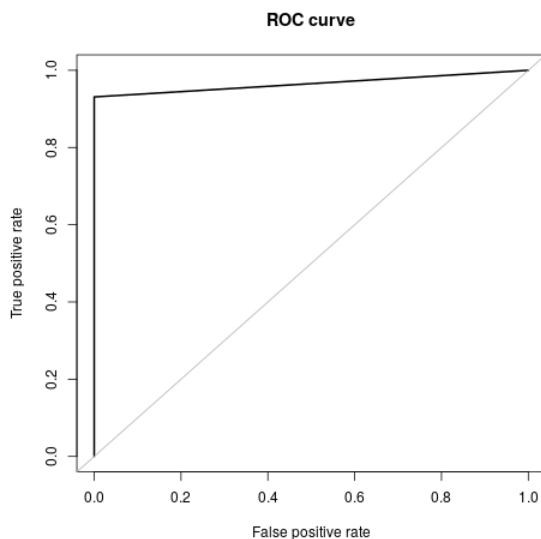
**Implementing Algorithm :**Find the auc value for oversampled data-set using logistic regression for imbalanced data.Re-Sample the data set with over method ,Set seed as 111 ,Omit the na's.If the predicted value is less than 0.5 make it 0 else make it 1 with variable name "y_pred_over". Find the auc value for under sampled data-set using logistic regression for imbalanced data. Re-Sample the data set with under method. Find the auc value for both sampled data-set

using logistic regression for imbalanced data , re-Sample the data set with both method and find the auc value for rose ,smote and sampled data-set using logistic regression for imbalanced data.Analyzing the best algorithm by comparing all the auc curve..
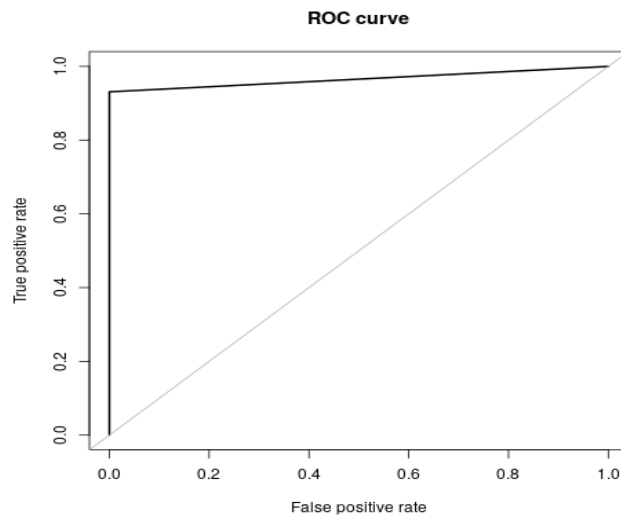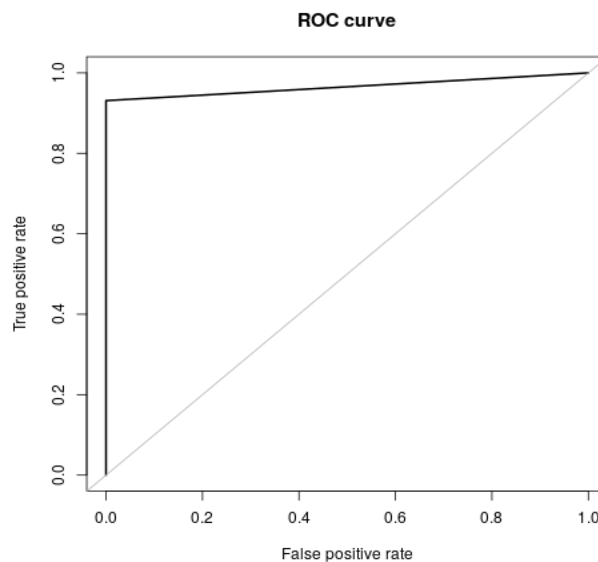
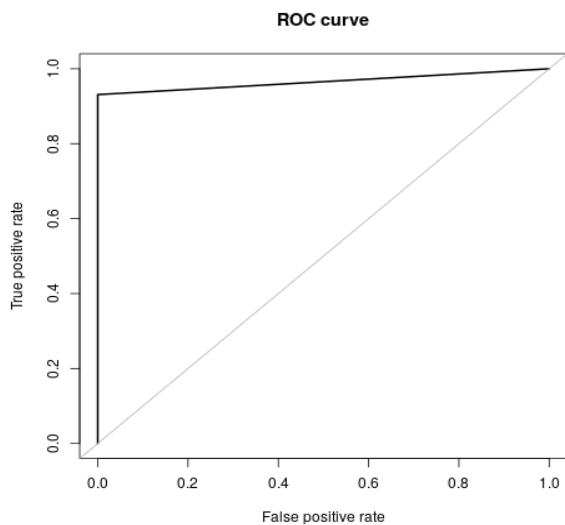**Fig 4 Under Sampling**



Area under the curve (AUC): 0.960

**Fig 5 Both Under and Over Sampling**



Area under the curve (AUC): 0.960

**Fig 6 Rose**



Area under the curve (AUC): 0.960

**Fig 3 Over Sampling**



Area under the curve (AUC): 0.960

Area under the curve (AUC): 0.972

**Fig 7 Smote**

## VI .CONCLUSION

Thus the Smote algorithm works far better than the other algorithms in predicting theCredit Card Financial Fraud.

## VII. FUTURE ENHANCEMENT

The future upgrade of this work is tobe based on new multiple models with varying access pattern needs attention to improve the effectiveness. Privacy preserving techniques applied in distributed environment resolves the security related issues preventing private data access.

**REFERENCES**

[1]"BLAST-SSAHA Hybridization for Credit Card Fraud Detection" (AmlanKundu, SuvasiniPanigrahi, ShamikSural,) IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, VOL. 6, NO. 4, OCTOBER-DECEMBER 2009.

[2] " Research on Credit Card Fraud Detection Model Based on Distance Sum"( Wen-Fang-YU, Na Wang)

[3]"Detecting Credit Card Fraud by Decision Trees and Support Vector Machines" (Y. Sahinand E. Duman)

[4]" Fraudulent Detection in Credit Card System Using SVM &amp; Decision Tree" (VijayshreeB. Nipane, Poonam S. Kalinge, DipaliVidhate, Kunal War, Bhagyashree P. Deshpande)

[5] "Credit Card Fraud Detection Using Decision Tree Induction Algorithm" SnehalPatiletal, International Journal of Computer Science and Mobile Computing, Vol.4 Issue.4, April-2015

[6]"Data Mining Techniques for Credit Card Fraud Detection: Empirical Study "(MarwanFahmi, AbeerHamdy, Khaled Nagati).

[7] Gupta, Shalini, and R. Johari. "A New Framework for Credit Card Transactions Involving Mutual Authentication between Cardholder and Merchant." International Conference on Communication Systems and Network Technologies IEEE, 2011

[8] Y. Gmbh and K. G. Co, "Global online payment methods: Full year 2016,"