

# CREDIT CARD FRAUD DETECTION USING MACHINE LEARNING

S.Madhuri<sup>1</sup>, N.Hiranmayee<sup>2</sup>, Y.Revathi<sup>3</sup>, S.S.D.Lavanya<sup>4</sup>, M.Kartheek<sup>5</sup>, A.ChandraNagaSai<sup>6</sup>

<sup>1</sup>Student, Department of Computer Science and Engineering, Sri Vasavi Engineering College, Pedatadepalli, Tadepalligudem, AndhraPradesh, India

<sup>2-6</sup> Department of Computer Science and Engineering, Sri Vasavi Engineering College, Pedatadepalli, Tadepalligudem, AndhraPradesh, India

\*\*\*

**Abstract** - The Project Credit Card Fraud Detection provides fundamental concepts for identifying financial fraud. Financial fraud is still a common and expensive tactic in the connected and digital world of today. Since majority of transactions are done using credit cards these days, the risks have grown throughout the current time. The size of data created by digital transactions is constantly growing, making manual fraud detection methods are laborious and inefficient. Credit Card fraud can be identified by Machine Learning (ML) algorithms, which analyze data with a diversity of methods to produce the most accurate results when identifying fraudulent transactions. This Research focuses on application of Random Forest as an effective tool for credit card fraud detection. Random Forest, an ensemble learning technique, belongs to family of decision tree-based methods is able to handle large datasets, non-linear relationships and feature analysis makes it well suited for this task.

**Key Words:** Credit Card, Fraud Detection, Machine Learning, Random Forest, Ensemble, Decision Tree.

## 1.INTRODUCTION

Maintaining the security and stability of the financial industry depends on identifying and stopping fraudulent transactions. ML based credit card fraud detection is a vital use of technology in the banking sector. In recent years, the rapid advancement of ML techniques has empowered organizations to enhance their ability to identify and prevent fraudulent activities. This approach leverages data-driven methods to analyze vast amounts of data and detect suspicious outlines that may indicate fake transactions. This is where the power of ML steps came into action as a vigilant guardian against the frauds. We are living in a world which is rapidly adapting digital payment systems. Illegal transaction is a most common and costly problem that affects individuals, businesses, and financial institutions worldwide. The impact of fraud extends beyond monetary losses, it erodes trust in financial systems, damages reputations, and can have far-reaching consequences for both individuals and organizations. ML algorithms, such as unsupervised as well as supervised learning, are employed to build models which can automatically classify transactions as legitimate or potentially fraudulent. These models rely on historical data and continuously adapt to evolving fraud schemes. Key components of this process include feature engineering, data

preprocessing, and model training, which enable the system to make accurate predictions. Algorithms of ML can analyze the vast amount of data and identify the patterns that may indicate the fraudulent activities. In this report, we will discover how ML is used for detecting frauds. ML in actuality, is a bunch of clever algorithms which are mostly used to find patterns in a data stream of any kind, to provide helpful information, includes the type of fraud committed, future patterns. To handle this problem, advanced data-driven techniques are being employed to detect and prevent fraudulent transactions. The Random Forest algorithm is one such effective ML method. For a diversity of ML applications, including regression, anomaly detection, and classification, Random Forest is a potent ensemble learning method. When it is implemented in Python using the Scikit-Learn package, it becomes a powerful and popular tool for creating predictive models. To increase the system's overall accuracy and resilience, ensemble learning combines the predictions of several ML models. The greatest applications of Random Forest are in classification tasks. Scikit-Learn's 'RandomForestClassifier' class is castoff for this purpose. It builds on collection of decision trees, where each tree votes on the class label, and the final prediction is determined by majority voting. Random Forest, when used with Scikit-Learn, provides an easy-to-implement, highly accurate, and robust result for several ML tasks. Its ability to handle a varied series of data types, its resistance to overfitting, and its capacity to handle high-dimensional data make it a widespread choice for many real-world applications in both academia and industry.

## 2. LITERATURE SURVEY

This Literature survey will summarize some preliminary research that was conducted by numerous writers on this relevant work, and we'll take some significant papers into account and continue to develop our work.

[1]. T.Singh, F.DiTroia, C.Vissagio (2015), proposed a research paper "SVM and malware detection". In this research, we test three advanced malware scoring techniques that have exposed potential in prior research, namely, Hidden Markov Models, Simple Substitution Distance, and Opcode Graph based detection. We then perform a careful robustness analysis by employing morphing strategies that cause each score to fail. We show

that combining scores using a SVM profits results that are suggestively additional robust than those attained utilizing somewhat of the individual scores.

[2]. **Samaneh Sorournejad, Zojah and Atani (2016)**, proposed a research article "A survey of Credit Card Fraud Detection Techniques". In this paper, after inspecting difficulties of credit card fraud detection, we seek to review the states in credit card fraud detection techniques, datasets and evaluation criteria. The pros and cons of fraud detection techniques are computed and compared. Furtherly, techniques used for fraud detection can be classified based on the type of data they handle. Some methods are tailored for numerical data, which includes quantitative values, while others are specialized in handling categorical data, which comprises non-numeric attributes like names or categories. This categorization helps in choosing the right method based on the data characteristics for effective fraud detection.

[3]. **Wedge, Canter and Rubio (2017)**, proposed a research "Solving the False Positives problem in fraud prediction" paper published. Our automated feature engineering-based method to lower false positives in fraud prediction is presented in this paper. The fraud estimate industry is plagued with false positives. Only one in five cases that are reported as fraud are thought to be fraudulent, and one in six consumers report having had a legitimate transaction declined in the previous year. We employ the Deep Feature Synthesis Algorithm to repeatedly derive behavioral features based on the card's historical data linked to a transaction in order to solve this issue. To build a classifier, we employ a random forest. Using data from a big international bank, we tested our ML model and contrasted it with their current system.

[4]. **Aditya Oza(2019)**, proposed a research paper " A Survey on Fraud Detection". This paper applies various ML techniques on Logistic regression and SVM to the problematic payments fraud identification using a labeled dataset containing payment transactions and demonstrates that these methods have a high accuracy rate and a low number of false positives after it arises to detecting fraud transactions.

[5]. **Gupta A., Lohani, M. C., & Manchanda, M. (2021)**, This is the era, where the plastic money concept is widely adapted all over the world, but each novel expertise has their own loopholes also. In this script numerous types of anomalies can be which can harm the stoner economically. These anomalies can be defined as frauds in financial sector. To descry these types of frauds, numerous ways and models are proposed by the experimenters. In this study the proposed work tries on implementing an automated prototype using different ML algorithms to identify these kinds of frauds, especially related to credit cards transactions. The proposed model applied four techniques used in ML, specifically, Random Forest, Logistic Regression, Naive Bayes and SVM on a very large dataset to recognize the fraud.

### 3. EXISTING SYSTEM

Financial fraud detection classification that use ML have made significant advancements but they even aspect numerous limitations and challenges. These limitations include:

**Overfitting:** Overfitting to training data is a problem for some complex machine learning models, which makes them less generic to new data. Overfitting may arise instance of feature spaces with high dimensions, which are typical in contemporary datasets.

**Imbalanced Data:** When dealing with imbalanced datasets, where one state has much less samples than the other, some models may struggle to correctly classify the minority class. Additional techniques may often needed to address this issue.

**Scalability:** A robust system is required to manage a high volume of transactions. Scalable remedies that can manage big datasets and analyze transactions in real-time are necessary due to trading that occurs frequently and the enormous number of financial transactions.

**Concept Drift:** Fraudsters constantly adapt to new strategies, which can lead to concept drift. Models will become less active as time passes they struggle to adapt to evolving fraudulent behavior.

**False Positives and Negatives:** Finding a middle ground between reducing the number of false positives (regular transactions reported as fraudulent) and false negatives (fraudulent transactions missed) is difficult. Overly cautious systems may inconvenience genuine users, while overly lenient ones may miss fraud.

**Unlabeled Data:** Anomaly detection-based systems often require labeled data for training. Creating such labeled datasets can be labor-intensive, and obtaining sufficient labeled samples of fraud can be problematic.

**Model Drift:** Model drift can affect ML replicas that are applied in fraud identification systems. Over time, the models may become less effective as fraud techniques change, necessitating ongoing retraining and adaptation.

### 4. PROPOSED SYSTEM

The Goal of suggesting Random Forest algorithm is to overcome the drawbacks and limitations of the existing system by utilizing the Random Forest's advantages. An overview of main components and methodology of the suggested system is given below:

**Data Collection:** Collecting the vast number of historical transactions are gathered. The data include amount of transactions, and the card holders information and the target class feature.

**Data Preprocessing:** The data that is collected from various sources, is cleaned and preprocessed to remove errors, inconsistencies, unwanted attributes, and duplicate records. It may also involve data transformation and normalization to make it suitable for data analysis.

**Feature Engineering:** Feature engineering is a critical and creative process in the era of ML. It involves selecting, transforming, and creating relevant input variables (features) for a predictive model. Effective feature engineering can significantly impact the model's performance, making it an essential step in the preprocessing pipeline.

**Handling unbalanced data:** Transactional data are highly unbalanced where legitimate transactions may dominate over fraudulent transactions. To build the effective model handling unbalanced data is necessary. It can be achieved by applying oversampling techniques(SMOTE).

**Splitting the dataset:** Dividing the dataset into sets for testing and training. The system's effectiveness is increased when a larger pool of transactions is included in the training set by means of opposed to the testing set.

**Model Selection:** Select the suggested algorithm Random Forest for its high predicted accuracy, capacity to handle unbalanced datasets and efficiency. As it is the ensemble method, it combines the advantages of many models. These will be drastic decrease in the false positives and negatives.

**Performance Evaluation:** The metric features like accuracy, precision, recall and f1score are taken into consideration. Additionally, the threshold notch of the aforementioned metrics is used to estimate the model's performance.

**Result Analysis:** The target variable depicts the output of our prediction. Value 1 of the target feature indicates the fraudulent transaction and legitimate transaction otherwise.

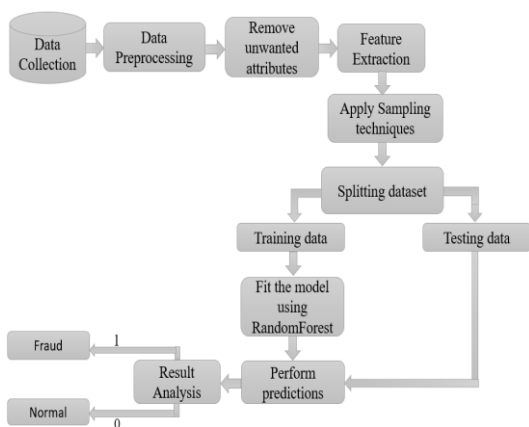


Fig -1: System Architecture

#### 4.1 PROPOSED DATASET

We have chosen a dataset from Kaggle website to into this module.

**Dataset URL:**

<https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>

Only the numerical input variables that have undergone Principal Component Analysis (PCA) transformation are included in the dataset. Using PCA, the principal components are functions V1, V2, V3..., V28. The transaction amount is represented by the feature "Amount," which can be utilized for cost-sensitive learning based on examples. The response variable, 'Class', takes value 0 in the absence of fraud and 1 in the case of fraud.

#### 4.2 PROPOSED RANDOM FOREST MODEL

Preparing dataset, training the model, performance evaluation and result analysis are the few activities in our proposed model. Let's go deeper into the procedure.

**Step1: Obtaining and Preparing datasets:**

Formulate the dataset in the appropriate directories with separate subdirectories for each class. The dataset should be divide into training and testing sets and hence it is possible to analyze the model's performance effectively.

**Step2: Preprocessing:**

Preprocessing is the critical step in the ML in identifying instances of financial fraud to get the cleaned data for effective training of the model and to enhance the predictive capacity.

**Step3: Feature Engineering:**

Feature engineering is a creative process in the era of ML and data science. It involves selecting, transforming, and creating relevant input variables (features) for a predictive model. Effective feature engineering can significantly influence a ML model's performance.

**Step4: The model Architecture's Personalization:**

Transactional datasets are highly unbalanced where legitimate transactions may dominate over fraudulent transactions. To build the effective model handling unbalanced data is necessary. It can be achieved by applying oversampling techniques (SMOTE).

**Step5: Training the model:**

A supervised ML algorithm Random Forest is applied to effectively and efficiently identify fraudulent transactions. Training a model on identify illegal transactions requires preprocessed data without any errors and fine tuning to build the model user friendly.

**Step6: Evaluation of Performance and results:**

The model performance is estimated grounded on metric scores like accuracy, precision, recall and f1-score. One can fine tune the model based on confusion matrix which depicts the false negatives and false positives if any.

**5. EXPERIMENTAL RESULTS**

The outcomes of experiment specify that Random Forest is a useful model for identifying credit card fraud. It possesses precision, accuracy, recall, and f1. To make the required adjustments for deployment in practice, how the model performs in the real world must be assessed.

Accuracy in the situation of financial fraud detection, is the proportion of correctly classified transactions to all transactions, divided into true positives and true negatives, representing legitimate and fraudulent transactions, respectively.

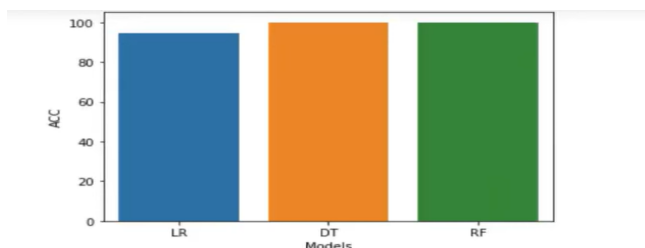
Precision, it is ratio of true positives to all positive predictions (true positives plus false positives). High precision is crucial for fraud detection because it denotes a low false alarm rate that is, transactions that are valid but are mistakenly reported as fraudulent.

Recall, also known as sensitivity or true positive rate, it is the proportion of actual positives (true positives plus false negatives) to true positives. It assesses how well the model can detect all fraudulent transactions. High recall is essential to preventing fraud from going unnoticed.

The harmonic average of recall and precision is the F1 score. It offers a harmony between these two measurements. When recall and precision are equally significant, the F1 score can be a helpful indicator.

**5.1 COMPARISON OF MODELS**

A bar plot is used to compare the accuracy of different models. Here we compared the accuracy of Random Forest with Logistic Regression and Decision Tree models of Machine Learning which are familiar for classification tasks. Random Forest gives the highest accuracy of 99.9% than other models.



**Fig -2: Comparison of models**

**5.2 CONFUSION MATRIX**

A confusion matrix is a fundamental visualization to measure the capability of a classification model. It shows the true positive, true negative, false positive, and false negative predictions, allowing you to see the model's accuracy, precision, recall, and F1-score. This matrix represents the actual positives, actual negatives, predicted positives and predicted negatives. True Positive, False Positive, True Negative and False Negative values are 55064, 0, 9 and 55003 respectively.

$$\begin{bmatrix} 55064 & 0 \\ 9 & 55003 \end{bmatrix}$$

**5.3 CLASSIFICATION REPORT**

This Classification report provides precision, recall, f1-score and support for a binary classification. It also provides overall accuracy, micro average and weighted average metrics.

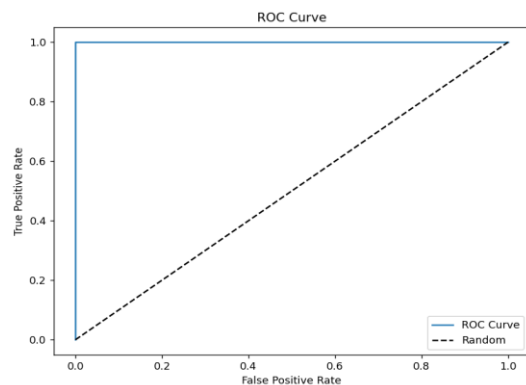
	precision	recall	f1-score	support
0	1.00	1.00	1.00	55073
1	1.00	1.00	1.00	55003
accuracy			1.00	110076
macro avg	1.00	1.00	1.00	110076
weighted avg	1.00	1.00	1.00	110076

**Fig -3: Classification Report**

**5.4 ROC (Receiver Operating Characteristic)**

The area under the ROC curve (AUC) is a familiar metric for summarizing the capability of a binary classification model.

particularly when assessing the trade-off between true positive and false positive rates at different threshold values.

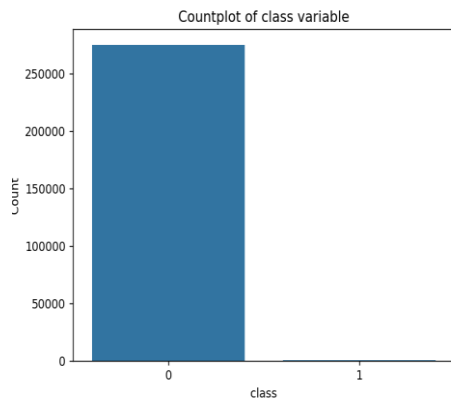


**Fig -4: ROC Curve**



## 5.5 Count plot

A bar plot is used to visually display the dispersal of projected fraud and non-fraud instances in Credit Card fraud detection with the Random Forest. It aids in threshold choice and model evaluation by providing a visual representation of model's performance and the ability to point out imbalances.



**Fig -5:** Count plot of target variable

## 6. CONCLUSION

Finally, our credit card fraud detection project employing a Random Forest model offers a powerful and effective approach to mitigating fraud risks with a remarkable accuracy off 99.9%. The ability of the Random Forest model is to precisely identify fraud that has been demonstrated. Its ensemble learning method reduces false positives and false negatives by combining multiple decision trees to produce reliable results. Despite being well-known for their ability to predict outcomes, Random Forest models also provide some interpretability, which aids fraud investigators in understanding why particular transactions are reported as possibly fraudulent. Monitoring constantly emerging frauds is a continuous process in real time.

## 7. REFERENCES

[1]. Abdulalem Ali, Shukor Abd Razak, "Financial Fraud Detection Based on Machine Learning" (2022).

[2]. Mosa M.M. Megdad, Bassem S. Abu-Nasser and Samy S. Abu-Nazer, "Fraudulent Financial Transactions Detection Using Machine Learning" (2022).

[3]. Gupta, A., Lohani, M. C., & Manchanda, M. "Financial fraud detection using naive Bayes algorithm in highly imbalance data set". *Journal of Discrete Mathematical Sciences and Cryptography* (2021).

[4]. Saheed, Y. K., Hambali, M. A., Arowolo, M. O., & Olasupo, Y. A. Application of ga feature selection on Naive Bayes, random forest and SVM for credit card fraud detection.

International Conference on Decision Aid Sciences and Application (DASA) (2020).

[5]. Jemima Jebaseeli T, Venkatesan R, Ramalakshmi K. Fraud detection for credit card transactions using random forest algorithm. Singapore: Springer; (2020).

[6]. Priyanka Purushu, Jongwook Woo, "Financial Fraud Detection adopting Distributed Deep Learning in Big Data".

[7]. Adepoju, O., Wosowei, J., lawte, S., & Jaiman, H. Comparative evaluation of credit card fraud detection using machine learning techniques. *Global Conference for Advancement in Technology (GCAT)* (2019).

[8]. Aditya Oza, proposed a research paper "A Survey on Fraud Detection" in 2019. This paper applies different ML techniques on Logistic regression and Support vector Machine to the problem of payments fraud detection.

[9]. "Credit Card Fraud Detection: A Realistic Modelling and a Novel Learning Strategy" published by *IEEE Transactions on Neural networks and learning systems* (2018).

[10]. Xuan S, Liu G, Li Z, Zheng L, Wang S, Jiang C. Random forest for credit card fraud detection. *IEEE 15th international conference on networking, sensing and control (ICNSC)* (2018).

[11]. Wedge, CanterandRubio, "Solving the False positives problem in fraud prediction" (2017).

[12]. Samaneh Sorournejad, Zojah and Atani, "A survey of Credit Card Fraud Detection Techniques" (2016).

[13]. T. Singh, F. DiTroia, C. Vissagio, "Support Vector Machines and malware detection" (2015).

[14]. "A comparative analysis of credit card fraud detection using neural network, decision tree and logistic regression" by Abraham, A., Mithun, N. P., & Soman, K. P.

[15]. "Detection of Financial Statement Fraud: A Review of the Literature" by Elliott, R. K., & Willingham, J. J.