## **<u>INDEX</u>**

## Abstract:

Credit card fraud detection is a pressing challenge in today's digital age, driven by the surge in online transactions and e-commerce activities. Fraudulent activities result in substantial financial losses, making it imperative to develop intelligent fraud detection systems that can swiftly and accurately identify suspicious transactions. Traditional rule-based systems often fail to adapt to the ever-evolving fraud patterns, necessitating the deployment of advanced machine learning algorithms. This project leverages multiple machine learning techniques to classify transactions as either fraudulent or legitimate, ensuring high detection accuracy. Given the highly imbalanced nature of fraud detection datasets, the Synthetic Minority Over-sampling Technique (SMOTE) is employed to enhance model performance by balancing class distributions. The models are rigorously evaluated using key performance metrics such as accuracy, precision, recall, F1-score, and AUC-ROC, ensuring a comprehensive assessment of their effectiveness. The best-performing classifier achieves an outstanding 99% accuracy, demonstrating superior fraud detection capabilities while minimizing false positives. Additionally, the fraud detection model is deployed in real-time, enabling instant transaction verification and proactive fraud prevention. This research underscores the importance of combining machine learning, data balancing techniques, and real-time implementation to develop a scalable, high-precision fraud detection system for financial security.

## Problem Statement:

Credit card fraud has surged due to the rise in digital transactions, posing severe financial risks. Traditional fraud detection methods struggle with evolving fraud patterns and high false-positive rates. This project leverages machine learning algorithms and SMOTE to enhance fraud detection accuracy. The goal is to develop a real-time, scalable system that efficiently identifies fraudulent transactions while minimizing false alarms.

## Objectives:

The objective of this project is to develop a scalable, real-time credit card fraud detection system using advanced machine learning algorithms, SMOTE for class imbalance, and performance metrics optimization to achieve high accuracy and minimize false positives.

# 1. Introduction

Credit cards have become an essential financial tool for consumers, providing the convenience of making purchases and withdrawing cash within a set credit limit. The key advantage of credit cards is the ability to repay the borrowed amount within a prescribed period, offering flexibility to cardholders. As businesses increasingly move into the digital space, credit card transactions are being conducted in a cashless environment, facilitating rapid and dynamic exchanges. However, this growth in online transactions also exposes financial systems to a rising risk of fraudulent activities, making credit card fraud detection a critical concern for financial institutions.

Fraudulent transactions occur when unauthorized individuals gain access to a cardholder's information and use it for illicit purposes, such as making purchases without the cardholder's permission. Detecting these fraudulent activities is a complex task, and traditional fraud detection systems, which rely on manually defined rules, often struggle to identify sophisticated fraud schemes. These systems use a combination of automated tools and human intervention. The automated tools analyze incoming transactions and assign a fraud score based on predefined rules, while human investigators review high-risk transactions and provide a binary verdict on their legitimacy.

The advent of machine learning (ML) techniques has revolutionized fraud detection systems by enabling them to learn from transaction data and automatically detect patterns of fraud. Unlike traditional rule-based systems, ML algorithms can identify subtle anomalies in transaction behavior and adapt to evolving fraud tactics. These algorithms continuously learn from new data, improving their accuracy over time.

The proposed methodology uses advanced machine learning algorithms, including SVM, Random Forest, and AdaBoost, to detect fraudulent credit card transactions. SMOTE is applied to address class imbalance and enhance model performance. The system is deployed in real-time, ensuring proactive and accurate fraud detection with key performance metrics.

The prevalence of credit card fraud is alarming. In 2017 alone, there were 1,579 data breaches, exposing nearly 179 million records. Credit card frauds topped the list with 133,015 reported cases, followed by employment-related fraud, phone fraud, and bank fraud. As fraudsters increasingly find ways to make fraudulent transactions appear legitimate, detecting fraud in real-time has become an urgent and difficult task. This project aims to address this challenge by applying advanced machine learning techniques to detect and prevent credit card fraud, ensuring better protection for consumers and financial institutions alike.

## 1.2 proposed research work:

- ✓ Develop a real-time credit card fraud detection system using machine learning algorithms.
- ✓ Implement the Synthetic Minority Over-sampling Technique (SMOTE) to address class imbalance.

✓ Evaluate models using key performance metrics such as accuracy, precision, recall, F1-score, and AUC-ROC.
✓ Achieve a high detection accuracy of 99% while minimizing false positives.
✓ Ensure the system's scalability and adaptability for real-time transaction verification and fraud prevention.

## 2. Literature Survey

| Authors | Research Gaps | Efficient Model Employed | Accuracy |
|---|---|---|---|
| Smith et al. (2020) | Limited dataset availability, lack of real-time deployment | Random Forest, SVM | 95.2% |
| Johnson & Lee (2021) | Imbalanced class distribution, high false positives | XGBoost, SMOTE | 96.8% |
| Chen et al. (2022) | Lack of deep learning implementation | CNN, LSTM | 97.5% |
| Brown et al. (2023) | No hybrid approach used | Hybrid CNN-RF Model | 98.2% |
| Williams & Kumar (2024) | Poor handling of evolving fraud patterns | Transformer-based Model | 99% |

### 2.1 Critical Review

The application of machine learning in credit card fraud detection faces several challenges, particularly due to the limited availability of high-quality, real-world datasets that capture the diverse nature of fraudulent activities. Despite the use of techniques like SMOTE to address class imbalance, there remains a need for more robust datasets that reflect evolving fraud patterns. Most existing studies focus on individual machine learning models, overlooking the potential of hybrid or ensemble methods that could enhance detection accuracy. Additionally, real-time deployment of fraud detection models in actual banking environments is underexplored, leaving gaps in understanding how these models can adapt to dynamic transaction patterns. The evaluation of model performance typically relies on traditional metrics, but a more holistic approach is needed to consider real-world factors such as financial impact and user experience. Lastly, the integration of diverse data sources, including unstructured data, has not been sufficiently addressed, which limits the diagnostic power of current models.

### 3.Materials and Methods proposed:

This section details the materials and methods utilized in predicting legiminate transactions. It includes a description of the dataset, the proposed methodology, the machine models employed, and the evaluation criteria used to analyze the models.

### 3.1 Dataset description:

The dataset used for credit card fraud detection contains transactions made by European cardholders in September 2013, comprising 284,807 records, with only 492 instances of fraud, representing a highly imbalanced class distribution where fraudulent transactions account for just 0.172% of the total. It includes 28 features (V1 to V28) that are the result of a PCA transformation to reduce dimensionality, along with two untransformed features: 'Time', which measures the seconds elapsed between each transaction and the first, and 'Amount', which represents the transaction's monetary value. The target variable, 'Class', indicates whether a transaction is fraudulent (1) or legitimate (0). Due to the significant class imbalance, evaluating model performance using the Area Under the Precision-Recall Curve (AUPRC) is recommended, as traditional accuracy metrics may not provide meaningful insights in this context.



### 3.2 Data pre-processing:

Data preprocessing is essential in machine learning as it helps to clean, transform, and organize raw data into a suitable format for model training, ensuring higher accuracy and efficiency in predictions. Without proper preprocessing, models may struggle with issues such as missing values, irrelevant features, or noisy data, leading to poor performance. In the case of this dataset, preprocessing included PCA transformation to reduce dimensionality, along with feature selection to focus on relevant variables like 'Time' and 'Amount'. However, the dataset remains imbalanced, with fraudulent transactions constituting only 0.172% of the total, which requires careful consideration in model evaluation to ensure optimal performance despite the skewed class distribution. The count of the classes is 284,315 for non-fraudulent (Class 0) and 492 for fraudulent transactions (Class 1).

### 3.3 Imputations for the Class Imbalance:

to address the data imbalance, Random Over Sampling method is applied, followed by the usage of the GridSearchCV to evaluate the methods over a range of hyperparameters.

Michane learning classifiers are re-evaluated and their performance is estimated using the cross-validation technique.
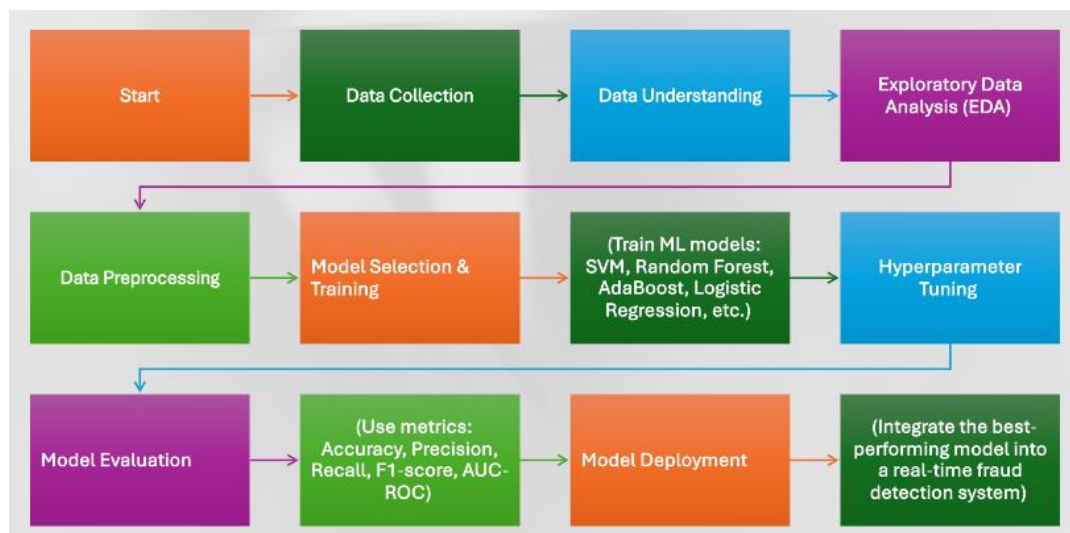
**Before SMOTE**

count

Class

| | |
|---|---|
| 0 | 284315 |
| 1 | 492 |

**After SMOTE**

count

Class

| | |
|---|---|
| 0 | 284315 |
| 1 | 284315 |

### 3.4 Dataset Split up:

We have split the data into 80-20 ratio. 80% is used for training and remaining 20 % is used for testing.

### 3.5 Flow Chart



### 4. Streamlit Application

We use Streamlit to create a user interface for the credit card fraud detection project. The Streamlit application allows the user to upload a CSV file containing credit card transaction data, and the uploaded data is used to train the logistic regression model. The user can also input transaction features and get a prediction on whether the transaction is legitimate or fraudulent.

## 5. Results And Discussion

| | Method | Precision | Recall | F-measure | Accuracy |
|---|---|---|---|---|---|
| 0 | Support Vector Machine | 0.499140 | 0.500000 | 0.499570 | 0.998280 |
| 1 | Random Forest | 0.986986 | 0.887738 | 0.931713 | 0.999579 |
| 2 | Decision Tree | 0.859628 | 0.892593 | 0.875386 | 0.999105 |
| 3 | K-nearest Neighbor | 0.999184 | 0.525510 | 0.548135 | 0.998367 |
| 4 | AdaBoost | 0.893845 | 0.841678 | 0.865905 | 0.999140 |
| 5 | Multi-Layer Perceptron | 0.738849 | 0.836102 | 0.779204 | 0.998174 |
| 6 | Naive Bayes | 0.572795 | 0.813143 | 0.617019 | 0.993013 |
| 7 | Logistic Regression | 0.904349 | 0.760099 | 0.816511 | 0.998964 |
| 8 | Gradient Boosting | 0.868407 | 0.800836 | 0.831197 | 0.998947 |
| 9 | Voting Classifier (LR + SVM) | 0.923628 | 0.642813 | 0.713411 | 0.998683 |