# Credit Card Fraud Detection

## Introduction

In the rapidly evolving landscape of digital transactions, financial fraud has become a critical issue for both consumers and financial institutions. The rise of online banking, e-commerce, and mobile payment systems has led to an increase in fraudulent activities, costing billions of dollars annually. Detecting and preventing such fraud is a significant challenge, given the sophisticated methods used by fraudsters. Machine learning (ML) has emerged as a powerful tool to address this challenge, offering the ability to analyze vast amounts of transaction data, identify patterns, and predict fraudulent activities.

This report presents a solution for detecting credit card fraud using advanced machine learning algorithms. The approach involves the use of synthetic data generation to address the class imbalance problem, combined with the application of various machine learning algorithms to predict potential fraud. The effectiveness of this approach is demonstrated through the achievement of high ROC-AUC scores and a significant increase in recall.

## Problem Statement

Financial transactions, especially credit card transactions, are highly imbalanced in nature, with fraudulent transactions representing only a small fraction of the total. This imbalance poses a significant challenge for machine learning models, which may become biased toward the majority class, leading to poor performance in detecting fraud. Therefore, an essential step in developing an effective fraud detection system is to address the class imbalance issue.

The objective of this project is to develop a robust fraud detection model that can accurately predict fraudulent transactions while minimizing false positives. The model should be capable of handling imbalanced datasets and should be able to generalize well to unseen data.

## Methodology

1. Data Collection and Preprocessing

The first step in the process involved collecting a dataset of credit card transactions. The dataset included various features such as transaction amount, time of transaction, and other behavioral patterns that could indicate fraudulent activity. Each transaction was labeled as either fraudulent or legitimate.

Given the highly imbalanced nature of the dataset, with legitimate transactions vastly outnumbering fraudulent ones, preprocessing was necessary to prepare the data for model training. The preprocessing steps included:

Data Cleaning: Removal of any missing or irrelevant data to ensure the dataset was clean and ready for analysis.
Feature Engineering: Creation of new features that could help the model better distinguish between fraudulent and legitimate transactions.

Data Scaling: Standardization of numerical features to ensure that all features contributed equally to the model's performance.

## 2. Addressing Class Imbalance with SMOTE

To address the class imbalance, the Synthetic Minority Over-sampling Technique (SMOTE) was employed. SMOTE is a widely used oversampling technique that generates synthetic instances of the minority class by interpolating between existing minority instances. This helps to balance the dataset by increasing the representation of the minority class, which in this case is the fraudulent transactions.

SMOTE works by selecting a random sample from the minority class and identifying its k-nearest neighbors. Synthetic instances are then created by interpolating between the selected sample and one of its neighbors. This approach helps to create a more balanced dataset, allowing the machine learning models to learn from both classes more effectively.

## 3. Model Selection

Several machine learning algorithms were employed to build the fraud detection model. Each of these algorithms has its strengths and was chosen to complement the others, providing a comprehensive approach to detecting fraud.

Random Forest: Random Forest is an ensemble learning method that operates by constructing multiple decision trees during training and outputting the mode of the classes as the prediction. It is known for its robustness and ability to handle large datasets with high dimensionality. Random Forest was chosen for its ability to handle both continuous and categorical data and its resilience to overfitting.

XGBoost: Extreme Gradient Boosting (XGBoost) is a powerful and efficient implementation of gradient boosting algorithms. It builds models in a stage-wise manner and optimizes for accuracy and speed. XGBoost is particularly effective in dealing with imbalanced datasets due to its ability to focus on hard-to-classify examples during the training process.

K-Means Clustering: K-Means is an unsupervised learning algorithm used for clustering data into k distinct clusters based on feature similarity. For fraud detection, K-Means was used to group similar transactions together, with the hypothesis that fraudulent transactions would form distinct clusters due to their differing characteristics compared to legitimate transactions.

DBSCAN: Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is another unsupervised learning algorithm that groups together points that are closely packed together, marking as outliers the points that are in low-density regions. DBSCAN is particularly useful for detecting anomalies, making it well-suited for identifying fraudulent transactions that do not conform to the patterns of legitimate transactions.

## 4. Model Training and Evaluation

The models were trained using the preprocessed dataset, with SMOTE applied to ensure balanced training data. The performance of each model was evaluated using various metrics, including precision, recall, F1-score, and the Receiver Operating Characteristic Area Under the Curve (ROC-AUC) score.

The ROC-AUC score is a critical metric for evaluating the performance of fraud detection models, as it measures the model's ability to distinguish between classes across different thresholds. A higher ROC-AUC score indicates better performance, with a score of 1 representing a perfect model.

**5. Results**
The performance of the models was as follows:

Random Forest: The Random Forest model achieved a ROC-AUC score of 0.90. While this is a strong performance, indicating that the model is effective at distinguishing between fraudulent and legitimate transactions, there was still room for improvement in terms of recall.

XGBoost: XGBoost performed slightly better than Random Forest, with a ROC-AUC score of 0.91. This improvement is likely due to XGBoost's ability to handle imbalanced datasets more effectively by focusing on harder-to-classify instances.

K-Means Clustering: K-Means, as an unsupervised algorithm, did not perform as well as the supervised models, with a lower ROC-AUC score. However, it was useful in identifying distinct clusters of fraudulent transactions, providing valuable insights into the nature of the fraud.

DBSCAN: The DBSCAN algorithm outperformed the other models with a ROC-AUC score of 0.93. This significant improvement can be attributed to DBSCAN's ability to detect outliers and anomalies effectively, which is crucial in identifying fraudulent transactions that deviate from normal patterns.

The DBSCAN model also achieved a 6% increase in recall compared to the Random Forest model, demonstrating its effectiveness in reducing false negatives and improving the detection of fraudulent transactions.


**Conclusion**
In conclusion, this report presents a solution for detecting credit card fraud using a combination of advanced machine learning algorithms and synthetic data generation techniques. The approach addresses the challenges of imbalanced datasets and demonstrates the effectiveness of models like Random Forest, XGBoost, K-Means, and DBSCAN in predicting fraudulent transactions. The high ROC-AUC scores achieved, particularly with DBSCAN, underscore the potential of these methods in safeguarding financial transactions against fraud.

Future work could focus on further refining these models, exploring additional features, and integrating real-time data to enhance the detection capabilities. As fraud tactics continue to evolve, the use of advanced machine learning techniques will be crucial in maintaining the security and integrity of financial systems.