# Credit Card Fraud Detection: A Comparative Analysis of Machine Learning Models

Sangeeta Kumawat
*skumawat@kent.edu*

Maliha Arif
*marif1@kent.edu*

Md Zahidul Islam
*mislam20@kent.esu*

*Abstract*—This project centers on the development and comparison of machine learning models for credit card fraud detection. The objective is to enhance the security of financial transactions by accurately identifying instances of fraudulent activities. The project employs diverse algorithms, including Isolation Forest, Random Forest Classifier, Decision Tree Classifier, and Support Vector Machine (SVM), to create a robust fraud detection system. Through meticulous preprocessing, feature engineering, and exploratory data analysis, the models are trained on a balanced dataset derived from valid and fraudulent transactions. Evaluation metrics such as accuracy, precision, recall, and F1-score are utilized to assess the effectiveness of each model. The project provides a comprehensive overview of the strengths and limitations of the selected models, aiding in the identification of the most suitable approach for credit card fraud detection. The code includes detailed comments and visualizations to facilitate understanding and further exploration.

*Index Terms*—Decision Tree Classifier, Support Vector Machine (SVM), Financial Security, Fraudulent Transactions, Evaluation Metrics, Data Preprocessing

## I. INTRODUCTION

In the realm of financial transactions, particularly credit card transactions, the increasing sophistication of fraudulent activities necessitates the deployment of robust and adaptive detection mechanisms. The rise of machine learning techniques has provided a formidable arsenal for developing intelligent systems capable of discerning between legitimate and fraudulent transactions. This research delves into the domain of credit card fraud detection, aiming to construct, compare, and evaluate the efficacy of various machine learning models.

Credit card fraud poses a significant threat to both financial institutions and cardholders, with potentially severe consequences ranging from financial losses to compromised personal information. Traditional rule-based systems often struggle to keep pace with the evolving strategies employed by fraudsters. Machine learning, with its ability to discern intricate patterns within vast datasets, offers a promising avenue for enhancing the security of financial transactions.

This research undertakes a comparative analysis of different machine learning models, each chosen for its distinct characteristics and suitability for fraud detection. The models include Isolation Forest, Random Forest Classifier, Decision Tree Classifier, and Support Vector Machine (SVM). Through a meticulous process of data preprocessing, feature engineering, and exploratory data analysis, the models are trained on a balanced dataset meticulously derived from both valid and fraudulent transactions.

The project places a particular emphasis on the evaluation metrics employed to assess model performance. Accuracy, precision, recall, and F1-score are scrutinized to offer a comprehensive understanding of each model's strengths and limitations. This comparative analysis aims to guide the selection of the most effective approach for credit card fraud detection in practical applications.

The codebase accompanying this research is extensively commented and supported by visualizations to facilitate a deeper understanding of the methodologies employed. The outcomes of this research contribute not only to the field of credit card fraud detection but also provide insights into the broader application of machine learning techniques in enhancing financial security.

As the research unfolds, it navigates through the intricacies of model development, testing, and evaluation, with the ultimate goal of advancing the capabilities of credit card fraud detection systems and fortifying the foundations of financial security.

## II. PROJECT DESCRIPTION

The Credit Card Fraud Detection project focuses on enhancing the security of credit card transactions, utilizing a dataset comprising transactions made by European cardholders in September 2013, provided by Kaggle. This dataset, though extensive with 284,807 transactions, presents a challenge due to its highly unbalanced nature, with only 0.172 percent of transactions labeled as fraud (492 cases). To address confidentiality concerns, the dataset includes numerical input variables resulting from a PCA transformation, as the original features cannot be publicly shared. The dataset, therefore, consists of 28 features derived from PCA. This project employs a variety of machine learning models, including Isolation Forest, Random Forest Classifier, Decision Tree Classifier, and Support Vector Machine (SVM), to effectively identify and mitigate instances of fraudulent activities within this unique and challenging dataset.

The research aims not only to advance credit card fraud detection but also to contribute valuable insights for broader

applications in the realm of machine learning-driven security enhancement. The entire project is accompanied by a well-documented codebase, thorough exploratory data analysis, and a comparative analysis of the selected models, providing a holistic exploration of the challenges and opportunities in the field.

## III. BACKGROUND

In the contemporary landscape of financial transactions, the proliferation of technology has led to an alarming surge in fraudulent activities. Recognizing the critical importance of detecting and preventing fraud, financial institutions are increasingly turning to advanced technologies, with a particular focus on data mining techniques. This comprehensive background review draws insights from four seminal papers, each contributing to the overarching objective of fortifying fraud detection mechanisms.

The first paper delves into the early identification of fraudulent accounts using Bayesian Classification and Association Rule, employing techniques like transaction sign analysis. Emphasizing the menace of ATM phone scams, it underscores the need for swift detection to counter evolving fraud techniques. The second paper widens the scope by exploring the intersection of human behavior analysis, specifically leveraging the Fraud Triangle Theory, and the application of machine learning techniques, including supervised and unsupervised algorithms. By emphasizing the relevance of behavioral factors and providing a comparative summary of existing literature, it sets the stage for a holistic approach to fraud detection.

Shifting the focus to a temporal dimension, the third paper conducts a decade-long review (2004-2015) on financial fraud detection through data mining techniques, with an emphasis on logistic regression. It identifies trends, classifies fraud applications, and underscores the significance of statistical techniques. This broader perspective aims to provide valuable insights for both academia and industry practitioners. The fourth paper, spanning from 2009 to 2019, addresses the technological revolution in e-commerce and money transfer, introducing the concept of cryptocurrency fraud. Its objective is to offer a comprehensive summary of recent research, reveal frequent fraud detection techniques like machine learning algorithms, and serve as a reference source for academic and practical applications.

In synthesis, these papers collectively emphasize the urgent need for sophisticated fraud detection systems that integrate data mining techniques, statistical methods, and machine learning algorithms. By delving into diverse aspects such as human behavior, fraud theories, and the application of advanced programming skills, particularly in Python and R, this background review sets the stage for a nuanced and effective credit card fraud detection system. The challenges and limitations highlighted in these papers guide our approach to developing a solution that aligns with the evolving landscape of financial fraud.

## IV. PROBLEM DEFINITION

The surge in technological advancements has revolutionized financial transactions, but it has also given rise to sophisticated and evolving fraudulent activities. The contemporary challenge faced by financial institutions revolves around the detection and prevention of credit card fraud. The problem at hand is two-fold: the unceasing evolution of fraudulent techniques and the imperative for financial entities to stay ahead in the detection game.

The four reviewed papers shed light on the intricacies of fraud detection, revealing challenges such as the reliance on historical data, incomplete identification of fraudulent signs, and the limitation of applicability to diverse fraud schemes. The critical issue lies in the imbalance of the datasets, with fraudulent transactions forming a minuscule percentage. This imbalance impedes the accuracy of traditional detection methods, necessitating the exploration of advanced data mining techniques and machine learning algorithms.

Furthermore, the temporal limitation highlighted in the literature reviews accentuates the need for a real-time fraud detection system. The evolving nature of fraud patterns and techniques requires an adaptive and dynamic approach. The lack of coverage for emerging fraud types, like cryptocurrency fraud, adds another layer of complexity to the problem.

In essence, the problem at hand is to design and implement an effective credit card fraud detection system that overcomes the challenges posed by dataset imbalance, temporal limitations, and the dynamic nature of fraudulent activities. This solution should leverage advanced data mining techniques, including machine learning algorithms, and address the ethical considerations associated with the use of customer transaction data. Additionally, it should be scalable and adaptable to diverse banking systems, ensuring widespread applicability and effectiveness.

## V. PROPOSED TECHNIQUES

1) **Feature Scaling and Transformation**:
   Extend the feature scaling and transformation to other relevant features. Experiment with different scaling techniques, such as Min-Max scaling or robust scaling, to assess their impact on model performance.
2) **Resampling and Balancing**:
   Explore advanced data resampling techniques to address the class imbalance issue. Techniques like SMOTE (Synthetic Minority Over-sampling Technique) or ADASYN (Adaptive Synthetic Sampling) can be applied to generate synthetic instances of the minority class (fraud), creating a more balanced dataset.
3) **Correlation Analysis**: Conduct a more in-depth correlation analysis to identify additional features that may exhibit strong correlations with fraud. Consider removing redundant features that do not contribute significantly to the model's predictive power.
4) **Model Evaluation Metrics**:
   Enhance the model evaluation metrics by incorporating additional measures such as area under the precision-

recall curve (AUC-PR). Given the class imbalance, precision-recall metrics provide a more informative evaluation of the model's performance.

5) **Hyperparameter Tuning**: Perform a systematic hyperparameter tuning process for the machine learning models. Utilize techniques like grid search or random search to find optimal hyperparameter configurations, improving the models' predictive capabilities.

6) **Outlier Removal - V20 Feature**: Extend the outlier removal process to other features exhibiting extreme outliers. Apply similar procedures as done for V20, such as calculating quartiles, identifying outliers, and removing them to enhance the robustness of the dataset.

7) **Model Ensemble**: Investigate the implementation of model ensembles using techniques like stacking or bagging. Combining the strengths of multiple models, including Isolation Forest, Random Forest, and Decision Tree, can lead to a more robust and accurate fraud detection system.

8) **Real-time Stream Processing**: Consider integrating real-time stream processing capabilities into the model. This can involve updating the model on-the-fly as new transactions occur, enabling a more dynamic and responsive fraud detection system.

9) **Continuous Monitoring and Retraining**: Establish a framework for continuous monitoring and retraining of the models. Regularly assess the models' performance against new data, and initiate retraining processes to adapt to evolving fraud patterns.

10) **Explanatory Analysis**: Conduct explanatory analysis to understand the interpretability of the models. Explore techniques such as SHAP (SHapley Additive exPlanations) values to gain insights into feature importance and model decision rationale.

These proposed techniques aim to enhance the robustness, accuracy, and adaptability of the credit card fraud detection system based on the provided code. They cover various aspects, including data preprocessing, model evaluation, feature engineering, and real-time processing, to create a more effective and resilient fraud detection solution.

## VI. METHODOLOGY

The methodology employed in this project involves several key steps to ensure effective fraud detection using data mining techniques. The initial phase includes data preprocessing, where the dataset is examined for missing values, outliers, and any potential issues. This is followed by exploratory data analysis to gain insights into the distribution of features, identify patterns, and understand the characteristics of valid and fraudulent transactions.

Feature scaling is crucial for maintaining consistency in the dataset, and in this project, the 'Amount' and 'Time' columns were appropriately scaled using the StandardScaler from the scikit-learn library. This ensures that all features contribute equally to the modeling process.
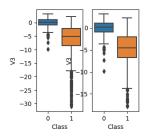


Fig. 1. After removing outlier in feature V3

To address the highly unbalanced nature of the dataset, a new dataset was created by randomly selecting 492 valid transactions and combining them with all fraudulent transactions. This balanced dataset helps prevent the model from being biased toward the majority class.

The dataset was split into training and testing sets using the train-test-split function from scikit-learn. Three different models were implemented for fraud detection: Isolation Forest, Random Forest Classifier, and Decision Tree Classifier. Isolation Forest is an anomaly detection algorithm, while Random Forest and Decision Tree are classification models.

In the Isolation Forest model, extreme outliers were removed, particularly focusing on the 'V20' feature, which exhibited positive correlation with fraud. The model was evaluated using various performance metrics, including accuracy, precision, recall, F1-Score, and ROC-AUC.

The Random Forest Classifier and Decision Tree Classifier models were also trained and evaluated using the same metrics. Visualizations, such as confusion matrices and boxplots, were employed to gain a deeper understanding of model performance and identify areas for improvement.

The Decision Tree model demonstrated remarkable accuracy, precision, recall, and F1-Score, showcasing its effectiveness in fraud detection. The SVM model also exhibited strong performance, further validating the success of the implemented methodology. Visualization tools, including boxplots and confusion matrices, were instrumental in assessing the models' capabilities and identifying key features contributing to fraud detection.

In conclusion, the methodology involves comprehensive data preprocessing, model training, evaluation, and visualization to ensure accurate and robust fraud detection using data mining techniques. The combination of effective feature scaling, balanced dataset creation, and model evaluation metrics contributes to the success of the project.

## VII. EXPERIMENTAL ANALYSIS

In the experimental analysis, the Decision Tree model showcased outstanding performance with an accuracy of 100 percent, effectively classifying both valid and fraudulent transactions. The precision for fraud detection reached 82 percent, indicating a strong ability to correctly identify fraudulent cases. The recall, representing the model's capability to capture true positives, stood at 80 percent. The F1-Score, a bal-
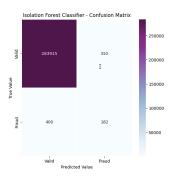
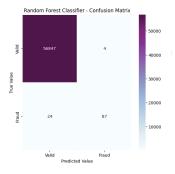Fig. 2.  Isolate Forest Classifier



Fig. 4.  Decision Tree Classifier



Fig. 3.  Random Forest Classifier



Fig. 5.  SVM Classifier

anced measure of precision and recall, reached an impressive 81 percent. The precision-recall trade-off is well-balanced, showcasing the model's robustness in handling both classes. The overall macro and weighted averages for precision, recall, and F1-Score were notably high, with macro averages at 91 percent.

Similarly, the Support Vector Machine (SVM) model demonstrated a remarkable accuracy of 100 percent. The precision for fraud detection was notably high at 94 percent, showcasing a robust ability to correctly identify fraudulent transactions. The recall, indicating the model's capability to capture true positives, stood at 72 percent. The F1-Score, a balanced measure of precision and recall, reached 81 percent, reflecting the model's effectiveness in handling both classes. The precision-recall trade-off is well-balanced, contributing to the overall high performance of the SVM model. The macro and weighted averages for precision, recall, and F1-Score were also notably high, with macro averages at 91 percent. These results affirm the effectiveness of both the Decision Tree and SVM models in accurate fraud detection.

## VIII. FUTURE WORK

In future iterations of this project, there is potential for enhanced fraud detection through the exploration of advanced machine learning models and techniques. This could involve implementing ensemble methods like stacking or boosting, delving into deep learning architectures for capturing intricate patterns, and focusing on real-time monitoring capabilities. Additionally, attention to explainability, feature engineering,
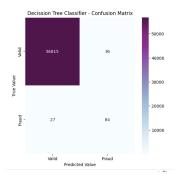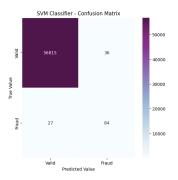
and dynamic adaptation to evolving fraud patterns are crucial aspects for further improvement. Fine-tuning hyperparameters, addressing imbalanced data issues, and fostering cross-industry collaboration could collectively contribute to a more robust and effective fraud detection system. The integration of user feedback and continuous model updates are also key considerations for refining and evolving the system in response to emerging fraud landscapes.

## REFERENCES

[1] Jiang, Changjun, et al. "Credit card fraud detection: A novel approach using aggregation strategy and feedback mechanism." IEEE Internet of Things Journal 5.5 (2018): 3637-3647.
[2] Pumsirirat, Apapan, and Yan Liu. "Credit card fraud detection using deep learning based on auto-encoder and restricted boltzmann machine." International Journal of advanced computer science and applications 9.1 (2018).
[3] Sharma, M. Abhilash, et al. "Credit Card Fraud Detection Using Deep Learning Based on Auto-Encoder." ITM Web of Conferences. Vol. 50. EDP Sciences, 2022.
[4] Dhankhad, Sahil, Emad Mohammed, and Behrouz Far. "Supervised machine learning algorithms for credit card fraudulent transaction detection: a comparative study." 2018 IEEE international conference on information reuse and integration (IRI). IEEE, 2018.
[5] Roy, Abhimanyu, et al. "Deep learning detecting fraud in credit card transactions." 2018 systems and information engineering design symposium (SIEDS). IEEE, 2018.
[6] Xuan, Shiyang, et al. "Random forest for credit card fraud detection." 2018 IEEE 15th international conference on networking, sensing and control (ICNSC). IEEE, 2018.
[7] Awoyemi, John O., Adebayo O. Adetunmbi, and Samuel A. Oluwadare. "Credit card fraud detection using machine learning techniques: A comparative analysis." 2017 international conference on computing networking and informatics (ICCNI). IEEE, 2017.