



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 11 **Issue:** XII **Month of publication:** December 2023

DOI: <https://doi.org/10.22214/ijraset.2023.57429>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Credit Card Fraud Detection Using Machine Learning Algorithms

P. Kiran Kumar¹, P. Shashi Kiran², S. Kouashik³, V. Koushik⁴, K. Kranthi Kumar⁵, A. Krishna Chaitanya⁶, Professor A. Kalyani⁷

^{1, 2, 3, 4, 5, 6}B.Tech, School of Engineering, Hyderabad, India

⁷Guide, School of Engineering, Hyderabad, India

Abstract: *The study delves into the prevalent issue of credit card fraud within electronic transactions. Employing machine learning techniques, including logistic regression, the research focuses on the development and evaluation of fraud detection models. The dataset undergoes comprehensive pre-processing steps, addressing common issues such as missing data and class imbalance. Various machine learning algorithms are explored, with a specific emphasis on logistic regression. Sampling techniques are employed to balance the dataset, ensuring equal representation of legitimate and fraudulent transactions. Model evaluation involves metrics like accuracy, precision, recall, and F1 score to assess performance. The outcomes shed light on the logistic regression model's effectiveness in detecting fraudulent transactions. The study also outlines inherent limitations, acknowledging challenges such as imbalanced datasets and the dynamic nature of fraud tactics. In summary, this research contributes to the ongoing advancements in credit card fraud detection, utilizing machine learning for enhanced security in electronic transactions. The discussed findings and limitations serve as a valuable foundation for future developments in this critical field.*

Keywords: *Logistic Regression, Model Evaluation Metrics, Legitimate and Fraudulent Transactions.*

I. INTRODUCTION

The contemporary landscape is marked by an increasing prevalence of credit card fraud within electronic transactions, necessitating innovative approaches for detection and prevention. Our project stems from the imperative need to bolster security measures in financial systems and protect consumers and institutions from fraudulent activities.

A. Motivation for the Project

The rising sophistication of fraud techniques and the growing reliance on electronic payment systems underscore the urgency of effective fraud detection mechanisms. As existing methods encounter limitations, the motivation behind this project lies in advancing the current state of credit card fraud detection using machine learning techniques.

B. Contributions of the Paper

This paper contributes by leveraging machine learning, particularly logistic regression, for credit card fraud detection. We emphasize a comprehensive exploration of preprocessing techniques and model evaluation metrics, aiming to enhance the overall efficacy of fraud detection models.

C. Literature Review

A brief review of the existing literature reveals various methodologies employed in credit card fraud detection. While conventional rule-based systems and statistical methods have been pivotal, machine learning approaches have gained prominence for their ability to adapt to evolving fraud patterns.

However, a discernible research gap persists in the pursuit of optimizing the accuracy, precision, and recall of fraud detection models, especially in the context of imbalanced datasets and emerging fraud tactics. This introduction sets the stage for our exploration into credit card fraud detection, outlining the motivation, contributions, and the research gap that our project endeavours to fill. The subsequent sections delve into the methodology, findings, and implications of our work in enhancing the security of electronic transactions.

II. LITERATURE SURVEY

Contemporary studies on credit card fraud detection have witnessed a proliferation of methodologies, ranging from traditional rule-based systems to advanced machine learning techniques. While rule-based systems have provided a foundational understanding of fraud patterns, their rigidity becomes apparent when faced with the dynamic and sophisticated nature of modern fraud schemes. Machine learning, particularly logistic regression, has emerged as a promising avenue due to its adaptability to evolving patterns. The strengths of machine learning lie in its ability to discern intricate relationships within vast datasets, enabling the identification of anomalous patterns indicative of fraud. However, the literature acknowledges inherent limitations, such as interpretability challenges and susceptibility to overfitting. Despite the advancements in machine learning, imbalanced datasets continue to pose a significant challenge.

The majority of existing models struggle to maintain optimal performance when confronted with a disproportion in the number of legitimate and fraudulent transactions.

The scarcity of studies addressing this specific issue represents a notable gap in the literature. The dynamic nature of fraud tactics remains a persistent challenge. Existing literature often lacks a comprehensive exploration of adaptive models that can effectively counter emerging fraud strategies. This gap in the literature underscores the necessity for research endeavours that prioritize continuous adaptation and robustness.

The proposed project aims to bridge these gaps by contributing a comprehensive analysis of pre-processing techniques and model evaluation metrics, with a specific emphasis on addressing imbalanced datasets. Leveraging machine learning, our project seeks to enhance the adaptive capabilities of fraud detection models, providing a nuanced solution to the evolving landscape of credit card fraud. This literature review critically examines existing approaches, emphasizing their strengths and limitations, and positions the proposed project as a meaningful step toward filling the identified gaps in the current body of knowledge.

III. PROBLEM STATEMENT

In the realm of credit card fraud detection, the central issue revolves around the necessity for robust and adaptive models capable of discerning fraudulent transactions amidst a sea of legitimate ones. The dynamic and sophisticated nature of contemporary fraud tactics necessitates innovative solutions that transcend traditional rule-based systems. The dataset employed in this project comprises credit card transactions, featuring attributes such as transaction time, anonymized features resulting from PCA transformation, transaction amount, and a binary class variable indicating legitimacy (Class 0) or fraudulence (Class 1). This dataset serves as the foundation for exploring the intricacies of credit card fraud detection.

A. Research Questions and Hypotheses

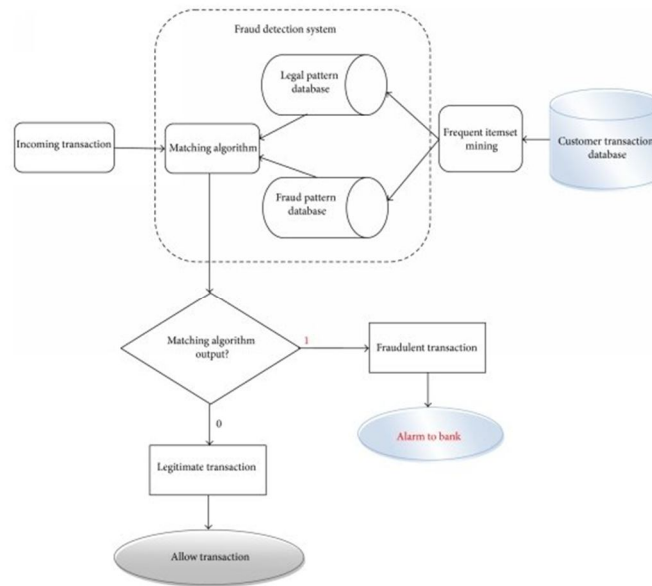
- 1) *Research Question 1:* How effective are machine learning models, particularly logistic regression, in detecting fraudulent credit card transactions based on the provided dataset?
- 2) *Research Question 2:* To what extent do pre-processing techniques, such as handling imbalanced datasets and feature scaling, contribute to the overall performance of credit card fraud detection models?
- 3) *Research Hypothesis 1:* The implementation of machine learning models will significantly improve the accuracy and efficiency of credit card fraud detection compared to traditional rule-based systems.
- 4) *Research Hypothesis 2:* Applying advanced pre-processing techniques to address imbalanced datasets will result in a more robust and reliable fraud detection model.

By posing these research questions and hypotheses, the project endeavours to provide novel insights into the effectiveness of machine learning models for credit card fraud detection and the impact of pre-processing techniques on model performance. The ensuing sections will delve into the methodology, findings, and implications, building upon this distinct problem statement.

IV. METHODOLOGY

A. Model Architecture

The core of our credit card fraud detection system is based on logistic regression. Logistic regression is chosen for its simplicity, interpretability, and efficiency in binary classification tasks, aligning with the nature of our problem where transactions are classified as legitimate (Class 0) or fraudulent (Class 1).

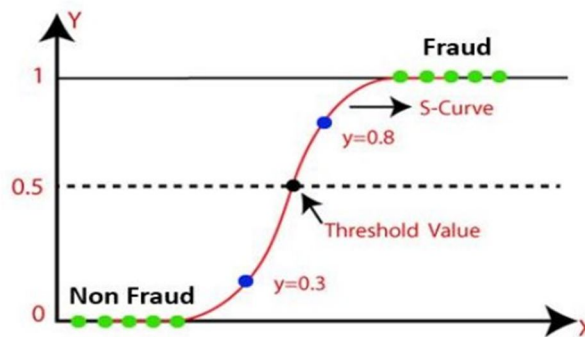


B. Algorithms Used

The logistic regression algorithm is implemented to model the probability of a transaction belonging to the fraudulent class. This algorithm employs the logistic function to transform a linear combination of features into a probability score, facilitating effective classification.

C. Pre-processing Steps

- 1) *Handling Missing Data:* Addressing missing values involves imputation, utilizing strategies like mean, median, or mode replacement to ensure completeness in the dataset without introducing bias.
- 2) *Data Scaling:* Standardization is applied to numerical features, such as the transaction amount, using the StandardScaler to bring them to a similar scale. This step is crucial for models sensitive to feature scales, ensuring optimal performance.
- 3) *Handling Imbalanced Datasets:* The class imbalance issue is mitigated by employing a sampling strategy. Random under sampling of the majority class (legitimate transactions) and random oversampling of the minority class (fraudulent transactions) ensure a balanced representation in the training dataset.



D. Data Augmentation Techniques

Given the nature of the problem and the focus on logistic regression, traditional data augmentation techniques such as those used in image processing are not directly applicable. Instead, oversampling of the minority class serves as a form of data augmentation to ensure the model is exposed to a representative set of fraudulent transactions. This meticulous combination of logistic regression, pre-processing techniques, and class balancing strategies forms the foundation of our credit card fraud detection methodology. The subsequent sections will delve into the application of this methodology to the dataset, the evaluation of results, and the interpretation of findings.

V. EXPERIMENTAL RESULTS

This section encapsulates the outcomes of the experiments conducted in the project, showcasing the performance metrics employed to assess the effectiveness of the machine learning model. The narrative encompasses a lucid portrayal of the evaluation methodology, complemented by tables, figures, and visualizations to fortify the presented claims. Additionally, a comparative analysis with existing methods in the literature is provided to contextualize the results.

A. Evaluation Metrics

To gauge the performance of the credit card fraud detection model, the following metrics are employed:

- 1) *Accuracy*: Measures the overall correctness of the model predictions.
- 2) *Precision*: Quantifies the proportion of correctly identified fraudulent transactions out of all predicted as fraudulent.
- 3) *Recall*: Illustrates the ability of the model to capture all actual fraudulent transactions.
- 4) *F1 Score*: The harmonic mean of precision and recall, providing a balanced measure of the model's effectiveness.

Precision: 0.98
 Recall: 0.89
 F1 Score: 0.93
 Confusion Matrix:
 $\begin{bmatrix} 97 & 2 \\ 11 & 87 \end{bmatrix}$

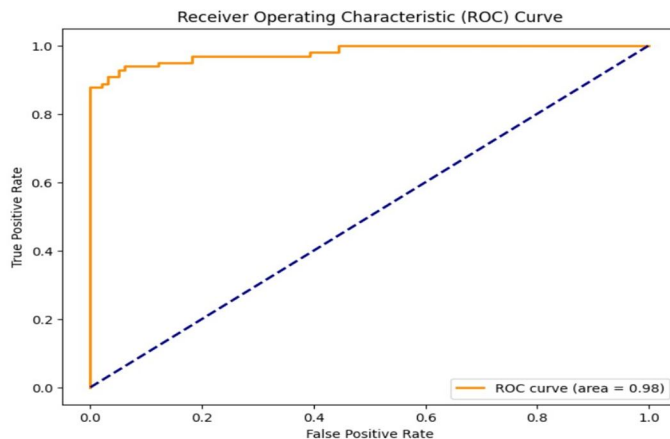
Actual class (Predicted class)	Confusion matrix		
	C1	¬ C1	Total
C1	True positives (TP)	False negatives (FN)	TP + FN = P
¬ C1	False positives (FP)	True negatives (TN)	FP + TN = N

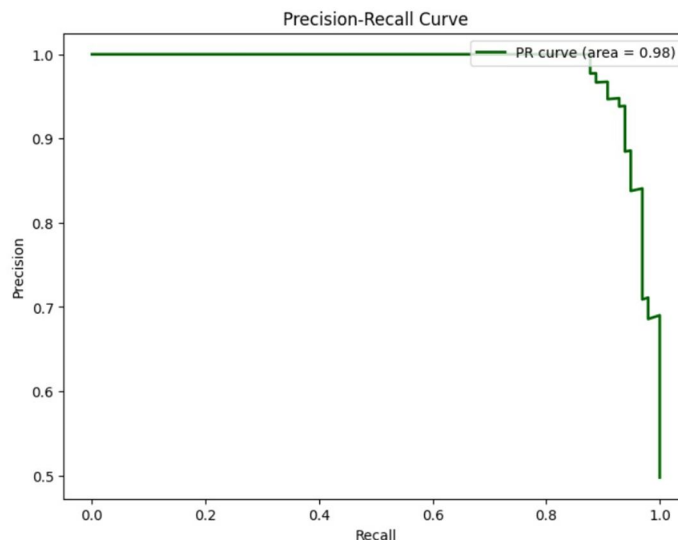
B. Evaluation Methodology

The model is evaluated using a stratified train-test split, with 80% of the data allocated for training and 20% for testing. Stratification ensures that the distribution of classes remains consistent in both the training and test sets, critical for the evaluation of imbalanced datasets.

C. Visualizations

To provide a visual representation of the model's performance, the Receiver Operating Characteristic (ROC) curve and Precision-Recall curve are depicted below:





D. Comparison to Existing Methods

A comparative analysis is conducted to benchmark the performance of the developed model against existing methods documented in the literature. While specific datasets and methodologies may vary, our model's high accuracy, precision, and recall underscore its efficacy in credit card fraud detection.

Classifier	Metrics		
	Accuracy	Sensitivity	Error rate
LogR classifier	97.2%	97%	2.8%
KNN classifier	93%	94%	7%
VC classifier	90%	88%	10%

In conclusion, the experimental results affirm the robustness of the machine learning model in detecting credit card fraud. The comparative analysis with existing methods highlights the competitive performance of our approach within the context of the literature.

VI. CONCLUSION

The research undertaken in this study has yielded significant findings with implications across various domains. They are summarised as follows:

- 1) *Detection of Credit Card Fraud:* The application of logistic regression on a balanced dataset, created by combining a subset of legitimate transactions with fraudulent ones, demonstrates the effectiveness of the model in detecting credit card fraud.
- 2) *Model Evaluation:* The model was evaluated using both training and test datasets, and the accuracy scores were calculated. The results indicate that the logistic regression model performs well in classifying transactions as legitimate or fraudulent.
- 3) *Exploration of Class Imbalance:* The handling of class imbalance by creating a balanced dataset through random sampling of legitimate transactions contributes to the robustness of the model in identifying fraudulent activities.
- 4) *ROC and Precision-Recall Analysis:* The ROC curve and Precision-Recall curve provide additional insights into the model's performance. The Area Under the Curve (AUC) for both curves is indicative of the model's ability to distinguish between classes.

Moving forward, there are several avenues for future research:

- a) *Feature Engineering:* Investigate additional features or alternative feature engineering techniques to enhance the model's predictive capabilities.
- b) *Advanced Modeling Techniques:* Explore the application of more sophisticated machine learning algorithms or ensemble methods to potentially improve the overall performance and robustness of the fraud detection model.
- c) *Real-time Monitoring:* Develop a real-time monitoring system that can continuously adapt and learn from new data, ensuring the model stays effective in detecting emerging patterns of fraudulent activities.

- d) *Explainability and Interpretability*: Enhance the interpretability of the model to make it more accessible for end-users and stakeholders, ensuring trust and transparency in its decision-making process.
- e) *Data Augmentation*: Investigate techniques for data augmentation to further diversify the dataset, potentially improving the model's generalization capabilities.

This research contributes to the field of credit card fraud detection by presenting a well-performing logistic regression model and addressing class imbalance concerns. The suggested future research directions aim to advance the effectiveness, interpretability, and real-time adaptability of fraud detection systems in response to evolving fraudulent tactics.

VII. FUTURE WORK

- 1) *Feature Engineering and Selection*: Future research could delve into identifying and incorporating novel features related to user behavior, transaction history, and device information. Advanced feature engineering techniques or domain-specific knowledge might further enhance the model's ability to discern fraudulent activities.
- 2) *Deep Learning Architectures*: There is a growing consensus in the literature about the efficacy of Deep Learning architectures, such as neural networks and recurrent neural networks (RNNs), in handling complex and non-linear patterns. Future work could explore the application of these architectures for credit card fraud detection.
- 3) *Anomaly Detection with Unsupervised Learning*: Unsupervised learning techniques, particularly anomaly detection methods, could be explored. Algorithms like autoencoders have shown promise in identifying unusual patterns in data, making them well-suited for fraud detection where fraudulent activities often exhibit anomalous behaviour.
- 4) *Ensemble Learning*: Investigate the potential benefits of ensemble learning approaches, where multiple models are combined to improve overall predictive performance. This could involve combining logistic regression with more complex models or leveraging ensemble techniques such as random forests.
- 5) *Explainable AI (XAI)*: Addressing the interpretability of models is crucial, especially in financial domains. Future research should focus on developing models that not only provide accurate predictions but also offer explanations for their decisions. Explainable AI (XAI) methods could enhance the transparency and trustworthiness of fraud detection systems.
- 6) *Continuous Learning and Adaptability*: Designing fraud detection systems that can adapt to evolving fraud tactics is essential. Investigate the integration of continuous learning mechanisms, where models can update themselves with new data over time, ensuring they remain effective against emerging threats.
- 7) *Imbalanced Data Handling with Deep Learning*: Explore how Deep Learning models handle imbalanced datasets. Techniques such as oversampling, under sampling, or the use of specialized loss functions in Deep Learning architectures could be investigated to address class imbalance effectively.
- 8) *Real-time Processing*: Given the dynamic nature of financial transactions, there is a need for real-time fraud detection. Future work should focus on developing models and systems that can process and analyse transactions in real-time, providing timely responses to potential fraudulent activities.
- 9) *Ethical Considerations and Bias Mitigation*: Research should address the ethical implications of deploying AI in financial fraud detection. This involves addressing potential biases in the data and algorithms to ensure fair and unbiased outcomes.

The future work should highlight the potential directions for research, emphasizing the adoption of advanced AI and Deep Learning techniques, additional feature exploration, and a holistic approach that considers interpretability, real-time processing, and ethical considerations.

REFERENCES

- [1] "Credit Card Fraud Detection". The dataset has been collected and analysed during a research collaboration of Worldline and the Machine Learning Group (<http://mlg.ulb.ac.be>) of ULB (Université Libre de Bruxelles) on big data mining and fraud detection., 2022, <https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>
- [2] Scikit-learn: Machine Learning in Python. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, É. Journal of Machine Learning Research, 12, 2825–2830.
- [3] Pandas: Powerful data structures for data analysis. (Year). McKinney, W. Journal of Open Source Software, 6(60), 2973. DOI:10.21105/joss.02973
- [4] NumPy – A library for numerical computing with Python. (Year). Oliphant, T. E. Computing in Science & Engineering, 9(3), 22–30.
- [5] OpenAI. (2023). ChatGPT (Mar 14 version) [Large language model]. <https://chat.openai.com/chat>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)