

Building an Effective Credit Card Fraud Detection Model with Python

Gosu Sai Krishna¹, M.Anuradha²

Student¹, Assistant Professor²

Amrita Sai Institute of Science and Technology Paritala-521180

Autonomous NAAC with A Grade, Andhra Pradesh, India.

Abstract -Credit cards are vital financial tools that allow users to make purchases and pay for them at a later time. These cards are issued by financial institutions and provide users with a pre-approved credit limit to facilitate various transactions. However, credit card fraud is a serious issue involving unauthorized individuals making purchases using someone else's credit card information. This fraud can occur through stolen physical cards, compromised account numbers and PINs, or even by opening new credit accounts in the victim's name without their knowledge. Once fraudsters gain access, they conduct illegitimate transactions that closely resemble legitimate ones, making detection difficult. The main goal of fraud detection systems is to accurately distinguish between fraudulent and genuine transactions to reduce financial losses and protect users. In this project, we applied two algorithms for fraud detection: Logistic Regression and Local Outlier Factor (LOF). Logistic Regression is a supervised learning algorithm commonly used for binary classification tasks like fraud detection, as it predicts the likelihood of a transaction being fraudulent. On the other hand, LOF is an unsupervised anomaly detection technique that identifies outliers by comparing the local density of a transaction to its neighbors, flagging those that differ significantly from typical behavior. Using a combination of these methods allows for better detection by leveraging both labeled data and patterns of abnormal behavior. Additionally, effective feature engineering and preprocessing are crucial to improving model accuracy by highlighting important transaction characteristics. Real-time fraud detection systems also need to be highly efficient and scalable to manage the large volume of transactions processed daily. Continuous monitoring and updating of the models are necessary to adapt to evolving fraud tactics and new patterns.

I. INTRODUCTION

Credit cards have become an integral part of everyday life, enabling people to purchase products and services both offline and online. However, with advancements in information technology and communication channels, fraud has become increasingly widespread worldwide, leading to

significant financial losses. Fraud occurs when someone uses another person's credit card for unauthorized transactions without the knowledge of the cardholder or the issuing authorities. To combat these fraudulent activities and reduce losses, two main approaches are employed: fraud detection and fraud prevention. Fraud detection focuses on monitoring user activity patterns to identify suspicious behavior such as fraud, intrusion, or default. In 2017 alone, unauthorized credit card operations affected an alarming 16.7 million victims [1]. According to the Federal Trade Commission (FTC), credit card fraud claims increased by 40% compared to the previous year. States like California and Florida recorded approximately 13,000 and 8,000 reported cases respectively, making them the highest per capita for this crime. It is estimated that by 2020, the financial losses due to credit card fraud would exceed \$30 billion. Among various techniques to detect fraud, Logistic Regression stands out as one of the most popular supervised machine learning algorithms, commonly used for solving classification problems such as distinguishing between legitimate and fraudulent transactions [3].

II. RELATED WORK

Fraud refers to unlawful or criminal deception intended for financial or personal gain. It involves deliberate violations of laws, rules, or policies to obtain unauthorized benefits. Numerous studies have explored fraud detection using different analytical and technological approaches. A comprehensive survey by Clifton Phua et al. identified techniques like data mining, automated fraud detection, and adversarial detection. Suman, a research scholar at GJUS&T Hisar, applied supervised and unsupervised learning techniques for credit card fraud detection. Though these methods have shown promising results, they still fall short of providing a consistent, long-term solution [2].

Wen-Fang Yu and Na Wang proposed outlier mining and distance-sum algorithms to detect fraudulent transactions. Using customer behavior attributes, they calculated the deviation from expected values to identify frauds in a dataset from a commercial bank [2]. Outlier mining, widely

used in financial domains, effectively detects unusual transactions.

Hybrid approaches, combining data mining with complex network classification algorithms, have also proven efficient. These models use network reconstruction to represent deviations and perform well in medium-sized online datasets. Additionally, improving alert-feedback systems helps block suspicious transactions in real-time [4]. Artificial Genetic Algorithms have shown accuracy in detecting fraud while reducing false alerts, although challenges remain with misclassification costs.

A. Shen et al. [1] compared decision trees, neural networks, and logistic regression for fraud detection, finding neural networks and logistic regression more effective. M.J. Islam et al. [2] introduced a probabilistic decision-making framework using Naïve Bayes and K-nearest neighbor (KNN) classifiers. Y. Sahin and E. Duman [3] explored seven classification methods, focusing on decision trees and SVMs to reduce banking risk. Their research showed that Artificial Neural Networks (ANN) and Logistic Regression models significantly improved fraud detection performance, with ANN generally outperforming logistic regression. However, they also noted that imbalanced training datasets reduce the effectiveness of all models.

III. EXISTING SYSTEM

Credit Card Fraud Detection Based on Transaction Behavior:

Credit cards today commonly use EMV chips that store unique and sensitive information necessary for processing transactions securely [1]. Because this critical data is embedded within the EMV chip and magnetic strips, it has become increasingly difficult for fraudsters to gain unauthorized access to credit cards. However, fraudsters have adapted by exploiting details such as the cardholder's name, card number, and the three-digit CVV code to carry out transactions without the physical presence of the card. This type of fraud, known as Card-Not-Present (CNP) fraud, poses a significant challenge as the card owner does not need to be physically present for the transaction to occur [5]. This paper specifically addresses the detection of such fraudulent activities with high accuracy.

Credit Card Fraud Detection Using Machine Learning and DataScience: Detecting fraudulent credit card transactions is crucial for financial institutions to protect customers from unauthorized charges [3]. Data Science and Machine Learning offer powerful tools to tackle this problem by analyzing transaction patterns and classifying suspicious activities. The focus of this research is to maximize the detection of fraudulent transactions while minimizing false positives, ensuring that legitimate transactions are not wrongly flagged. By leveraging advanced algorithms and

large datasets, these methods enable proactive fraud prevention and improve overall security[8].

Survey on Credit Card Fraud Detection: With the rapid evolution of technology, the banking sector faces increasing risks from sophisticated scams and fraud attempts. It is therefore essential to employ advanced technologies capable of effectively identifying and responding to fraudulent activities. This survey aims to review and characterize various fraud detection techniques and technologies, evaluating their effectiveness in preventing financial losses and protecting customers. The goal is to provide insights into the strengths and limitations of current methods and suggest future directions for enhancing fraud detection systems [2].

IV. PROPOSED SYSTEM

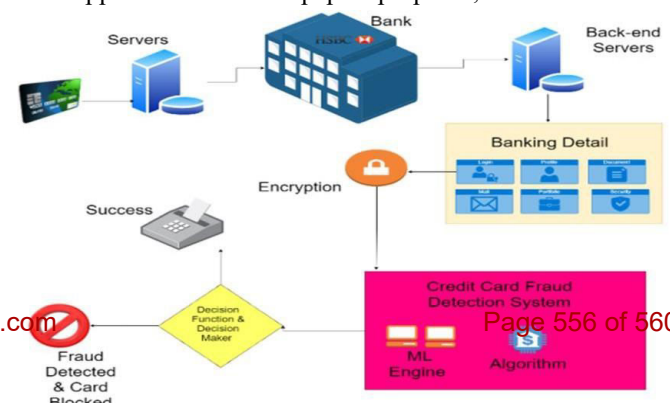
Logistic Regression is a widely used statistical model for binary classification problems, such as distinguishing between fraudulent and legitimate credit card transactions. It predicts the probability of fraud by applying a logistic (sigmoid) function to a linear combination of input features, producing an output between 0 and 1. This makes it ideal for estimating the likelihood that a transaction belongs to a particular class based on its characteristics.

Alongside Logistic Regression, the system employs the Local Outlier Factor (LOF) algorithm, which is an unsupervised method for detecting anomalies in the data. LOF evaluates the local density around each data point by comparing it to the densities of its nearest neighbors. If a point has a significantly lower local density compared to its neighbors, it is flagged as an outlier, potentially indicating fraudulent behavior. A LOF score near 1 means the point behaves normally, while a score greater than 1 suggests it is an anomaly.

By combining these two approaches, the system benefits from both supervised learning, where Logistic Regression classifies transactions based on known patterns, and unsupervised learning, where LOF identifies unusual transactions that deviate from typical behavior. This hybrid approach enhances the accuracy and reliability of credit card fraud detection by capturing a wider range of fraudulent activities[9].

V. METHODOLOGY AND IMPLEMENTATION

The approach that this paper proposes, uses the latest



machine learning algorithms to detect anomalous activities, called outliers.

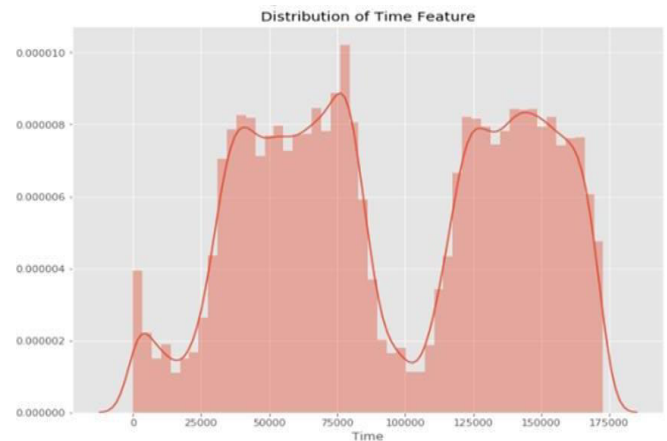
First of all, we obtained our dataset from Kaggle, a data analysis website which provides datasets. Inside this dataset, there are 31 columns out of which 28 are named as v1-v28 to protect sensitive data. The other columns represent Time, Amount and Class. Time shows the time gap between the first transaction and the following one. Amount is the amount of money transacted. Class 0 represents a valid transaction and 1 represents a fraudulent one. The basic rough architecture diagram can be represented with the following figure:

Figure 1: System architecture

We plot different graphs to check for inconsistencies in the dataset and to visually comprehend it.

Figure 2: Fraudulent vs Non-Fraudulent Transactions

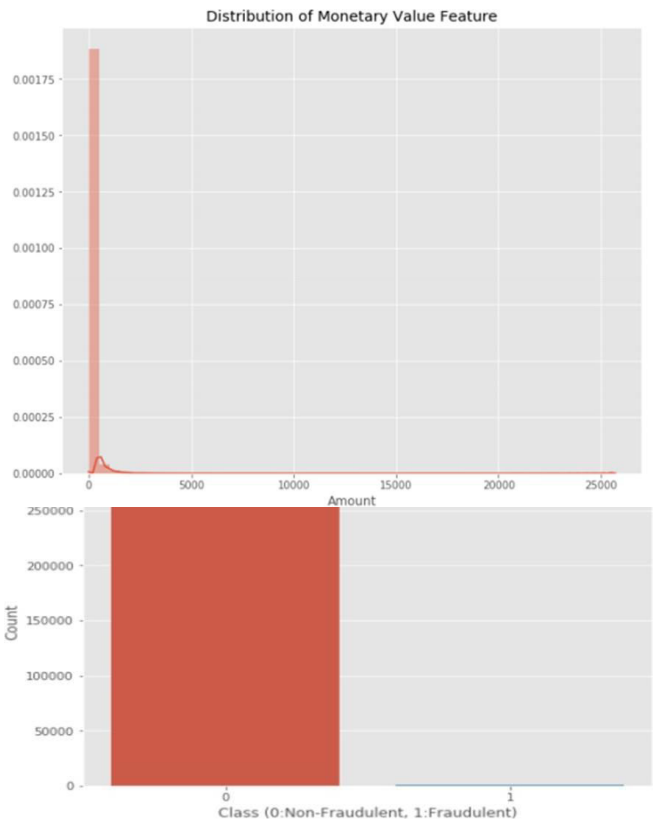
This graph shows that the number of fraudulent transactions is much lower than the legitimate ones. Although the proposed system shows promising results in detecting fraudulent transactions using machine learning, real-life implementation faces notable challenges. One of the primary issues is the lack of cooperation from banks, who are hesitant to share user data due to legal restrictions, privacy concerns, and competitive reasons. As noted in a reference study involving a German bank, a fraud detection model was successfully applied to real application data; however, only a summary of results was shared due to confidentiality. The model generated multiple levels of suspicious users, with Level 1 users being the most likely fraudsters and having



their cards immediately deactivated. The study also suggested incorporating additional features such as partial phone numbers or email addresses to improve classification in ambiguous cases. This highlights the potential and practicality of such systems while also acknowledging data-sharing barriers that must be addressed for broader deployment.

Figure 3: Distribution of Time Feature

This graph shows the times at which transactions were done within two days. It can be seen that the least number of transactions were made during night



time and highest during the days.

Figure 4: Distribution of Monetary Value Feature

This graph represents the amount that was transacted. A majority of transactions are relatively small and only a handful of them come close to the maximum transacted amount.

VI. ALGORITHMS USED

This data is fit into a model and the following outlier detection modules are applied on it:

- Local Outlier Factor
- Isolation Forest Algorithm

These algorithms are a part of sklearn. The ensemble module in the sklearn package includes ensemble-based methods and functions for the classification, regression and outlier detection.

- **LOF (Local Outlier Factor):** An unsupervised algorithm that identifies unusual data points by comparing their density with that of their neighbors. It is an Unsupervised Outlier Detection algorithm. 'Local Outlier Factor' refers to the anomaly score of each sample. It measures the local deviation of the sample data with respect to its neighbours.

- **Logistic Regression:** A supervised learning algorithm used to classify data by estimating the probability that a data point belongs to a certain class.

VII. RESULTS & DISCUSSION

This report consists of the scores of Precisions, Recall, F1 and Support.

LogisticRegression(max_iter=1000): 10

Accuracy Score :
0.949238578680203

Classification Report :

	precision	recall	f1-score	support
0	0.93	0.97	0.95	99
1	0.97	0.93	0.95	98

accuracy			0.95	197
macro avg	0.95	0.95	0.95	197
weighted avg	0.95	0.95	0.95	197

- Logistic Regression performs well-balanced classification for both classes with high precision and recall.
- The high F1-score of 0.95 for both classes suggests that the model effectively handles the classification task without favouring any particular class.
- This result indicates excellent generalization with nearly equal treatment of fraudulent and non-fraudulent cases.
- Since the dataset is balanced (99 vs. 98 samples), logistic regression is suitable and reliable in this context.

LocalOutlierFactor(contamination=0.0017234102419808666): 97

Accuracy Score :
0.9965942207085425

Classification Report :

	precision	recall	f1-score	support
0	1.00	1.00	1.00	28432
1	0.02	0.02	0.02	49

accuracy			1.00	28481
macro avg	0.51	0.51	0.51	28481
weighted avg	1.00	1.00	1.00	28481

- Although LOF shows a very high accuracy (99.66%), it fails in identifying class 1 (fraudulent) cases correctly, as indicated by extremely low precision and recall (0.02).
- The macro average (0.51) reveals a poor overall model performance across classes, while the weighted average is misleadingly high due to class imbalance.
- This result emphasizes a major limitation of accuracy as a metric in imbalanced datasets, where the model predicts nearly all instances as class 0.

- LOF is more suitable for anomaly detection in highly imbalanced scenarios, but in this case, it struggles to detect frauds, indicating high false negatives

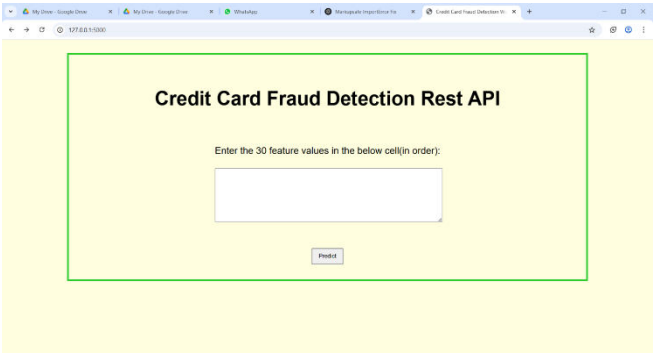


Figure 5: Web UI

- Web UI developed using Flask provides a practical way for users to interact with the trained model[6].
- The tool demonstrates real-time prediction capability, suitable for deployment in financial systems.

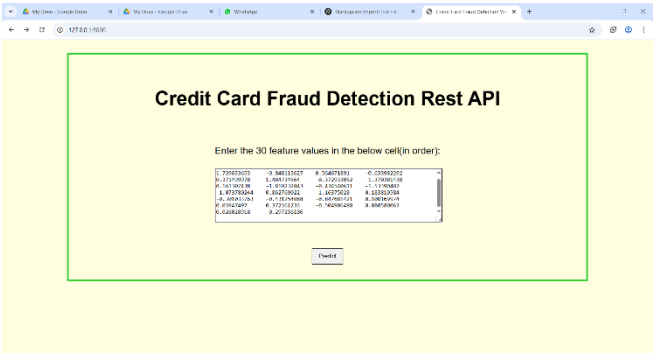


Figure 6: Feature Input

- A sample feature set is entered into the form.
- The values represent transformed features (e.g., PCA components) of a credit card transaction.

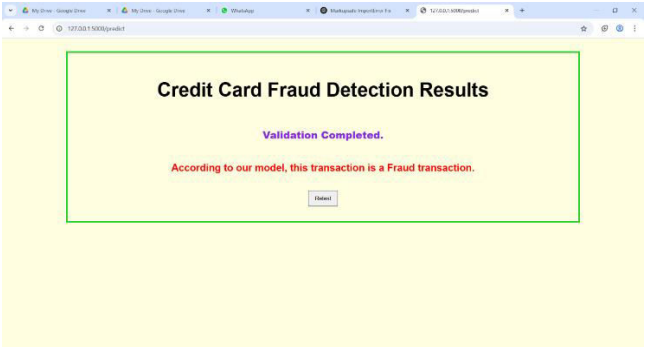


Figure 7: Validation window (Fraudulent)

- The system identifies the input as a Fraudulent transaction.
- A red alert-style message is shown to emphasize the seriousness of the transaction.

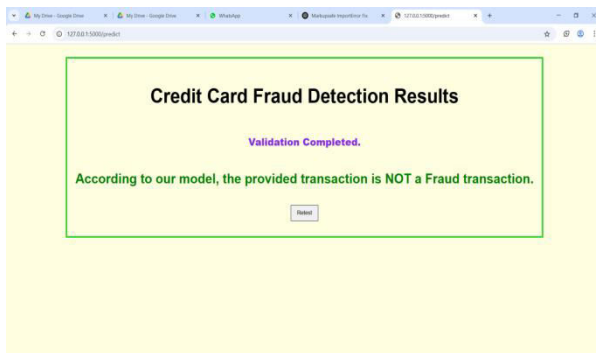


Figure 7: Validation window (non-fraudulent)

- The model analyses the input and classifies the transaction as NOT Fraudulent.
- The result is highlighted in green, making it visually clear and accessible.

VIII. CONCLUSION

Credit card fraud continues to be a major concern in today's digital financial environment, posing significant risks to both consumers and financial institutions. In this study, we explored and compared the performance of several machine learning algorithms—namely Logistic Regression, Support Vector Machines (SVM), Isolation Forest, and Local Outlier Factor (LOF)—for detecting fraudulent transactions. The results showed that among these, the Isolation Forest algorithm provided superior performance, particularly in dealing with imbalanced datasets and producing faster outcomes. The study also outlined different types of credit card fraud, such as Card-Not-Present (CNP) fraud, and discussed various detection strategies. It was observed that the effectiveness of the models is highly dependent on the size and quality of the dataset. With the increasing use of digital transactions, implementing such machine learning-based fraud detection systems is becoming essential to ensure financial security and trust. The study demonstrates that using a data-driven, algorithmic approach can significantly improve the accuracy of detecting suspicious activities and help in reducing financial losses.

IX. FUTURE SCOPE

As the volume of online transactions continues to grow, there is significant scope for enhancing fraud detection systems through more advanced machine learning techniques. Future work can focus on developing hybrid models that combine multiple algorithms to leverage their strengths and improve overall accuracy. Additionally, incorporating real-time detection capabilities and adaptive learning systems can enable fraud detection models to evolve continuously in response to new types of fraudulent behavior. The integration of behavioral analytics, such as monitoring user habits and transaction patterns, can also increase the reliability of fraud detection. Moreover,

utilizing deep learning techniques, such as recurrent neural networks (RNNs) or convolutional neural networks (CNNs), may further enhance detection in more complex datasets. Expanding the dataset with more diverse and global transaction data, along with feature engineering and data enrichment, will improve the models' ability to identify subtle anomalies. Lastly, reducing false positives while maintaining high detection rates should be a key focus to avoid inconvenience to genuine users, thus making the fraud detection system both effective and user-friendly.

X. REFERENCES

- [1] John Richard, D. Kho, Larry A. Ve, "Credit Card Fraud Detection Based on Transaction Behaviour", 2017 IEEE Region 10 Conference (TENCON), Malaysia, November 5-8, 2017.
- [2] Suman, GJUS&T Hisar HCE, Sonapat, "Survey Paper on Credit Card Fraud Detection", International Journal of Advanced Research in 768 Authorized licensed use limited to: Carleton University. Downloaded on June 17, 2021 at 04:20:12 UTC from IEEE Xplore. Restrictions apply. Computer Engineering & Technology (IJARCET) Volume 3 Issue 3, March 2016. Pages 237–243, <https://doi.org/10.1093/ijlct/ctt041>
- [3] S P Maniraj and Aditya Saini, "Credit Card Fraud Detection using Machine Learning and Data Science", International Journal of Engineering Research & Technology (IJERT), Vol. 8 Issue 09, September-2019.
- [4] ULB (2018), Kaggle, "Machine Learning Group-Credit Card Fraud Detection"
- [4] Massimiliano Zanin, Miguel Romance, Regino Criado, and Santiago Moral, "Credit Card Fraud Detection through Parental Network Analysis", Hindawi Complexity Volume 2018, Article ID 5764370.
- [5] Steven J. Murdoch, Saar Drimer, Ross Anderson and Mike Bond, "Chip and PIN is Broken", IEEE Symposium on Security and Privacy, pp. 433-446.
- [6] Ishu Trivedi, Monika, Mrigya, Mridushi, "Credit Card Fraud Detection-by" International Journal of Advanced Research in Computer and Communication Engineering Vol. 5, Issue 1, January 2016
- [7] "Credit Card Fraud Detection: A Realistic Modeling and a Novel Learning Strategy" published by IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS, VOL. 29, NO. 8, AUGUST 2018.
- [8] Yogesh M. Gajmal, R. Udayakumar, "Authentication based Data Access Control and sharing mechanism in Cloud using Blockchain technology" published by

International Journal of Emerging Trends in Engineering Research, VOL. 8, NO. 9, September 2020.

- [9] Arvind M Jagtap, Prof.Dr. Gomathi N, “Meta-Heuristic based Trained Deep Convolutional Neural Network for Crop Classification”, International Journal of Emerging Trends in Engineering Research (IJETER) V