# FRAUD SHIELD | SECURING CREDIT CARD TRANSACTIONS

White Paper

APRIL 6, 2024

BELLEVUE UNIVERSITY

Madhavi Ghanta DSC 680

## Business Problem

Businesses and financial institutions face a significant challenge as credit card fraud becomes more common. Effective solutions are required if we are to immediately identify and stop fraudulent transactions. With the help of machine learning, this research hopes to create a reliable fraud prediction system that can identify authentic credit card transactions from fraudulent ones. The objectives include raising security, cutting down on false positives, boosting customer satisfaction, and keeping up with the rapidly changing landscape of fraud techniques. Financial stability, consumer confidence, and the general health of the sector depend on this endeavour.

## Background /History

The detection of credit card fraud has developed in tandem with the rise in credit card usage. It switched from using manual processes to electronic authorization using magnetic stripes. Rule-based systems came into being as scammers changed to take advantage of weaknesses. Systems increasingly use behavioural analytics to detect fraud, thanks to recent breakthroughs in machine learning and real-time analytics. To ensure that fraud detection techniques continue to evolve, cooperation and data sharing are essential in the fight against changing fraud strategies.

## Datasets

The public-domain datasets utilized in this experiment were obtained from Kaggle (Shenoy, K. (2020, August 5)). From January 2019 to December 2020, this dataset includes both authentic and fraudulent credit card transactions together with stimulated credit card data.
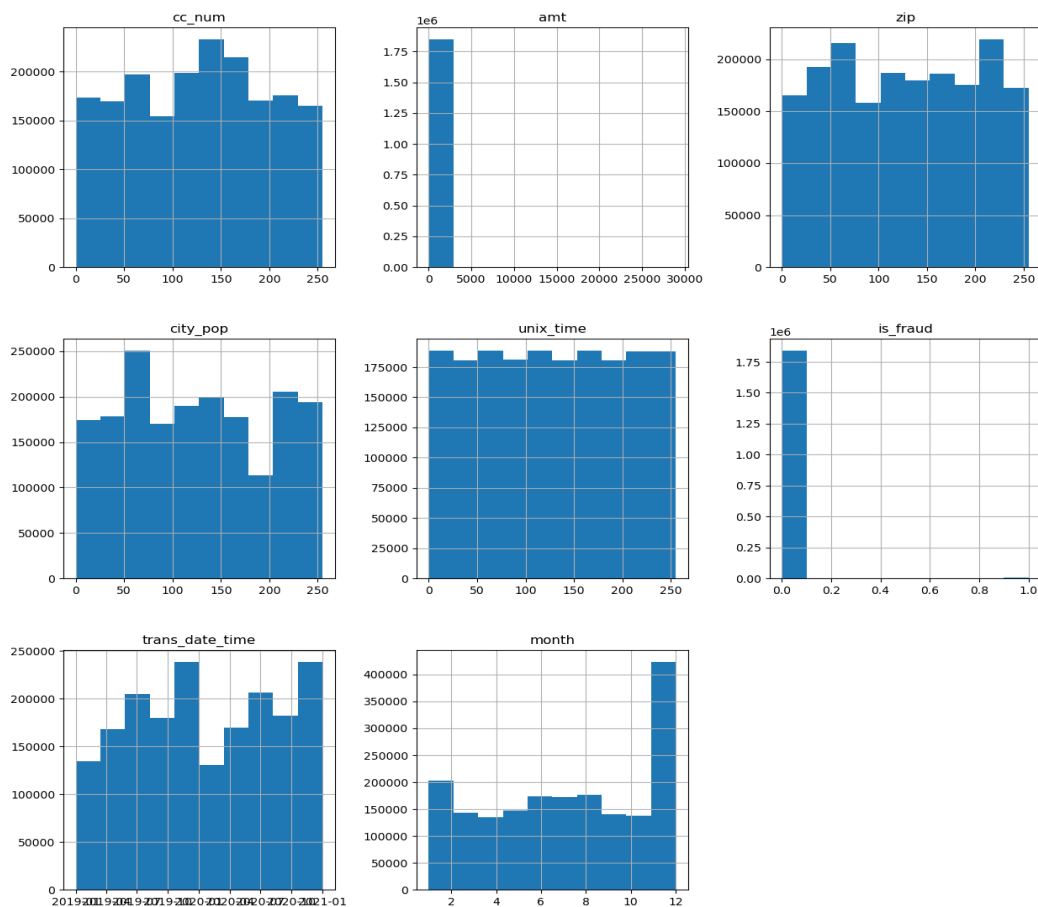
## Data Preparation

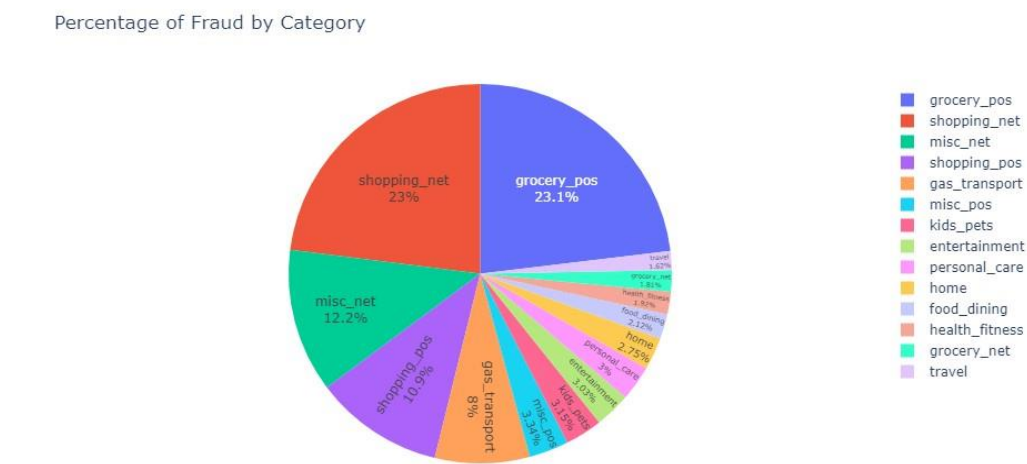The following steps were performed to prepare the data for modelling.

i.        Checked for null rows/columns in the data.

ii.       Performed check for duplicates.

iii.     Converted datetime string to a datetime datatype.

iv.     Added a month column for visualization.

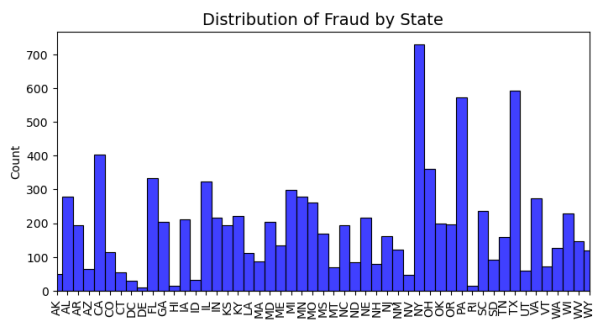v.      Dropped columns ('Unnamed: 0' and 'trans_date_trans_time')

## Visualizations

## Numeric Variables Distribution:
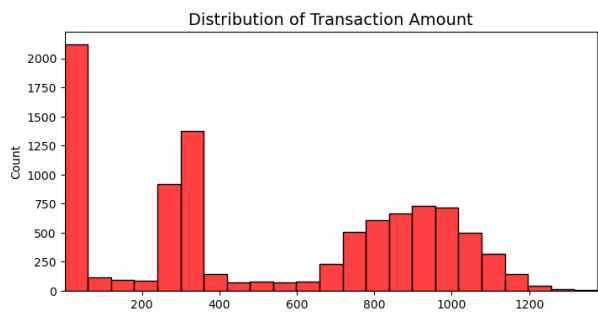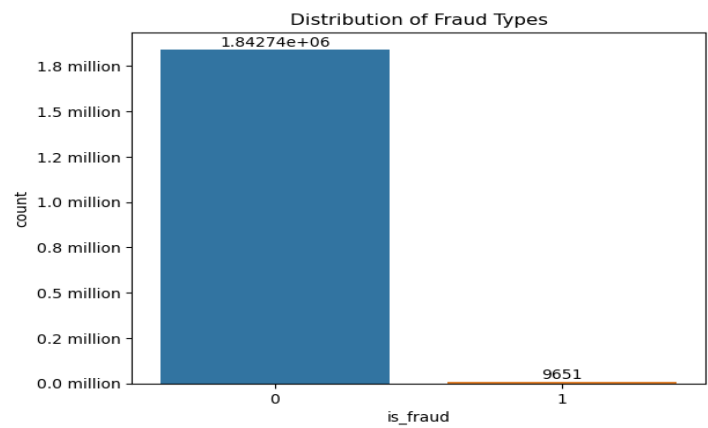
# Percentage of Fraud by Category

Percentage of Fraud by Category



# Distribution of Amount & State



Distribution of Transaction Amount

Distribution of Fraud by State

# Distribution of Fraud Types



Distribution of Fraud Types

To construct an effective model, it was essential to address the dataset's imbalance. To achieve this, the Synthetic Minority Oversampling Technique - SMOTE (Appendix (ii)) was employed.

**Methods**

With the balanced dataset, the subsequent step was to partition the data into training and testing datasets.

The following models were developed, with their respective outcomes recorded.

i.     Logistic Regression: Interpretability, computational efficiency, and adaptability for smaller datasets are advantages of using a Logistic Regression model for credit card fraud detection. It delivers clear probability estimates, serves as a baseline, and sheds light on the effects of features. Performance on intricate, non-linear data patterns, however, might be constrained by its linear nature.

ii.    Random Forest: There are several advantages to using a Random Forest model for credit card fraud detection, including the ability to handle imbalanced data, handle ensemble learning, be robust against noise and outliers, and capture intricate patterns. It is a good option due to its versatility, ease of tuning, and feature importance analysis.

iii.   Gradient Boosting: Because of its high prediction accuracy, ensemble nature, which lessens overfitting, and capacity to handle complicated patterns in the data, Gradient Boosting is an effective technique for detecting credit card fraud. Insights into feature importance are also provided, which helps find pertinent variables for fraud detection.

When it comes to managing credit card fraud detection jobs, Random Forest and Gradient Boosting are two excellent ensemble algorithms. While Gradient Boosting successively reduces the mistakes of earlier trees, Random Forest separately constructs a variety of trees.

iv. Neural Networks: Because they can identify complex non-linear patterns in data, neural network (NN) models are advantageous in the detection of credit card fraud. Their proficiency in feature learning from unprocessed data enables them to identify intricate and dynamic fraud behaviors.

### Analysis

Since the dataset consists of categorical features, it is essential to represent the categorical data in a numerical format. For this, an encoding technique (Label Encoder) was implemented. Additionally, the StandardScaler preprocessing technique was used to standardize or normalize numerical features in the dataset.

The models were then built, and the outcomes were recorded as follows.

### Random Forest Classifier

| Fraud Type | Precision | Recall | F1-Score | | Accuracy | ROC-AUC Score |
|---|---|---|---|---|---|---|
| 0 – non-Fraud | 1.00 | 1.00 | 1.00 | | 99.75% | 0.86 |
| 1 - Fraud | 0.79 | 0.72 | 0.75 | | | |

## Logistic Regression

| Fraud Type | Precision | Recall | F1-Score |
|---|---|---|---|
| 0 – non-Fraud | 1.00 | 0.94 | 0.97 |
| 1 - Fraud | 0.06 | 0.77 | 0.12 |

| Accuracy | ROC-AUC Score |
|---|---|
| 93.99% | 0.85 |

## Gradient Boosting

| Fraud Type | Precision | Recall | F1-Score |
|---|---|---|---|
| 0 – non-Fraud | 1.00 | 1.00 | 1.00 |
| 1 - Fraud | 0.50 | 0.89 | 0.64 |

| Accuracy | ROC-AUC Score |
|---|---|
| 99.4% | 0.94 |

## Neural Network Model

| Fraud Type | Precision | Recall | F1-Score |
|---|---|---|---|
| 0 – non-Fraud | 1.00 | 0.99 | 0.99 |
| 1 - Fraud | 0.26 | 0.75 | 0.38 |

| Accuracy | ROC-AUC Score |
|---|---|
| 98.7% | 0.87 |

## Conclusion:

In summary, this project focused on developing and implementing effective credit card fraud detection models. We explored different machine-learning techniques, including logistic regression, random forests, gradient boosting, and neural networks.

Through data preprocessing, feature engineering, and model evaluation, we aimed to create accurate and efficient fraud detection systems. The models demonstrated promising results, with some outperforming others in specific areas.

| Model | Random Forest Classifier | Logistic Regression | Gradient Booster | Neural Network |
|---|---|---|---|---|
| **Accuracy** | 99.75% | 93.99% | 99.4% | 98.7% |

With an accuracy of 99.75%, the Random Forest Classifier is the best-performing model for this dataset. However, the dataset is highly imbalanced, with a significantly larger number of non-fraudulent transactions. This can impact the model's performance and interpretation of metrics.

While all models demonstrated high accuracy and impressive performance in identifying non-fraudulent transactions, there is a trade-off between precision and recall for detecting fraudulent transactions.

## Assumptions

Simulated data often assumes an imbalanced class distribution, with a small proportion of transactions representing fraud. This project assumes that features used for modelling are independent or have certain dependencies that mimic real-world relationships.

## Limitations

Credit card fraud detection systems have the risk of false positives, where legitimate transactions are mistakenly flagged as fraud, and false negatives, where some fraudulent transactions go undetected. Imbalanced data and model complexity can affect performance, while the threat of adversarial attacks and the need for data privacy are ongoing concerns.

## Challenges

Ensuring security for simulated data and achieving real-world applicability could be areas of concern. Simulated datasets might not fully capture the complexities of actual fraud scenarios, potentially limiting the model's effectiveness in real-world situations.

## Future Uses / Additional Applications

The applications of fraud detection techniques are expanding across various industries and sectors as organizations seek to protect themselves from evolving threats and optimize their operations. For example, fraud detection techniques can be extended to other payment methods, such as mobile wallets, digital currencies (cryptocurrencies), and peer-to-peer payment systems.

## Recommendations

- Incorporating behavioural analysis to detect anomalies in customer transaction behaviour over time.
- Ensuring Fraud detection systems comply with financial regulations and data privacy laws.

## Implementation Plan

The implementation plan for credit card fraud detection involves defining clear objectives, collecting, and preparing transaction data, selecting appropriate machine learning

models, addressing class imbalance, evaluating model performance, and integrating it into a real-time system with continuous monitoring and alerts.

## Ethical Assessment

Several ethical considerations are crucial for this project:

i.  Data Privacy: Despite the simulation, safeguarding privacy is essential. Ensuring simulated data does not resemble real customer information avoids accidental exposure of PII (Appendix (i)).

ii.  Informed Consent: Transparency builds trust and addresses concerns. Understanding if consent was obtained for real-based simulations, even if anonymized, adds ethical integrity.

iii.  Intent and Use: Ethical utilization of simulated data is paramount, prohibiting any malicious or harmful intent.

iv.  Data Security: Robust security measures should be applied to the simulated dataset, treating it with the same importance as real customer data.

v.  Stakeholder Implications: The viewpoints of stakeholders like credit card issuers, customers, and regulators to ensure ethical alignment should be considered.

## References

Shenoy, K. (2020, August 5). *Credit Card Transactions Fraud Detection Dataset*. Kaggle. https://www.kaggle.com/datasets/kartik2112/fraud-detection?select=fraudTrain.csv

## Appendix

i.      PII stands for Personally Identifiable Information. It refers to any data that can be used to identify a specific individual. Protecting PII is crucial to safeguard individuals' privacy and prevent identity theft, fraud, and unauthorized access to sensitive information.

ii.     SMOTE - Synthetic Minority Oversampling Technique (SMOTE), generates synthetic instances for the minority class, thus balancing the dataset.

## Questions

1. Is credit card fraud defined?

2. What is simulated credit card data?

3. Why do we need to simulate the credit card data?

4. What is a Neural Network Model? Why was this selected for the fraud detection project?

5. What is the difference between the Random Forest and Gradient Boost model?

6. Which model performs better between the Random Forest and the Gradient Boost models?

7. Will these models withstand vast datasets?

8. Can we provide real-time fraud detection using these implemented models?

9. Will these implemented models work with data from various credit card firms?

10. Can these models be extended for fraud detection in any other sector?