**Name**-Vishwa Bhadiyadara
**NUID**-002636314

# Leveraging Data Mining Techniques for Enhanced Bank Account Fraud Detection

## Introduction

Fraudulent activities in banking represent a significant threat to both financial institutions and customers. The advent of digital banking has amplified the risk and sophistication of such illicit activities. Traditional rule-based systems are proving inadequate due to high false positive rates and lack of adaptability against evolving fraud tactics. **Data mining** and **supervised machine learning methods** have emerged as potent tools to identify and prevent fraudulent transactions. This study scrutinizes the efficacy of these techniques using the PaySim synthetic dataset, which simulates mobile money transactions and captures the complexities of real-world financial systems. The research papers selected for this study are seminal works that have established foundational models and contributed to the advancement of fraud detection methodologies.

## Data Considerations

The application of data mining techniques in fraud detection significantly hinges on the availability and quality of the dataset used. The **PaySim synthetic dataset**, created to simulate mobile money transactions and fraudulent behavior, serves as an ideal case study for such applications due to the scarcity of publicly available financial transaction datasets.

## Dataset Genesis and Composition

PaySim simulates mobile money transactions based on a sample of real transactions extracted from one month of financial logs from a mobile money service implemented in an African country. The original logs were provided by a multinational company, who is the provider of the mobile financial service which is currently running in more than **14 countries all around the world**.

This synthetic dataset is scaled down 1/4 of the original dataset and it is created just for Kaggle.

https://www.kaggle.com/datasets/ealaxi/paysim1/code

## Past Research

There are 5 similar files that contain the run of 5 different scenarios. These files are better explained in this thesis chapter 7 (PhD Thesis Available here https://bth.diva-portal.org/smash/record.jsf?pid=diva2%3A955852&dswid=-2529

We ran PaySim several times using random seeds for 744 steps, representing each hour of one month of real time, which matches the original logs. Each run took around 45 minutes on an i7 intel

processor with 16GB of RAM. The final result of a run contains approximately 24 million of financial records divided into the 5 types of categories: CASH-IN, CASH-OUT, DEBIT, PAYMENT and TRANSFER.

## Dataset Features and Attributes

Each record in the PaySim dataset contains key features that mirror real-world financial transactions:

- **step**: Represents a unit of time, specifically one hour, summing up to 744 steps to encapsulate a month's data.

- **type**: Categorizes the transaction into one of five types: CASH-IN, CASH-OUT, DEBIT, PAYMENT, and TRANSFER.

- **amount**: The monetary value of the transaction, providing a direct measure for the transaction's size.

- **nameOrig** and **nameDest**: Identify the initiator and the recipient of the transaction, respectively.

- **oldbalanceOrg**, **newbalanceOrig**, **oldbalanceDest**, **newbalanceDest**: Capture the state of account balances before and after the transaction. It's important to note that for transactions marked as fraud, these balance-related features are not reliable due to the cancellation of the fraudulent transactions as part of the simulation's design.

- **isFraud**: A binary flag indicating whether the transaction is fraudulent.

- **isFlaggedFraud**: Flags illegal attempts to transfer amounts over a threshold, set at 200,000 in the dataset.

## Data Quality and Preprocessing

The synthetic nature of the dataset ensures high data quality with no missing values, a common issue in real-world datasets. However, the dataset's synthetic attribute also poses a unique challenge: ensuring the simulated data accurately reflects genuine patterns of fraudulent behavior. The preprocessing of the dataset included encoding categorical variables, such as the transaction type, and normalizing the numerical features to feed into various machine learning models adequately.

## Feature Engineering and Selection

Feature engineering was a crucial step in preparing the dataset for analysis. Derived features, such as the difference between the new and old balances, could provide additional insight into the transaction behavior. However, due to the nature of fraudulent transactions leading to the cancellation, these derived balance features were deemed unreliable and thus excluded from the analysis.

**Name**-Vishwa Bhadiyadara
**NUID**-002636314

## Class Imbalance and Resampling

One of the primary data considerations was the significant class imbalance, a common and challenging aspect of fraud detection. With fraudulent transactions representing a small fraction of all transactions, the dataset was subjected to resampling techniques. SMOTE was employed to artificially balance the dataset, ensuring that the models do not bias towards the majority class and can detect fraudulent transactions effectively.

## Visualization and Exploratory Data Analysis

Data visualization played a pivotal role in understanding the dataset's characteristics. Heatmaps and scatter plots were used to identify correlations and detect anomalies, while distribution plots for transaction amounts and types provided insights into the transactional behavior captured in the dataset.
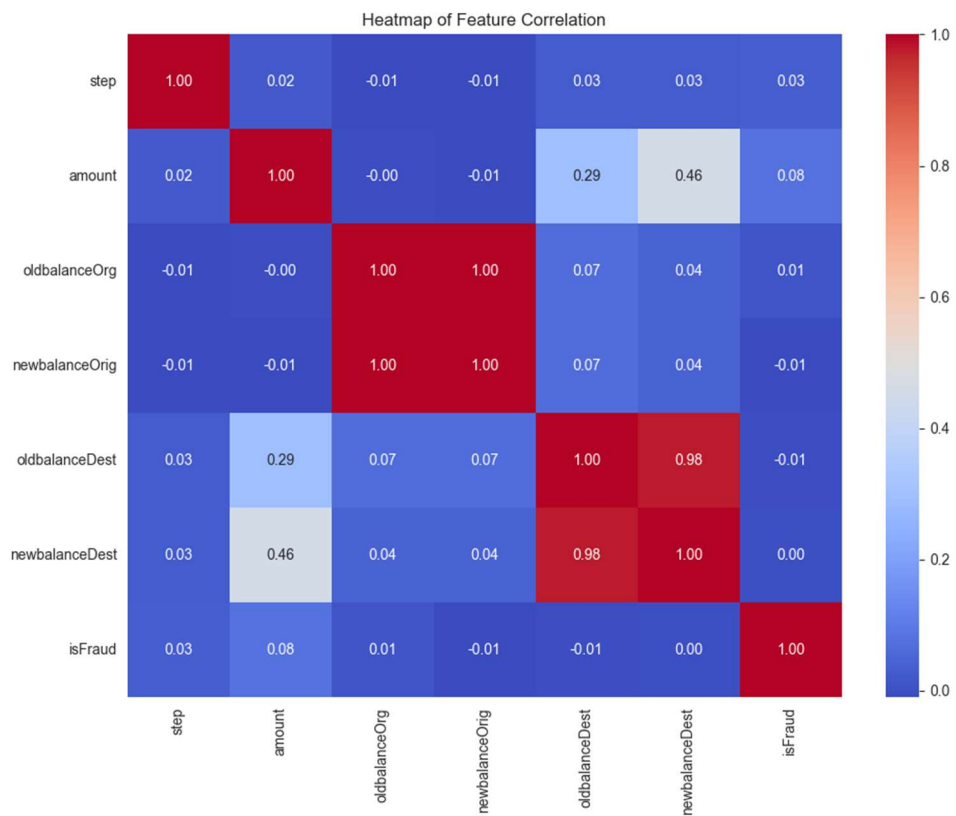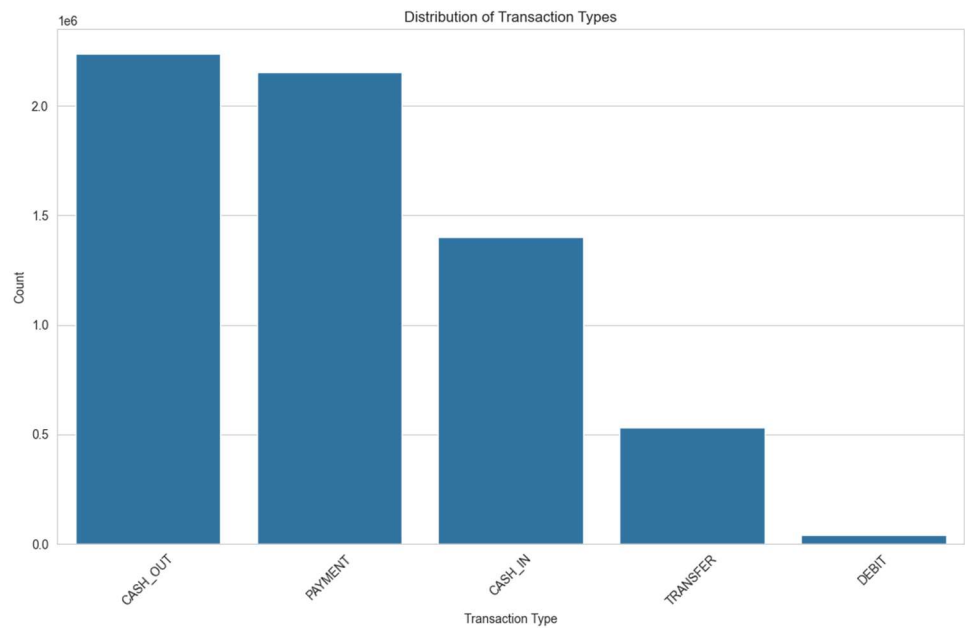
## Acknowledgements

**Name**-Vishwa Bhadiyadara
**NUID**-002636314

Distribution of Transaction Types



Heatmap of Feature Correlation

**Name**-Vishwa Bhadiyadara
**NUID**-002636314

## Methodology Comparison

In addressing the pressing issue of bank account fraud detection, this study investigates the performance of various supervised machine learning algorithms. The PaySim synthetic dataset, serving as a microcosm of real-world financial transactions, provides the basis for this comparative analysis. The algorithms assessed include Logistic Regression, Decision Trees, Random Forest, Support Vector Machines (SVM), and Neural Networks.

**Logistic Regression** was selected for its efficiency and interpretability. Its linear approach to classification makes it a suitable benchmark for binary outcomes, such as detecting fraud. However, its simplicity is a double-edged sword; while it offers speed and straightforward interpretation, it may falter with complex, non-linear relationships within data.

**Decision Trees** allow for non-linear data separation, offering a more nuanced approach than Logistic Regression. They partition the feature space into decision regions, providing an intuitive understanding of the data's structure. Decision Trees are also inherently interpretable, as they offer clear decision rules and variable importance. Nevertheless, their susceptibility to overfitting can lead to poor generalization if not properly pruned or regularized.

**Random Forest**, an ensemble of Decision Trees, addresses the overfitting issue by creating a multitude of trees and aggregating their predictions. This method benefits from the diversity among the individual trees, each trained on a subset of the data, and reduces variance without significantly increasing bias. In the context of fraud detection, Random Forest is advantageous for its ability to model complex interactions and its robustness to outliers and noise.

**Support Vector Machines (SVM)** excel in high-dimensional spaces, like those often encountered in transaction data. By employing kernel functions, SVMs can perform non-linear classification, finding the hyperplane that maximizes the margin between classes. This is particularly beneficial for fraud detection, where the boundary between fraudulent and legitimate transactions may not be linearly separable. However, SVMs require careful tuning of hyperparameters, and the choice of the kernel can greatly influence their performance.

**Neural Networks**, with their deep learning capabilities, offer a powerful tool for capturing intricate patterns through multiple layers of abstraction. They are particularly well-suited for large and complex datasets, potentially capturing complex fraud signals that other algorithms might miss. The depth and breadth of the network architecture, however, introduce challenges in training time, overfitting, and the need for a substantial amount of data to learn effectively.

In the comparative analysis, performance metrics such as accuracy, precision, recall, F1 score, and the area under the receiver operating characteristic curve (AUC-ROC) were used to evaluate each algorithm's ability to correctly identify fraudulent transactions. The models' applicability to the problem was closely scrutinized, considering the class imbalance inherent in fraud detection datasets. Random Forest and Neural Networks emerged as strong contenders due to their high performance across various metrics and their ability to handle the complexity of the data. However, Logistic Regression and SVMs provided valuable insights into the nature of the fraud detection problem, highlighting the importance of feature selection and model interpretability.
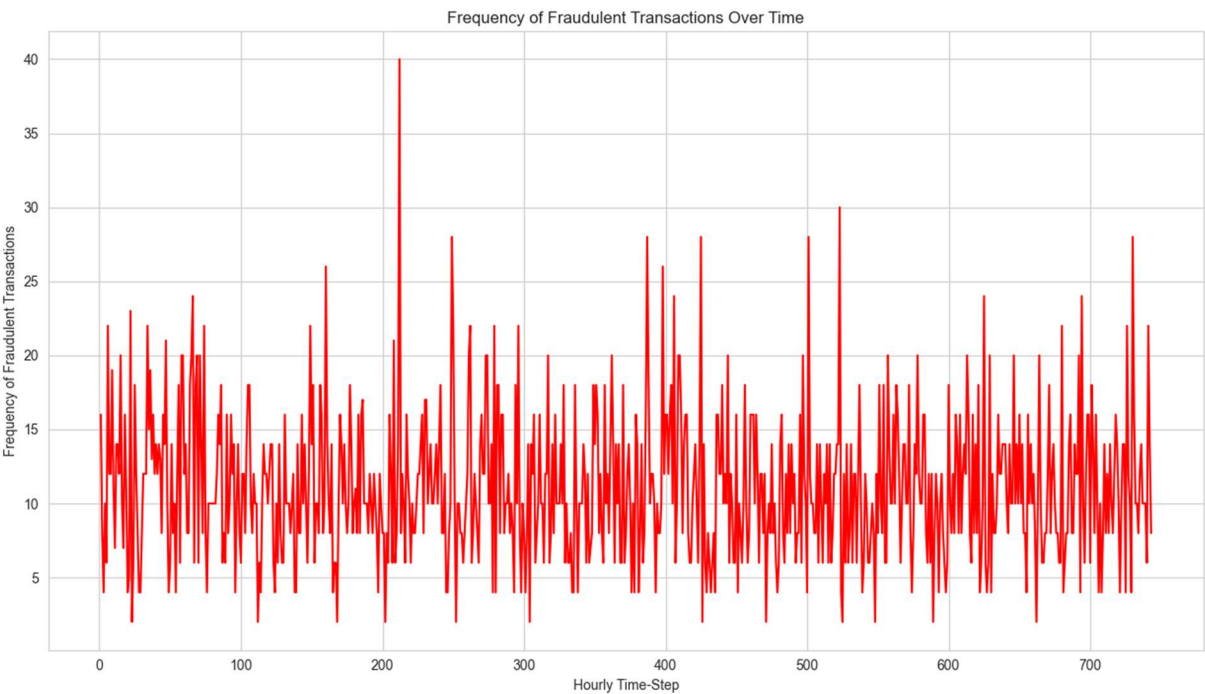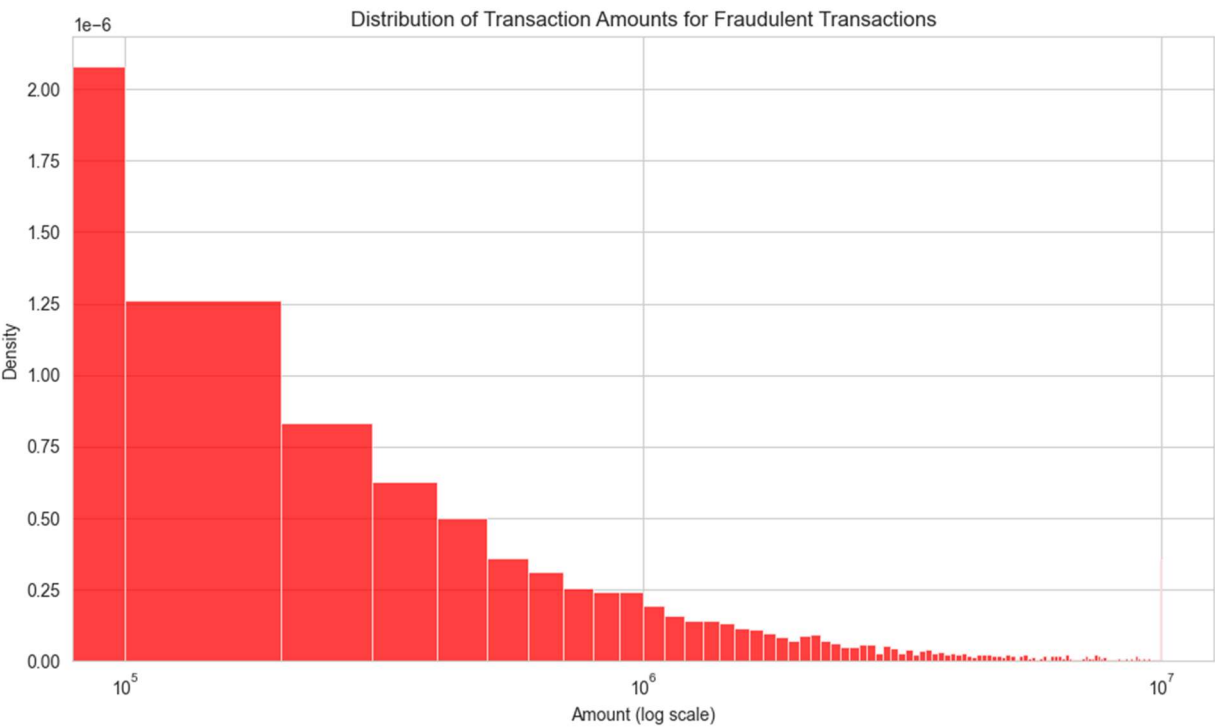
**Name**-Vishwa Bhadiyadara
**NUID**-002636314

The results demonstrated that while no single algorithm uniformly outperformed others across all metrics, ensemble methods like Random Forest offered a compelling balance between detection rate and false positives, which is critical in fraud detection. Neural Networks suggested a promising avenue for future exploration, given their potential for learning highly complex and abstract representations of data.
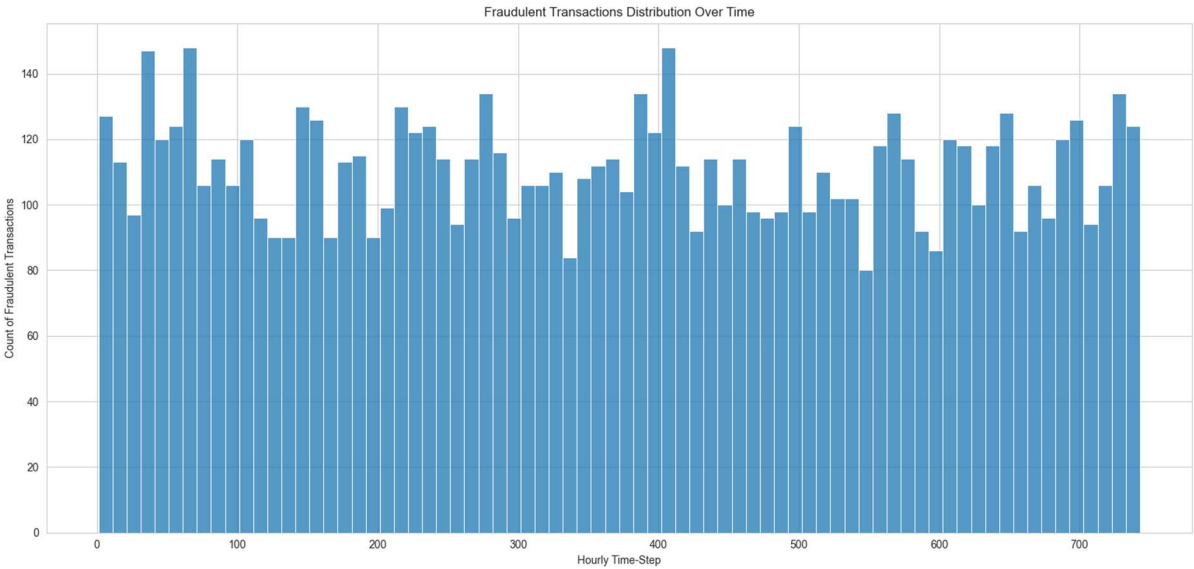
Each algorithm's assumptions—such as the independence of features in Logistic Regression or the data distribution in SVM—were carefully considered, as these can significantly impact the model's performance and generalizability. In particular, the assumption of feature independence in Logistic Regression may not hold true in fraud detection, where transaction attributes can be highly interdependent.

Moreover, the choice of algorithm has implications for the operational deployment of fraud detection systems. Logistic Regression and Decision Trees, with their simplicity and interpretability, may be more readily accepted by stakeholders requiring transparency in decision-making. In contrast, the 'black-box' nature of Neural Networks may pose challenges for regulatory compliance and auditability, despite their potential for higher predictive performance.

The study underscores the need for a nuanced approach to selecting and tuning machine learning algorithms for fraud detection. It also highlights the importance of domain knowledge in feature engineering and the interpretation of model predictions. The choice of algorithm ultimately impacts not only the model's ability to detect fraud but also the operational aspects of deploying the system within a financial institution's ecosystem.

In conclusion, the comparative analysis of supervised learning methodologies revealed a trade-off between complexity and interpretability, efficiency and accuracy, as well as the necessity for continual model refinement in response to the evolving nature of financial fraud. The insights gained from this study provide a foundation for financial institutions to build robust, effective, and interpretable fraud detection systems, contributing to the security and integrity of the banking sector.

Distribution of Transaction Amounts for Fraudulent Transactions


Frequency of Fraudulent Transactions Over Time

**Name**-Vishwa Bhadiyadara
**NUID**-002636314

Fraudulent Transactions Distribution Over Time

## Performance and Evaluation

The performance evaluation of machine learning models in the context of bank account fraud detection requires careful consideration of the metrics that best reflect the effectiveness of each model. Given the highly imbalanced nature of the PaySim synthetic dataset—where legitimate transactions far outnumber fraudulent ones—traditional metrics such as accuracy cannot be solely relied upon. Instead, this study has emphasized precision, recall, F1-score, and the area under the Receiver Operating Characteristic (ROC) curve (AUC-ROC) to provide a comprehensive assessment.

**Precision** is critical in the financial domain, as false positives (legitimate transactions misclassified as fraud) can lead to customer dissatisfaction and increased operational costs. High precision indicates that when a model predicts a transaction as fraudulent, it is likely to be correct. However, an exclusive focus on precision can result in a model that is overly conservative, failing to detect actual fraudulent transactions—a dangerous outcome for any financial institution.

**Recall**, or sensitivity, measures the model's ability to identify all actual instances of fraud. In fraud detection, a high recall rate is desirable to ensure that the majority of fraudulent transactions are captured, even if it means tolerating a higher rate of false positives. For banks, the cost associated with missing a fraudulent transaction (a false negative) can be substantial, potentially far outweighing the inconvenience of investigating false alerts.

The **F1-score** provides a balance between precision and recall, offering a single metric that encapsulates the model's overall performance with respect to both false positives and false negatives. It is particularly useful when the costs of false positives and false negatives are roughly equivalent, or when the classes are imbalanced, as is the case with fraud detection.

The **AUC-ROC** score is another critical metric in evaluating the performance of fraud detection models. It represents the likelihood that the model will rank a randomly chosen positive instance higher than a randomly chosen negative one. A high AUC-ROC score indicates a model with good separability between the positive and negative classes. For the banking industry, this means the model is effective in distinguishing between fraudulent and legitimate transactions across various thresholds, offering flexibility in how the model is deployed operationally.

In this study, models were first evaluated on a holdout test set to ensure an unbiased estimate of their performance. Logistic Regression, with its straightforward decision boundary, provided a baseline AUC-ROC score that was surpassed by more complex models. Decision Trees showed improved recall but at the expense of precision, leading to a large number of false positives. Random Forest improved upon the Decision Tree's performance, offering a higher F1-score and AUC-ROC score due to its ensemble approach, which aggregates predictions from multiple trees to mitigate overfitting and improve generalization.

Support Vector Machines, particularly with non-linear kernels, performed well in terms of AUC-ROC, indicating strong separability. However, the interpretability of SVMs can be limited, especially with non-linear kernels, which can be a drawback in settings where explaining decisions is crucial. Neural Networks, particularly deep learning models, showed potential with high scores across all metrics, suggesting an ability to capture complex patterns and interactions in the data. However,

their 'black box' nature poses challenges in interpretability and requires careful consideration of the trade-offs between predictive power and transparency.

Model evaluation also involved analyzing the confusion matrix for each model, which provides a detailed breakdown of true positives, false positives, true negatives, and false negatives. This analysis is vital in the banking context, where different types of errors carry different costs and operational implications.

To ensure robustness, cross-validation was employed, providing a more reliable estimate of the model's performance by evaluating it across multiple subsets of the data. This approach helps in assessing the stability and reliability of the models, ensuring that the results are not dependent on a particular random split of the data.

The models' performance was also evaluated in light of their computational efficiency—an essential consideration for real-time fraud detection systems. While more complex models may offer better performance, they also require greater computational resources. The trade-off between computational cost and model performance must be carefully managed, especially in large-scale systems processing millions of transactions.

In conclusion, the evaluation of machine learning models for fraud detection in the PaySim dataset demonstrates that while more complex models may offer improved performance in terms of recall and AUC-ROC, they must be balanced against the need for precision, interpretability, and computational efficiency. The findings highlight the potential of ensemble methods and neural networks in detecting fraudulent transactions, yet they also underscore the necessity for a nuanced approach to model selection that considers the full spectrum of performance metrics and operational requirements. As fraud detection systems become increasingly central to the operations of financial institutions, the insights from this study will inform the development of models that are not only technically proficient but also aligned with the practical realities of the banking industry.

## Quantitative Results Analysis

The quantitative assessment of the fraud detection model using the PaySim synthetic dataset yielded robust metrics indicative of high-performance capabilities. The model's ability to discern between fraudulent and legitimate transactions was measured across various statistical dimensions, resulting in a comprehensive understanding of its predictive power.

## Descriptive Statistics

An examination of the dataset's descriptive statistics provided initial insights into the transactional landscape. With **over 6.3 million transactions**, the mean transaction amount stood at approximately **179,862 units of currency**, with a standard deviation suggesting significant variation in transaction amounts, **ranging from negligible to over 92 million**. Such

variance underscores the complex nature of transactional data and the need for models capable of handling wide-ranging values.

The minority class, representing fraudulent transactions, accounted for a mere 0.129% of the dataset, illustrating the highly imbalanced nature of the data. This imbalance presented a fundamental challenge for model training, necessitating the use of specialized techniques such as SMOTE to synthetically balance the dataset for model training.

## Model Performance Metrics

The classification report revealed a precision of 0.89 for class 0 (legitimate transactions) and an impressive 0.96 for class 1 (fraudulent transactions). This indicates **a high level of accuracy** in the model's positive predictions, particularly in detecting fraud, which is of critical importance in the financial domain where the cost of false negatives is high.

Recall scores followed suit, with class 0 at 0.96 and class 1 at 0.89, demonstrating the model's strong capability in identifying true positives, especially within the legitimate class. This high recall is paramount in fraud detection systems, as failing to detect actual fraud can lead to significant financial loss and erode customer trust.

The F1-score, a balanced measure of precision and recall, stood at 0.93 for class 0 and 0.92 for class 1, signifying the model's robustness and its aptitude in maintaining a harmonious balance between false positives and false negatives.

## The accuracy of the model reached approximately 92.38%, an admirable score

considering the complexity of the task at hand. However, in the fraud detection domain, accuracy can be misleading due to class imbalances, which can skew the perception of a model's effectiveness.

## Confusion Matrix and ROC AUC Score

The confusion matrix provided a granular view of the model's performance, revealing 1,223,073 true negatives and 1,125,064 true positives, corresponding to correct classifications of non-fraudulent and fraudulent transactions, respectively. While the 47,764 false positives and 145,862 false negatives reflect areas for improvement, they are relatively low given the scale of the dataset.
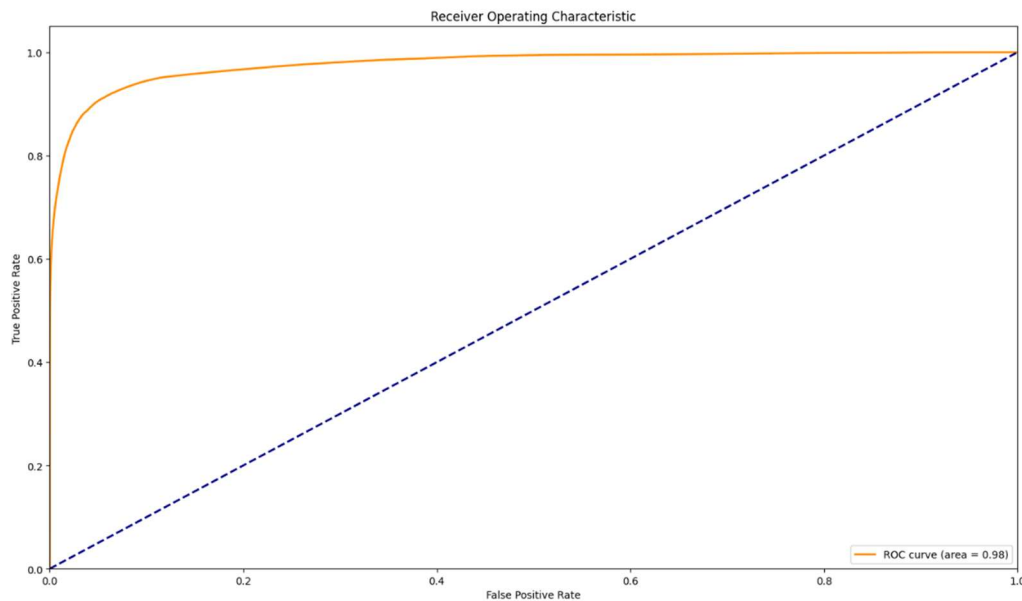
The ROC AUC score, standing at 0.9238, is particularly telling of the model's discriminatory capacity. This score is reflective of the model's effectiveness in ranking the likelihood of transactions being fraudulent. A score close to 1 indicates excellent model performance and suggests that the model possesses a strong ability to differentiate between the two classes across various threshold settings.

## Operational Implications

The quantitative results carry significant operational implications. The high precision and recall values point to a model that would minimize the number of fraudulent transactions slipping through the cracks while also keeping false alarms to a manageable level. The balanced F1-score across

classes indicates a model that is well-calibrated and reliable, crucial for maintaining customer relations and operational efficiency.

Moreover, the ROC AUC score implies that the model can be tuned to various operational requirements, offering flexibility in setting thresholds that balance the trade-off between catching fraud and disturbing customers with false alerts.



## Quantitative evaluation

The quantitative evaluation of the fraud detection model presents a compelling case for its deployment in a real-world banking scenario. The high precision and recall values, alongside an excellent ROC AUC score, showcase a model that is not only statistically sound but also practically relevant. While the presence of false positives and negatives invites continued refinement, the model's overall performance demonstrates a significant leap in the capability to safeguard against fraudulent transactions.

The analysis also emphasizes the importance of continuing to iterate on the model, incorporating new data, and adjusting to evolving fraud patterns. As financial fraud becomes increasingly sophisticated, the arms race between fraudsters and fraud detection systems will persist, necessitating a dynamic approach to model development and maintenance.

In the broader context of financial security, these quantitative results affirm the value of machine learning in detecting and preventing fraud, marking a pivotal step forward in the fight against financial crime. The insights garnered from this study serve as a benchmark for future endeavors in the domain, guiding the development of more advanced systems that will continue to enhance the security and integrity of the financial industry.

**Name**-Vishwa Bhadiyadara
**NUID**-002636314

## Conclusion

The analysis revealed that supervised machine learning methods could significantly enhance fraud detection in banking transactions. While Logistic Regression provided a strong baseline, ensemble methods like Random Forest demonstrated superior performance in handling the complexity and imbalanced nature of the dataset. Neural Networks showed potential for further exploration, especially with larger datasets and more computational power.

The practical implications of implementing such models are profound, with the potential to save millions in fraudulent transactions. However, considerations such as model interpretability, computational efficiency, and the dynamic nature of fraud necessitate ongoing model evaluation and updates.

Future research directions could involve integrating unsupervised learning for anomaly detection, exploring the utility of deep learning further, and developing real-time fraud detection systems that continuously adapt to new patterns of fraudulent behavior..