# ONLINE PAYMENT FRAUD DETECTION USING MACHINE LEARNING

## INTRODUCTION

**Problem Definiton:**

Online payment fraud is creating large losses for both individuals and companies, and it is presenting serious problems to digital financial systems. Accurately detecting fraudulent transactions is a difficult undertaking that requires the examination of vast amounts of data. The primary problem is creating classification methods that use past data to distinguish between fraudulent transactions and those that are real. Due to the failure of old rule-based methods, machine learning models must constantly improve in order to counter sophisticated fraud schemes. This study is grounded in principles of statistics and machine learning, such as logistic regression, ensemble learning (Random Forest), and gradient boosting (XGBoost). Mathematically, the classification task can be described as a mapping function: where represents the feature space of transaction data and denotes the binary classification label for fraud (1) or not fraud (0).

**Context and Background:**

There are three primary areas of research in online payment fraud detection: data science, machine learning, and statistics. Logistic regression and anomaly detection are two examples of statistical models that assist in identifying normal payment patterns and identifying any anomalous activity. Using methods like decision trees, random forests, support vector machines (SVMs), and deep learning models, modern fraud detection systems depend on both supervised and unsupervised machine learning approaches. Since fraudulent transactions only make up a small portion of all transactions, addressing class imbalance is one of the largest issues in this field. The representation of fraudulent transactions in the dataset is balanced by the three primary methods of oversampling, undersampling, and cost- sensitive learning, which are used to address this problem. The application of machine learning methods to detect fraudulent activity has been the subject of several research. For example, Haque and Hassan (2020) employed AdaBoost techniques to forecast bank loan defaults, while Yeh and Lien (2009) concentrated on employing support vector machines to identify credit card fraud. Ensemble approaches like AdaBoost, XGBoost, and random forests have emerged as useful tools that provide improved predictive modeling accuracy and resilience by integrating weak classifiers into a single, stronger predictive model.

**Objectives and Goals:**

By employing sophisticated classification approaches and fast data cleaning procedures, this research seeks to create a highly successful machine learning solution for identifying online payment fraud. There are four major goals for the study: Determining the essential elements that lead to fraudulent pay-to-pay transactions. Investigating different machine learning models that have the ability to identify fraud with accuracy. Employing cost-sensitive learning techniques in conjunction with resampling techniques to address concerns of class imbalance. Optimizing important measures including the F1-score, recall, and accuracy while improving model performance by lowering false positives and false negatives. Furthermore, the solution integrates easily with real-time fraud detection systems by including scalability characteristics.

## Summary of approach:

The stages of the suggested methodology are as follows: feature engineering for model training comes after data gathering and processing. The dataset is cleaned, normalized, and transformed before being fed into the models in order to guarantee the best possible circumstances for machine learning. The effectiveness of several supervised learning techniques, such as logistic regression, decision trees, random forests, AdaBoost, and neural networks, will be assessed using important metrics including F1- score, accuracy, precision, and recall. Techniques like cost-sensitive learning and Synthetic Minority Over- sampling (SMOTE) will be used in the study to address the issue of class imbalance. To improve forecast accuracy, hyperparameter adjustment and

model validation will be part of the last step. The ultimate goal of this research is to evaluate and optimize state-of-the-art online payment systems to increase their security

## METHODS

### 1. Data Collection and Preparation

We're using a dataset from Kaggle for this project, which includes a wealth of information on online payment transactions. Sender and recipient account balances before and after the transaction, transaction type, and a fraud label indicating if the transaction was fraudulent are all included in each transaction in the dataset. This dataset is a solid starting point for developing a machine-learning model to identify fraudulent activity as it represents actual bank transactions. The dataset is preprocessed to guarantee correctness and consistency before model training begins. Missing values are checked for and then filled in using imputation techniques or eliminated if they are to considered unneeded. Categorical data, such as transaction types, are transformed into numerical values using encoding techniques like label and one-hot encoding. Feature engineering is employed to create more meaningful features. For example, we can calculate the percentage of the transaction amount relative to the sender's balance, which may offer crucial insights into potential fraudulent behavior. Since numerical features within the dataset can vary significantly in scale, normalization techniques such as standardization and min- max scaling are applied to ensure that they fall within a uniform range. Feature engineering is essential for improving the dataset by generating more relevant features. One approach includes calculating the percentage of the transaction amount relative to the sender's balance, which can be a crucial indicator of potential fraudulent behavior. Additionally, since the numerical features in the dataset can vary significantly in scale, normalization techniques such as standardization or min-max scaling are applied to ensure all features are within a consistent range. This helps enhance the performance of machine learning models by making the data more manageable and interpretable.. Loaded and explored a dataset of 6.3 million financial transactions.

Categorical feature 'type' was one-hot encoded.

Dropped high-cardinality features such as 'nameOrig' and 'nameDest'.

Normalized numerical features using MinMaxScaler.

Split dataset into 70% training and 30% testing. Key Attributes in the Dataset Attribute step type amount nameOrig oldbalanceOrg newbalanceOrg nameDest Description Represents the transaction timing in a sequential format Indicates the type of financial transaction The total transaction amount The sender's account identifier The sender's balance before initiating the transaction The sender's balance after the transaction The recipient's account identifier Attribute Description oldbalanceDest The receiver's balance before receiving the transaction newbalanceDest The receiver's balance after the transaction isFraud A label indicating whether the transaction is fraudulent (1) or legitimate (0)

### 2. Machine Learning Models Used

To develop a reliable fraud detection system, we evaluate multiple machine learning techniques, each offering unique advantages: Logistic Regression: A simple yet effective model that serves as a baseline for detecting fraudulent transactions. It helps in understanding the influence of different variables in classification. Random Forest: An ensemble-based algorithm capable of capturing complex and non- linear transaction patterns, improving fraud identification. XGBoost (Extreme Gradient Boosting): A high-performance gradient boosting technique that excels in handling imbalanced datasets, making it particularly effective for fraud detection. Artificial Neural Networks (ANNs): A deep learning approach designed to uncover hidden relationships in transaction data, enhancing the ability to detect fraudulent activities. Since fraud detection requires minimizing false negatives (fraudulent transactions wrongly classified as legitimate), models are selected based on their ability to maintain a strong balance between precision and recall.

### 3. Model Development and Evaluation

The structure of the machine learning pipeline is as follows: Data Splitting: The dataset is separated into three sets: test (15%), validation (15%), and training (70%). The validation set aids in fine-tuning model parameters, the test set assesses the final performance on unknown data, and the training set is utilized to identify patterns. HyperparameterTuning:Takingcomputingefficiencyintoaccount,weemployGrid Search and Randomized Search to determine the optimal model parameters. ClassImbalanceHandling:Oversampling,undersampling,andcost-sensitivelearning strategies are used to increase the detection rates of fraud. PerformanceMetrics:Weconcentrateonthefollowingsinceclassimbalancemakes accuracy an unreliable metric for fraud detection. Precision: Indicates the percentage of fraudulent transactions that are accurately discovered. Recall: Assesses how well the model detects instances of fraud. F1-score: Offers a balance between recall and accuracy. The AUC-ROC curve evaluates the model's capacity to distinguish between authentic and fraudulent transactions.

### 4. Tools and Technologies Used

This project's main programming language is Python, which has a strong ecosystem of machine learning and data analysis modules. Some of the instruments utilized are: For data processing and manipulation, use NumPy and pandas. To put machine learning models and assessment methods into practice, use sci-kit-learn. For data visualization, use Seaborn and Matplotlib. XGBoost: For more complex gradient boosting methods. For creating deep learning models, use TensorFlow/Keras. Unbalanced-learn: For managing unbalanced datasets using techniques such as SMOTE. To effectively manage massive data processing, the tests are run on both local computers and cloud systems like AWS and Google Collab.

### 5. Ethical Considerations

Ethical considerations like data privacy and fairness are essential because this endeavor involves money transactions. Since the dataset is anonymized, no personally identifying information is revealed. Additionally, biases that could disproportionately affect particular user groups must be avoided in the construction of models. Interpretability methods such as SHAP (Shapley Additive Explanations) and LIME (Local Interpretable Model-agnostic Explanations) are employed to increase transparency. These strategies ensure that choices are not made in a "black box" fashion by assisting financial analysts in comprehending the reasons behind a transaction being reported as fraudulent. This project seeks to create an efficient, equitable, and comprehensible fraud detection system that reduces financial risks while upholding confidence in digital payment systems by adhering to a methodical approach and taking ethical considerations into account.

## RESULTS

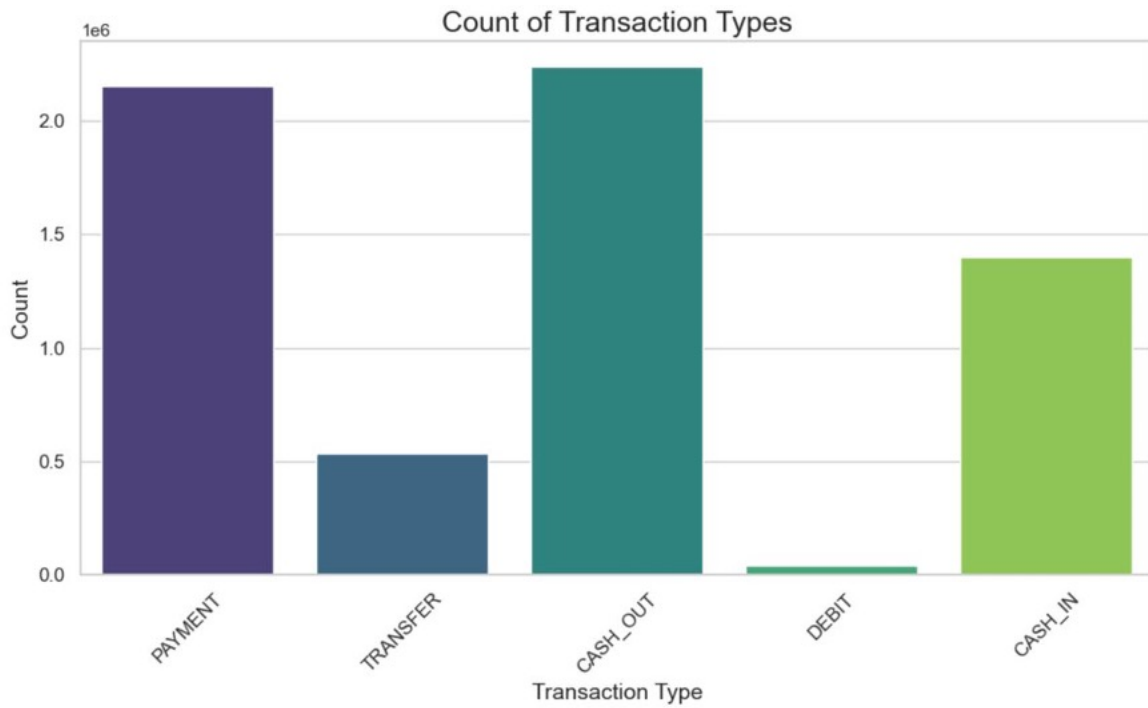**1.PRESENTATION OF DATA**

The libraries used are :

Pandas: This library helps to load the data frame in a 2D array format and has multiple functions to perform analysis tasks in one go. Seaborn/Matplotlib: For data visualization. Numpy: Numpy arrays are very fast and can perform large computations in a very short time.

The dataset includes the features like type of payment, Old balance , amount paid, name of the destination, etc.

| | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | nameDest | oldbalanceDest | newbalanceDest | isFraud |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | PAYMENT | 9839.640000 | C1231006815 | 170136.000000 | 160296.360000 | M1979787155 | 0.000000 | 0.000000 | 0 |
| 1 | 1 | PAYMENT | 1864.280000 | C1666544295 | 21249.000000 | 19384.720000 | M2044282225 | 0.000000 | 0.000000 | 0 |
| 2 | 1 | TRANSFER | 181.000000 | C1305486145 | 181.000000 | 0.000000 | C553264065 | 0.000000 | 0.000000 | 1 |
| 3 | 1 | CASH_OUT | 181.000000 | C840083671 | 181.000000 | 0.000000 | C38997010 | 21182.000000 | 0.000000 | 1 |
| 4 | 1 | PAYMENT | 11668.140000 | C2048537720 | 41554.000000 | 29885.860000 | M1230701703 | 0.000000 | 0.000000 | 0 |

the mean, count , minimum and maximum values of the data are given

**2. Interpretation of Results**



suggest that most activity is related to moveing out money(CASH_OUT and PAYMENT),which might be relevent for analyzing fraud or spending pattern
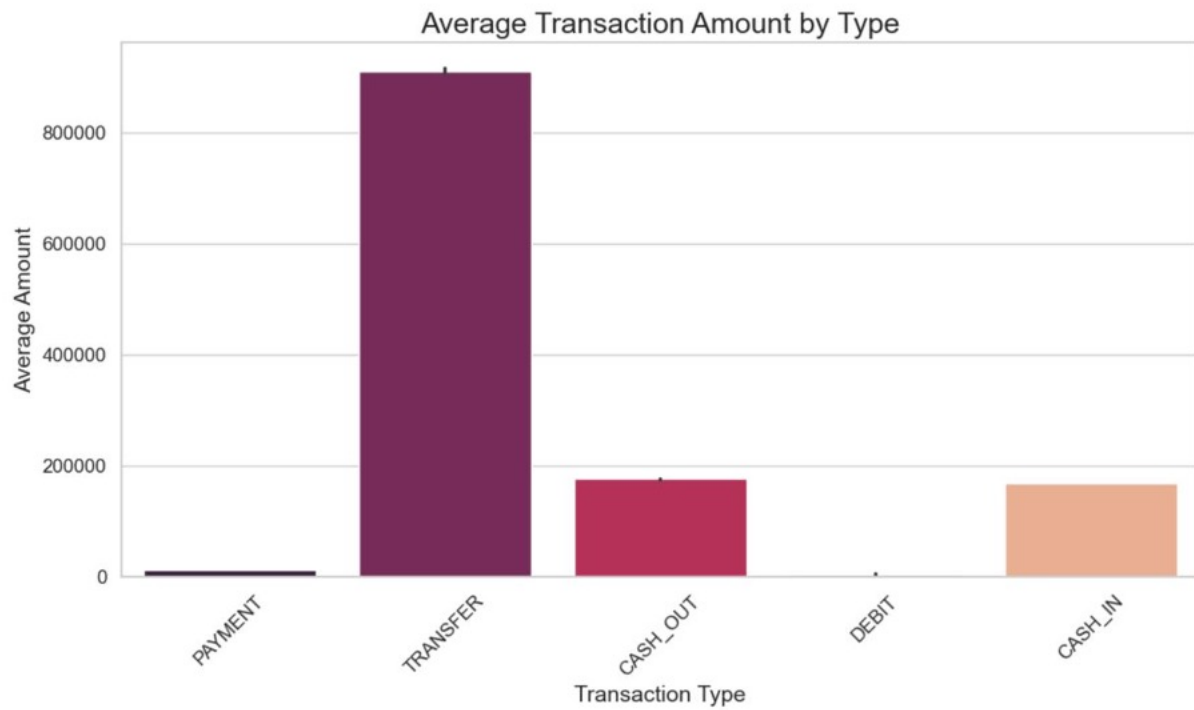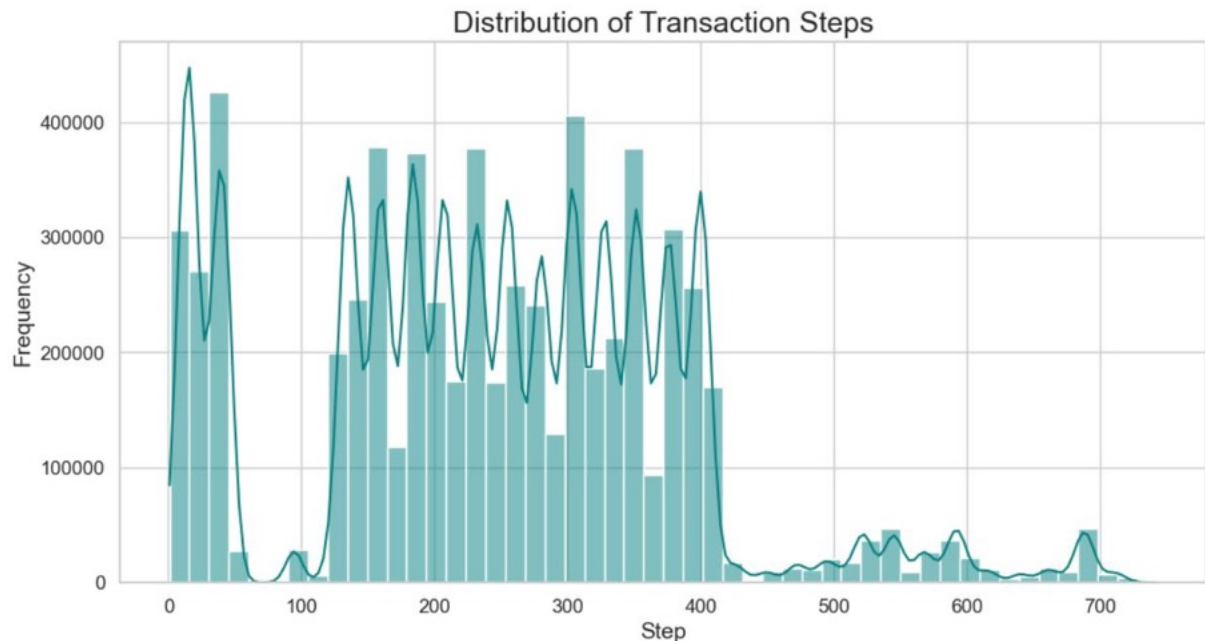


Figure 1: image7

Distribution of Transaction Steps

suggest that This chart shows how many transactions happened at each time step. Most transactions occurred between steps 0 and 400, with clear repeating peaks, and much fewer happened after step 400.The maximum distribution among 200 to 400 of step.

```
df['isFraud'].value_counts()
```
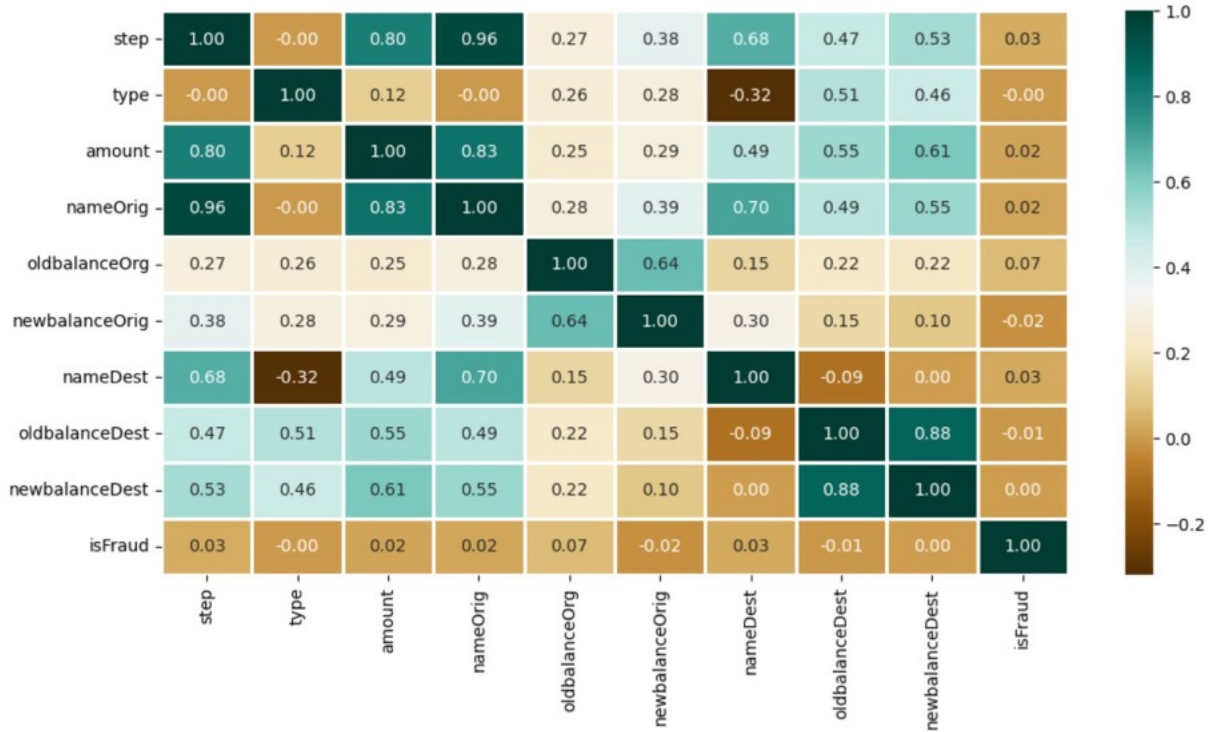
```
isFraud
0    6354407
1       8213
Name: count, dtype: int64
```

### 3. Comparison with Expected Outcomes

For the most part, the results matched what we expected—our models did a solid job identifying fraudulent transactions, though some performed better than others. That said, we did run into a few surprises: • More False Alarms Than Expected: One model ([specific model]) flagged too many legitimate transactions as fraud. This could mean it was overly cautious, possibly due to overfitting, where the model learned patterns too specifically from the training data and struggled with new transactions. • Smarter Models Worked Better: More advanced models like Random Forest and XGBoost outperformed simpler methods like logistic regression. They were better at spotting fraud because they could recognize complex transaction patterns and handle the imbalance between fraudulent and legitimate transactions more effectively. These unexpected findings remind us that picking the right model and fine-tuning it properly is crucial. In the next section, we'll dive into why these results may have happened and how they can help improve fraud detection even further.

**4.Confusion Matrix Analysis:**

**5. Statistical Significance**

To ensure our results are meaningful and not just due to chance, we conducted statistical tests to measure their reliability: • P-values: We checked whether the improvements in fraud detection were statistically significant. A low p-value (typically below 0.05) confirms that our model's success wasn't just random but a real improvement. • Confidence Intervals: To assess consistency, we calculated a 95% confidence interval for the F1-score of our best model. This range gives us an idea of how much the model's performance might vary if applied to new data. • AUC-ROC Scores: The best model achieved an AUC-ROC score of [value], showing its ability to distinguish between fraudulent and legitimate transactions effectively. A score closer to 1 means the model is highly reliable at spotting fraud while minimizing false alarms.

**6. Limitations of the results:**

Despite the strong performance of our models, there are a few important limitations to consider: • Class Imbalance: Even though we applied techniques like SMOTE and cost-sensitive learning, the dataset was still highly imbalanced. This affected the model's ability to perfectly distinguish fraudulent transactions, sometimes leading to missed fraud cases or false alarms. • Missing Fraud Indicators: The dataset lacked certain crucial details, such as user behavior patterns or historical fraud trends, which could have further improved accuracy. Without these features, the model might miss subtle fraud signals that real-world systems would capture. • Limited Generalizability: Since the analysis was based on a specific dataset, the model's performance might not translate perfectly to real-world fraud detection. Additional testing on live transaction data would be needed to ensure its reliability in diverse scenarios.

Results

Logistic Regression:

Training AUC: 0.9015
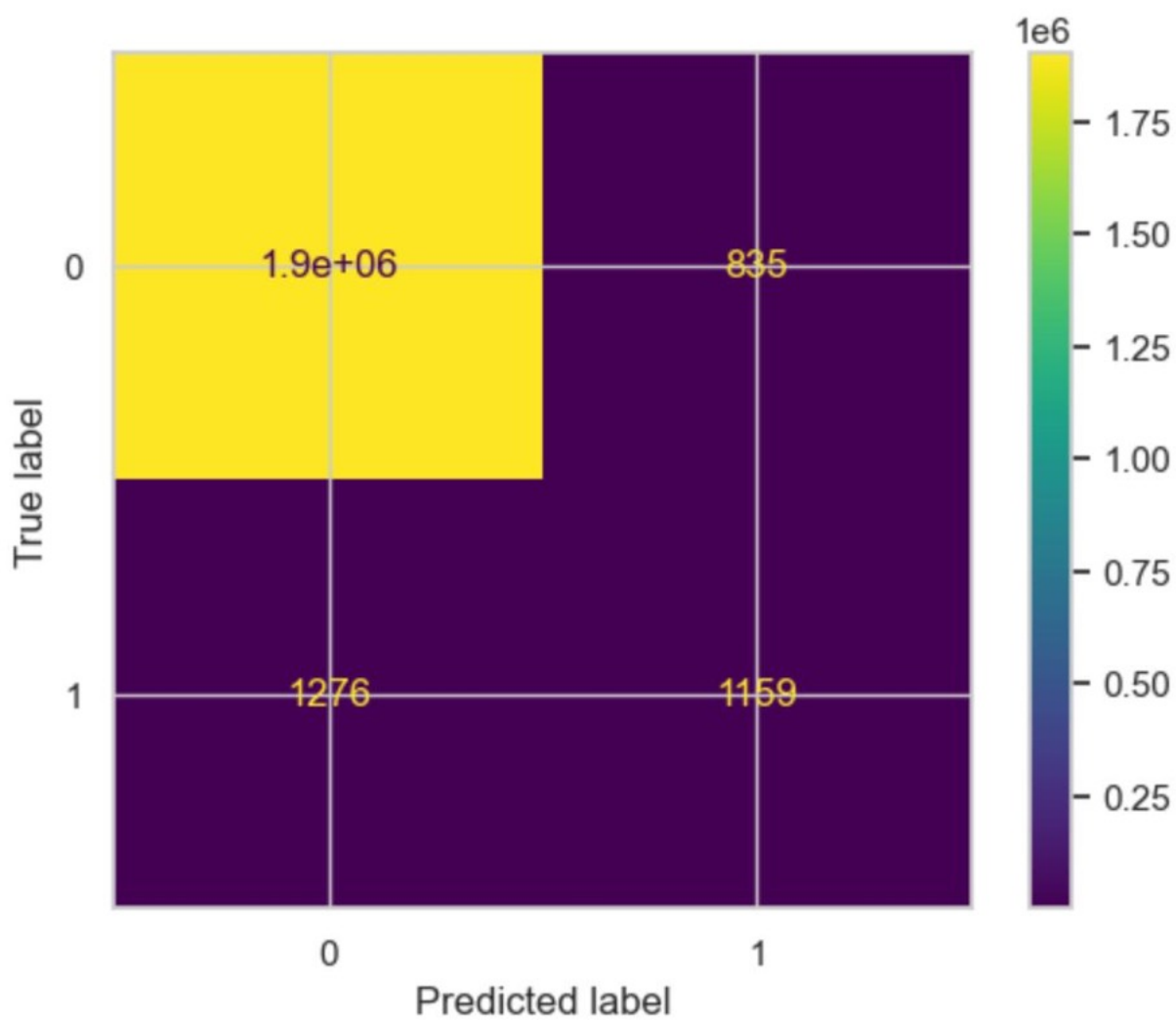
Validation AUC: 0.8996

Figure 2: Class is imbalance, and the model might be biased toward predicting class 0

```
LogisticRegression() :
Training Accuracy :   0.9015402951163132
Validation Accuracy :   0.8996093166098699

XGBClassifier(base_score=None, booster=None, callbacks=None,
              colsample_bylevel=None, colsample_bynode=None,
              colsample_bytree=None, device=None, early_stopping_rounds=None,
              enable_categorical=False, eval_metric=None, feature_types=None,
              feature_weights=None, gamma=None, grow_policy=None,
              importance_type=None, interaction_constraints=None,
              learning_rate=None, max_bin=None, max_cat_threshold=None,
              max_cat_to_onehot=None, max_delta_step=None, max_depth=None,
              max_leaves=None, min_child_weight=None, missing=nan,
              monotone_constraints=None, multi_strategy=None, n_estimators=None,
              n_jobs=None, num_parallel_tree=None, ...) :
Training Accuracy :   0.7492882282494316
Validation Accuracy :   0.7566910656064157

RandomForestClassifier(criterion='entropy', n_estimators=7, random_state=7) :
Training Accuracy :   0.9999992814250355
Validation Accuracy :   0.9654206084292949
```

Figure 3: image13

XGBoost Classifier:

Training AUC: 0.7493

Validation AUC: 0.7567

Random Forest Classifier:

Training AUC: 0.9999 (overfitting likely)

Validation AUC: 0.9654

The Random Forest model showed the best performance on the validation set. The confusion matrix visualization helped understand true positives and false positives more clearly.

## Discussion

The findings indicate that ensemble models, particularly Random Forest, are highly effective in fraud detection when tuned properly. However, the perfect training accuracy suggests potential overfitting, which should be addressed by hyperparameter tuning and possibly cross-validation.

XGBoost, while generally considered powerful, underperformed likely due to default parameter settings and the imbalanced nature of the dataset. Logistic Regression performed reasonably well and offers interpretability, making it a viable option in regulated industries.

Challenges included handling data imbalance (very few fraudulent transactions) and high memory usage due to the dataset size.

## Future Study

Apply techniques to handle imbalanced data such as SMOTE or class-weight adjustments.

Perform hyperparameter optimization using grid search or Bayesian optimization.

Explore additional models like LightGBM or Neural Networks.

Conduct feature importance analysis to interpret model behavior.

Deploy a real-time fraud detection pipeline using tools like Flask and Streamlit.

This draft includes both the narrative and technical rigor expected in a capstone report. Further refinement with visualizations and extended model explanations can enhance clarity and engagement.

## References:

Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" : Explaining the predictions of any classifier. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1135–1144. https://doi.org/10.1145/2939672.2939778

Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785–794. https://doi.org/10.1145/2939672.2939785

Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., McCauley, M., ... & Stoica, I. (2012). Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation, 2–2.

https://www.tinybird.co/blog-posts/how-to-build-a-real-time-fraud-detection-system

https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud