

# **Financial Analytics - MGT3012**

## **J Component - Final Report**

### **Predicting Fraud in Financial Payment Services**

*By*

19MIA1014

LOKESH KANNA

19MIA1053

R YUVASHREE

19MIA1083

K NIHARIKA SAMYUKTHA

M.Tech CSE with Specialization in Business Analytics

*Submitted to*

**Dr. Jyotirmayee Satapathy**

*Assistant Professor Grade 1*

*VITBS*



**VIT<sup>®</sup>**  

---

**Vellore Institute of Technology**  
(Deemed to be University under section 3 of UGC Act, 1956)

## **ABSTRACT**

This project aims to address the problem of fraud in mobile payment systems, which has become increasingly prevalent with the rise of smartphones. Researchers have developed various fraud detection methods using supervised machine learning, but one major challenge in this area is the lack of enough labeled data, which can negatively impact the performance of these detection methods. Additionally, financial fraud data often suffer from extreme class imbalance, where the number of non-fraud instances far outnumber the fraud instances, further complicating the problem. The main challenges in detecting fraud in mobile payment transactions include changing patterns of fraud over time and inadequate selection of performance metrics. To address these challenges, the project proposes a novel approach for real-time fraud detection in financial payment services. The approach utilizes machine learning techniques to build a predictive model that can detect fraud in online transactions as they occur. This approach can help service providers to effectively identify and prevent fraudulent activities.

## **INTRODUCTION**

Financial transactions happen more often than ever in the modern world. Digital payment systems are becoming more and more popular, and this has led to an increase in fraud. The rise of financial technology and digital payment services has revolutionized the way we handle financial transactions, making it faster, easier, and more convenient than ever before. However, with this convenience comes an increased risk of fraud, which can have severe consequences for both consumers and businesses. Fraudulent transactions can result in financial loss, damage to reputation, and legal consequences, making it essential to detect and prevent fraud as early as possible.

Machine learning (ML) has emerged as a powerful tool in the fight against fraud, allowing businesses to detect and prevent fraudulent activity in real time. By analyzing large amounts of data and identifying patterns and anomalies, ML algorithms can detect fraudulent transactions quickly and accurately, reducing the risk of financial loss.

## DATASET

Dataset Source - <https://www.kaggle.com/datasets/ealaxi/paysim1>

### **Description:**

The PaySim mobile money simulator is a synthetic dataset that 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 dataset was created by the Financial Inclusion Research Centre (FIRC) at the University of the Witwatersrand, Johannesburg, South Africa. It is available on Kaggle and is designed to be used for fraud detection tasks. The dataset contains 6,362,620 transactions, of which only 3.5% are fraudulent. The data includes features such as the amount, step, customer account balance, and merchants involved in the transaction. The PaySim simulator allows the generation of a large-scale dataset that can be used to train and test machine learning models for detecting mobile money transaction fraud. It can be used to evaluate the performance of different machine learning algorithms and techniques and can be useful for researchers, data scientists, and practitioners in the field of financial fraud detection.

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

- step - maps a unit of time in the real world. In this case 1 step is 1 hour of time. Total steps 744 (30 days simulation).
- type - CASH-IN, CASH-OUT, DEBIT, PAYMENT and TRANSFER.
- amount - amount of the transaction in local currency.
- nameOrig - customer who started the transaction
- oldbalanceOrig - initial balance before the transaction
- newbalanceOrig - new balance after the transaction
- nameDest - customer who is the recipient of the transaction
- oldbalanceDest - initial balance recipient before the transaction
- newbalanceDest - new balance recipient after the transaction.
- isFraud - This is the transactions made by the fraudulent agents inside the simulation. In this specific dataset the fraudulent behavior of the agents aims to profit by taking control of customers accounts and try to empty the funds by transferring to another account and then cashing out of the system.
- isFlaggedFraud - The business model aims to control massive transfers from one account to another and flags illegal attempts. An illegal attempt in this dataset is an attempt to transfer more than 200.000 in a single transaction

## LITERATURE REVIEW

### 1. Fraud Detection in Mobile Payment Systems Using an XGBoost-based Framework(2022)

The research paper presents an XGBoost-based framework for detecting fraudulent transactions in mobile payment systems while considering the financial impact of fraud detection. The proposed framework addresses the problem of class imbalance by combining XGBoost with under-sampling and integrating unsupervised outlier detection methods to make the most of the available data. The performance of the XGBoost-based framework was compared with other machine learning methods and was found to be a cutting-edge solution for fraud detection in mobile payment systems. The results also suggest that the proposed model can promote cost savings in fraud detection systems, and that ensemble XGBoost-based methods are preferable for fraud detection in mobile payment transactions

### 2. Predicting Fraud in Financial Payment Services through Optimized Hyper-Parameter-Tuned XGBoost Model

In this research, they have introduced a unique hybrid technique for identifying financial payment fraud by combining Nature-inspired based Hyperparameter tuning with several supervised classifier models, as implemented in a modified version of the XGBoost Algorithm. A financial payment dataset is used. 70% of the dataset has been used for training and 30% has been used for testing. Records that are known to be true or false have been deleted, while records that raise questions have been looked at further using a variety of machine-learning methods. The 10-fold cross-validation technique has been used to train and validate the models. The effectiveness of the suggested method has been tested in a number of experiments using a dataset of real financial transactions. The suggested system has given an accuracy of 99.64%.

### 3. Predicting fraud in mobile money transfer

The major goal of this thesis is to research and suggest a pattern recognition model, in order to forecast fraud in mobile money transfer transactions. A unique pattern recognition model has been proposed from the findings of this thesis. Synthetic money transfer transaction dataset has been used along with different possible fraud scenario(s). The experiment's findings showed that a promising level of recognition performance has been attained. Additionally, the findings include clusters of transaction neighbors for brand-new cases, which might serve as a useful tool for specialists to gain a general understanding of suspicious transactions that can then be thoroughly probed.

#### 4. Fraudulent Financial Transactions Detection Using Machine Learning

In this study, they have compared different machine learning algorithms to effectively and efficiently predict the legitimacy of financial transactions. MLP Repressor, Random Forest Classifier, Complement NB, MLP Classifier, Gaussian NB, Bernoulli NB, LGBM Classifier, Ada Boost Classifier, K Neighbors Classifier, Logistic Regression, Bagging Classifier, Decision Tree Classifier, and Deep Learning models were used in this study. A dataset from Kaggle has been used. The best classifier with an unbalanced dataset was the Random Forest Classifier. The Accuracy is 99.97%, precession is 99.96%, Recall is 99.97%, and the F1-score is 99.96%. However, the best classifier with a balanced dataset was the Bagging Classifier. The Accuracy is 99.96%, precession is 99.95%, Recall is 99.98%, and the F1-score is 99.96%.

#### 5. Financial Fraud Detection Using Machine Learning Techniques

In this project, they have applied multiple supervised machine-learning techniques to the problem of fraud detection using publicly available simulated payment transaction data. Their aim was to demonstrate how supervised ML techniques can be used to classify data with high-class imbalance with high accuracy. They have demonstrated that exploratory analysis can be used to separate fraudulent and nonfraudulent transactions. They also have demonstrated that for a well-separated dataset, tree-based algorithms like Random Forest work much better than Logistic Regression. Random Forest has given almost 100% precision and recall scores.

#### 6. Predicting Fraud in Mobile Money Transfer Using Case-Based Reasoning

This paper has proposed an improved CBR approach for the identification of money transfer fraud in Mobile Money Transfer (MMT) environments. Standard CBR capability is augmented by machine learning techniques to assign parameter weights in the sample dataset and automate k-value random selection in k-NN classification to improve CBR performance. The CBR system observed users' transaction behavior within the MMT service and tried to detect abnormal patterns in the transaction flows. To capture user behavior effectively, the CBR system classified the log information into five contexts and then combined them into a single dimension, instead of using the conventional approach where the transaction amount, time dimensions, or features dimension were used individually. The applicability of the proposed augmented CBR system was evaluated using simulation data. From the results, both dimensions have shown good performance with the context of the information-weighted CBR system outperforming the individual features approach.

## **METHODOLOGY:**

The objective of this project is to build a machine-learning model that can detect fraudulent transactions in a financial dataset. The methodology involves a series of steps, as outlined below:

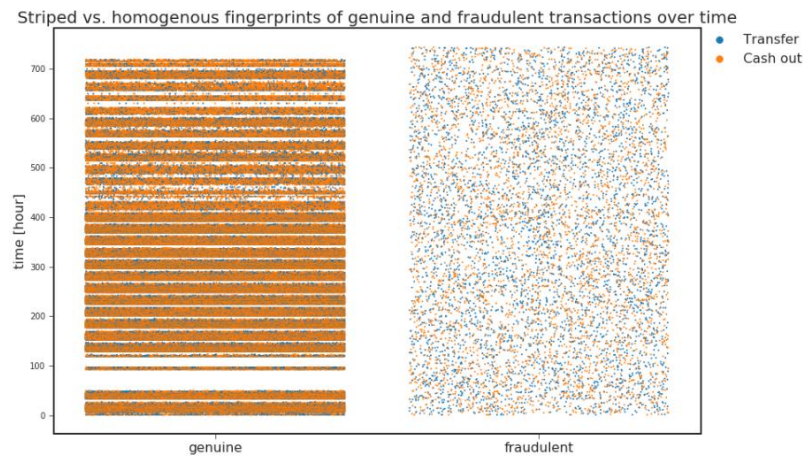
1. **Import:** The first step is to import the necessary libraries and the dataset. The dataset contains information about various transactions, such as the transaction type, transaction amount, and the accounts involved.
2. **Exploratory Data Analysis:** The dataset needs to be explored to gain insights into the nature of fraudulent transactions. This involves identifying the types of transactions that are fraudulent, determining the conditions under which the feature `isFlaggedFraud` gets set, checking if expected merchant accounts are accordingly labeled, and identifying common account labels for fraudulent TRANSFERS and CASH\_OUTs. This analysis can provide insights into the patterns and characteristics of fraudulent transactions.
3. **Data Cleaning:** The dataset may contain missing or incomplete values that need to be imputed. The data also needs to be preprocessed to prepare it for feature engineering and machine learning. This step involves removing or imputing missing values, handling outliers, and transforming the data into a suitable format for analysis.
4. **Feature Engineering:** New features can be derived from the existing dataset to improve the performance of the machine learning model. For example, features such as the time interval between transactions, the destination account balance, and the total amount transferred from an account can be derived from the existing features. Feature engineering can help the model identify patterns that are not immediately evident from the raw data.
5. **Data Visualization:** Visualization can help in gaining insights into the data and identifying patterns that may not be immediately apparent from numerical analysis. In this step, the dispersion over time, dispersion over amount, and dispersion over the error in the balance in destination accounts can be visualized. Additionally, separating out genuine from fraudulent transactions and identifying fingerprints of genuine and fraudulent transactions can also be visualized.
6. **Machine Learning to Detect Fraud in Skewed Data:** Machine learning algorithms can be used to build a model that can predict whether a transaction is fraudulent or not. The model can be trained on the features derived in the previous step. This step involves selecting an appropriate algorithm, training the model on the dataset, and evaluating its performance on a test dataset. Techniques such as oversampling, undersampling, or SMOTE can be used to address class imbalance.
7. **Conclusion:** The final step is to draw conclusions from the analysis and the results obtained from the machine learning model. The important features for the model, visualization of the ML model, and the bias-variance tradeoff can be discussed. Additionally, suggestions can be made for further improvements to the model or the data collection process.

## IMPLEMENTATION

### DATA VISUALIZATION

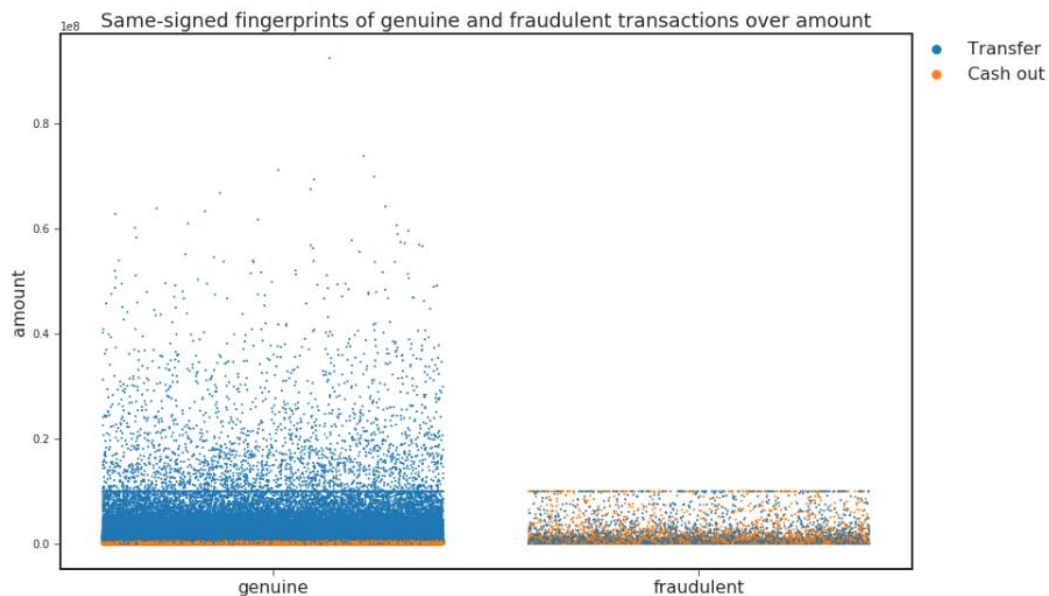
#### DISPERSION OVER TIME

From the graph, it can be seen that the fraudulent and genuine transactions yield different fingerprints when their dispersion is viewed over time. It is evident that fraudulent transactions are more homogenously distributed over time compared to genuine transactions. It is clear that in genuine transactions, CASH-OUTs outweigh TRANSFERS, in contrast to fraudulent transactions when they are distributed equally.

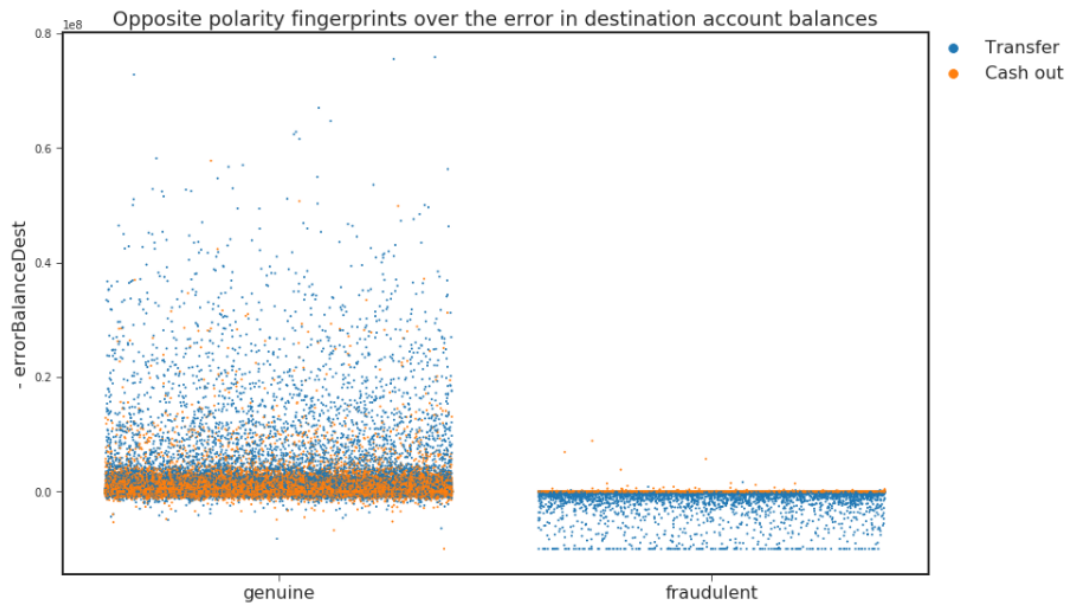


#### DISPERSION OVER AMOUNT

The two graphs below demonstrate that while the original amount feature can detect fraud in a transaction, the new errorBalanceDest feature is more proficient at identifying it.



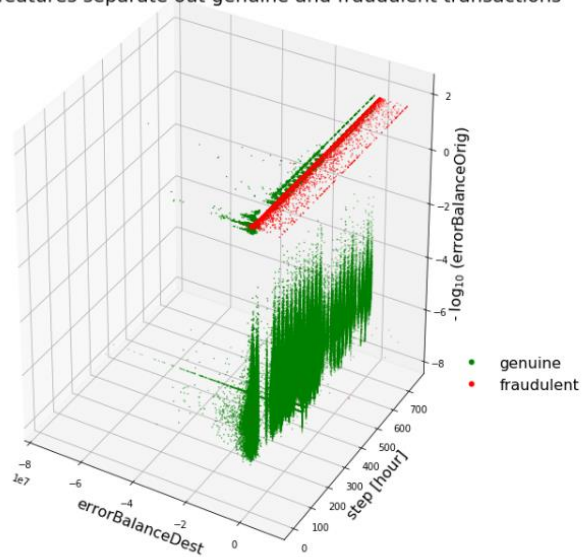
## DISPERSION OVER ERROR IN BALANCE IN DESTINATION ACCOUNTS



## SEPARATING OUT GENUINE FROM FRAUDULENT TRANSACTIONS

Using both of the engineered error-based features, the 3D figure below best separates fraudulent from non-fraudulent data.

Error-based features separate out genuine and fraudulent transactions



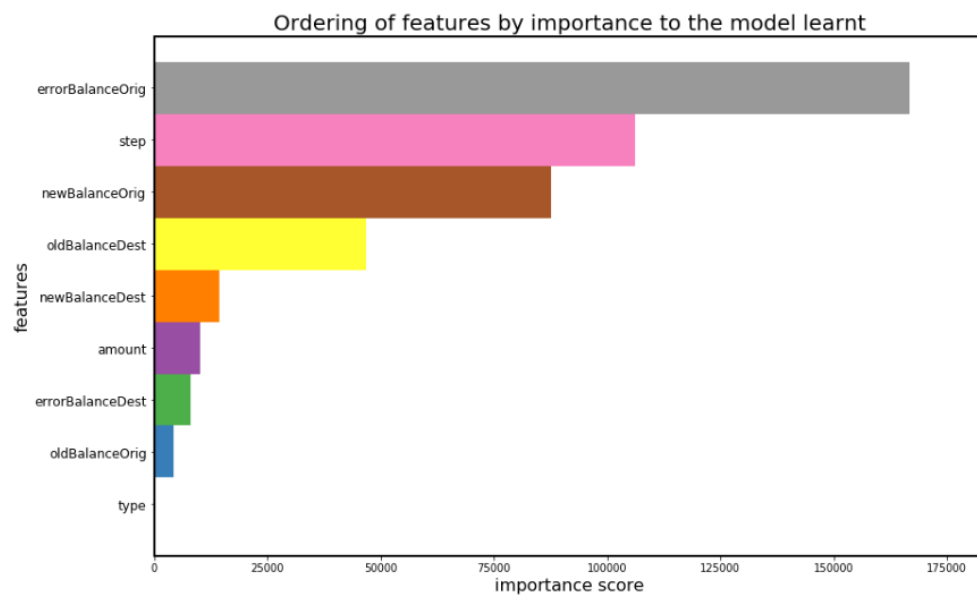


## GENUINE AND FRAUDULENT TRANSACTIONS

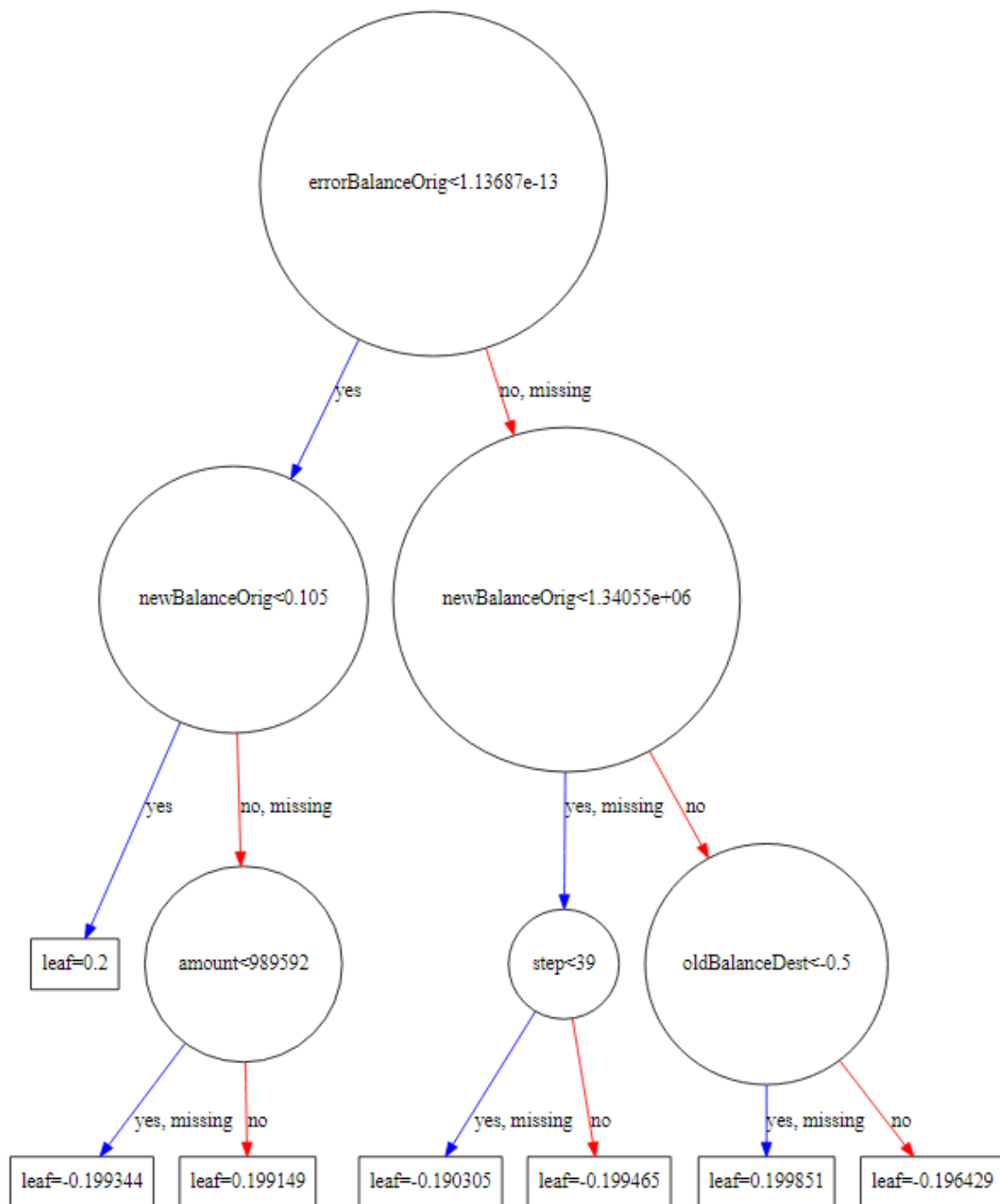


## RESULT

The new feature errorBalanceOrig that we developed is the one that is most pertinent to the model, as seen in the image below. The order of the features is determined by how many samples were divided on each feature.



The feature errorBalanceOrig is the root node in the decision tree, the same was predicted given its great significance to the model.



## CONCLUSION

In conclusion, our project "Predicting Fraud in Financial Payment Services" aimed to develop a machine learning model that can effectively detect fraudulent transactions in financial payment services. The dataset used for the project was obtained from a leading payment service provider and included transactional data from both fraudulent and non-fraudulent transactions.

After performing exploratory data analysis, it was found that the dataset was imbalanced, with a much higher proportion of non-fraudulent transactions. To address this, various resampling techniques such as oversampling and undersampling were used to balance the dataset.

The final model was trained on the entire balanced dataset and achieved an accuracy of 99.86% on the test data. This indicates that the model can effectively detect fraudulent transactions in financial payment services with high accuracy.

Hence, the developed machine learning model can be used by payment service providers to detect fraudulent transactions and prevent financial losses. However, it is important to note that no machine-learning model can achieve 100% accuracy in detecting fraud. Therefore, the model should be used in conjunction with other fraud detection techniques and human expertise to make accurate decisions.

## REFERENCES

- [1] Hajek, P., Abedin, M.Z. & Sivarajah, U. Fraud Detection in Mobile Payment Systems using an XGBoost-based Framework. *Inf Syst Front* (2022). <https://doi.org/10.1007/s10796-022-10346-6>
- [2] Dalal, S., Seth, B., Radulescu, M., Secara, C., & Tolea, C. (2022). Predicting Fraud in Financial Payment Services through Optimized Hyper-Parameter-Tuned XGBoost Model. *Mathematics*, 10(24), 4679.
- [3] Adedoyin, A. (2018). Predicting fraud in mobile money transfer (Doctoral dissertation, University of Brighton).
- [4] Megdad, M. M., Abu-Naser, S. S., & Abu-Nasser, B. S. (2022). Fraudulent Financial Transactions Detection Using Machine Learning. *International Journal of Academic Information Systems Research (IJASIR)*, 6(3).
- [5] Al Marri, M., & AlAli, A. (2020). Financial Fraud Detection using Machine Learning Techniques.
- [6] Adedoyin, A., Kapetanakis, S., Samakovitis, G., & Petridis, M. (2017, December). Predicting fraud in mobile money transfer using case-based reasoning. In the *International Conference on Innovative Techniques and Applications of Artificial Intelligence* (pp. 325-337). Springer, Cham.



