

A report on Fraud Detection in Mobile Payment system
using an

XGBoost based framework

INTRODUCTION

Mobile payments allow users to deposit, withdraw, spend, transfer, and send money using their phones. Almost 300 services exist globally, with high usage in Sub-Saharan Africa and Asia. In 2020, mobile payments totaled \$767 billion with 1.2 billion registered users. Statistics reports its enormous potential during COVID-19 to increase promptness and efficiency while minimizing face-to-face contact. Commercial and technical factors are crucial for their future growth. Regarding cost efficiency, mobile payments in developing countries tend to be high volume but low value. Security is a technical concern due to inadequate legal frameworks and enforcement mechanisms. Therefore it is necessary to control fraudulent activities to increase smooth and efficient mobile payment methods. The rise of mobile payments increases the risk of fraud. Security measures are being investigated to prevent attacks and develop countermeasures.

Detecting fraudulent transactions is key, especially with mobile payment services. Automated detection systems are needed for immediate prevention. The challenges in detecting fraud in mobile payments include extreme class imbalance, changing fraud patterns, and inadequate performance metrics. This results in a poor user experience for legitimate customers as their transactions may be rejected. second challenge reduces performance and efficiency, requiring updates for machine learning models to meet objectives. For the challenge, mobile payment providers face difficulty in selecting the appropriate balance between false positives and false negatives in detecting fraud. Studies have found high accuracy in supervised and deep learning methods. A major issue is the extreme class imbalance of transactions, with a dominance of legitimate transactions. This leads to poor classification performance on the minority class of fraudulent transactions. Two approaches have been used to address the issue, The first relies on under-sampling to balance the dataset (Pambudi et al., 2019). Limitation: Approach loses important information in discarded transactions, reducing accuracy. Attempt made to isolate fraud with unsupervised outlier detection (Buschjäger et al., 2021). However, a thorough assessment of Machine learning methods and their integration for better detection results has not been published. To address problems with extreme class imbalance in mobile payment transactions, we propose enhancing XGBoost by including a data sampling component. In financial applications, filtering unusual observations is necessary to ensure reliability and prevent malicious use. This is useful for detecting financial fraud attempts, as their behavior differs from normal transactions (Bernard et al., 2021). Outlier detection methods process real-time data to identify organized crime groups accurately and with fewer false positives than traditional supervised learning. Outlier detection methods are effective in detecting credit card, banking, and health insurance fraud. However, these methods have seen limited use in detecting financial fraud despite suggestions that they deserve more attention due to the limitations of supervised

algorithms. The scarcity of outlier detection methods may be due to the challenge of identifying fraudulent activity when it overlaps with legitimate behavior in noisy datasets. Detecting outliers in the financial domain is difficult due to challenges such as a lack of efficient methods and differences in behavior between legitimate and fraudulent scenarios. Secondly, unsupervised learning is favored due to scarce labeled data. Thirdly, legit behavior shifts with time, and fraudsters disguise their acts.

This study mainly focuses on three aspects:

- To create a new method of detecting fraud in mobile payment systems by combining XGBoost with class-balancing adjustments and unsupervised outlier detection. This approach will be effective in detecting fraud in scenarios where there is a class imbalance.
- In addition, the study proposes a new way to evaluate the performance of mobile payment fraud detection systems. Unlike traditional measures, proposed method takes into account both the cost savings from correctly identifying fraudulent transactions and the losses incurred from falsely identifying legitimate transactions as fraudulent.
- To demonstrate the effectiveness of the fraud detection framework, they use the PaySim dataset, which contains over 6 million mobile payment transactions. the results show that the framework not only outperforms existing fraud detection methods in terms of accuracy but also saves providers of mobile payment systems a significant amount of money.

Literature Review:

From evaluating and assessing the literature review in the study the literature lacks a thorough assessment of the latest machine learning-based techniques that address the issue of class imbalance through under-sampling methods. Additionally, there has been a lack of attention given to hybrid semi-supervised approaches that combine supervised learning and unsupervised outlier detection methods. Furthermore, evaluations of fraud detection performance in mobile payment systems have only used standard performance metrics, without considering the financial impact of fraud detection.

METHODOLOGY

Proposed fraud detection model:

There are two main fraud detection model given in the study:

- Extreme Gradient boosting (XGBoost) method
- XGBOD Method – Extreme Gradient Boosting Outlier Detection Method

Extreme Gradient Boosting Method:

Enhanced with random under-sampling, it is proposed to use both the supervised learning capacity and robustness of XGBoost, a cutting-edge machine learning algorithm, and the data sampling component to solve the class imbalance problem inherent in mobile payment transaction data.

Under sampling for handling class imbalance problem

The high imbalance between legitimate and fraud classes makes fraud detection a challenging work. Taking into consideration the importance of class imbalance in financial fraud detection, enormous methods have been implied to improve the classification performance of supervised learning methods.

Over-sampling creates artificial instances in minority class but can lead to over fitting and ignores majority class. Under sampling is better due to increasing financial fraud data. RUS method used in this study controls the number of samples selected from original data. It's a non-heuristic approach that randomly selects a subset from majority class, enabling sampling of heterogeneous data efficiently (Haixiang et al., 2017).

Extreme Gradient Boosting

XGBoost uses gradient boosted decision trees to build additive models incrementally, minimizing overall error. This creates an ensemble of base learners that outperform individual classifiers. The model is regularized to control over fitting and improves robustness to noise through stochastic gradient boosting. The objective function minimizes error (Chen & Guestrin, 2016).

$$\text{obj}(1)(t) = \sum (y_i - \hat{y}_{(t-1)} + f_t(x_i))^2 + \sum \Omega(f_t)$$

Extreme Gradient Boosting Outlier Detection Model

XGBOD is a semisupervised ensemble approach that combines numerous unsupervised outlier detection methods with an XGBoost classifier to discover anomalies. Unsupervised outlier detection, feature selection, integrating the outlier score matrix with original features, and training an XGBoost classifier on an enhanced feature space are all steps in the process. The balance selection method is utilised to maintain diversity and accuracy by selecting the most relevant TOS features. It uses a discounted accuracy function to select the subset of p most relevant TOS based on accuracy, similarity, and a discounting factor for similarity.

Machine Learning Methods for Comparative Evaluation

The Machine Learning Methods used for Comparative Evaluation are :

- Supervised Learning Methods for Imbalanced data
- Outlier Detection Method

Supervised Learning Methods for Imbalanced data:

k-nearest neighbor classifier -

The k-nearest neighbour (k-NN) approach is a non-parametric classifier that compares training and test cases. It is an instance-based strategy that uses a majority vote to classify an instance by examining its k most-similar instances, often based on Euclidean distance. This straightforward method has been shown to be effective in detecting highly unbalanced credit card fraud. k-NN can be employed efficiently in unsupervised outlier identification mode in financial fraud detection, when fraud examples are believed to be far from legitimate class samples.

Support Vector Machine –

The Support Vector Machine technique is an effective classifier for detecting fraud in high-dimensional data. It calculates the best hyperplane to maximise the margin between classes, which is represented by support vectors. Kernel functions are used to address nonlinear interactions by mapping to a new feature space where linear separation is possible.

Random Forest –

Random Forest (RF) is a machine learning technique that generates a number of trees and ensures the generalisation error converges to a specified limit by using several decision trees trained separately on distinct input samples. RF has a non-differentiable decision boundary, and random feature selection is utilised to separate the nodes in each tree, making the RF classifier more resilient to noise. Because of its hierarchical structure, RF is especially useful in detecting financial fraud when the class distribution is uneven. The good performance of RF on financial fraud detection jobs is well documented.

Outlier Detection Method

This section discusses unsupervised machine learning methods for outlier detection. The methods aim to represent legitimate data using clusters and assign a score to unseen instances for comparison against a decision boundary. The evaluation in this study includes proximity-based, linear model-based, ensembling, and neural network-based methods for outlier detection.

Proximity Based Methods:

Proximity-based approaches detect outliers by analysing the neighbourhood of each data instance and do not require prior knowledge of data distribution. They do not, however, scale well for high-dimensional data. LOF and k-NN are two examples, whereas CBLOF and ABOD address the curse of dimensionality. HBOS assumes feature independence and is computationally efficient, however due to density estimation restrictions, it fails to detect local outliers.

Linear Model- Based Methods:

Linear model-based methods create decision boundaries to distinguish genuine class instances from the remainder of the data space. OCSVM builds a hyperplane in high-dimensional space by minimising structural risk and allowing certain instances to fall outside the boundary to minimise overfitting. MCD combines a multivariate location and scale estimator with a resilient clustering technique to minimise the determinant of the covariance matrix for each cluster, and outlier scores are generated using the Mahalanobis distance. However, excessive overlap across clusters can cause convergence issues.

Ensembling Methods

Isolation Forest computes an isolation score to find outliers from the rest of the data samples. It assumes that outliers can be isolated closer to the tree's root and computes the isolation score using the average path length. LODA is a collection of weak learners represented by one-dimensional histograms that approximate the likelihood of random data projections. It is robust to both a large number of samples and missing data, allowing it to detect anomalous samples in real time.

Neural Network-Based Methods

Feature learning is used by neural network-based approaches such as autoencoders (AEs) and variational autoencoders (VAEs) to reduce dimensionality and detect outliers. AEs can be trained to learn valid behaviour and compute a reconstruction error as an outlier score, but VAEs use joint data distribution and generative elements for robust disentangled representation learning. GANs have also been used in unsupervised outlier identification, with multi-objective generative adversarial active learning (MO-GAAL) using GANs to identify fraudulent outliers from valid data. MO-GAAL has proven to be effective in detecting financial fraud.

CONCLUSION

The study used the Paysim Dataset, which included 6,362,620 mobile transactions, out of which 8,213 were fraudulent, and the data was randomly partitioned into a 3:1 ratio of training to testing data. This approach used data from real world which was pre-processed and engineered to generate predictions of fraudulent transactions, which were then analyzed with XGBoost algorithm.

The study's findings revealed that the XGBoost algorithm was extremely effective in detecting fraudulent transactions, with an accuracy rate of 98.4%. In terms of accuracy, precision, and recall, the approach outperformed other machine learning algorithms such as random forests and support vector machines. The XGBoost algorithm's ability to handle imbalanced datasets is one of its primary characteristics. The number of fraudulent transactions in fraud detection is often substantially lower than the number of genuine transactions, resulting in an imbalanced dataset. Despite the class imbalance, the XGBoost algorithm is capable of handling this imbalance and making correct predictions.

The XGBoost algorithm's precision and recall were likewise high, indicating that it was able to detect the majority of fraudulent transactions while minimising false positives. This is significant because false positives can be costly, resulting in customer unhappiness and financial losses.

The performance of the XGBoost algorithm was also evaluated in comparison to that of random forests and support vector machines, two other machine learning algorithms. Both of these methods were exceeded in terms of accuracy, precision, and recall by the XGBoost algorithm. This demonstrates how the XGBoost algorithm is excellent at identifying fraudulent transactions

in mobile payment systems. Overall, the study's findings offer convincing proof of the XGBoost-based framework's ability to identify fraudulent transactions in mobile payment systems. The algorithm has the ability to increase the security and integrity of mobile payment systems due to its high accuracy, precision, and recall levels.

It is significant to highlight that when implementing such a system, the technique has some constraints that must be taken into mind. To increase the accuracy of the system, it is necessary to address two major limitations: the potential for false positives and the requirement for continual monitoring and updating of the model.

The findings of the paper "Fraud Detection in Mobile Payment Systems Using an XGBoost-based Framework" show that machine learning-based methods can enhance fraud detection in mobile payment systems. The XGBoost algorithm fared better than other machine learning algorithms in terms of accuracy, precision, and recall when it came to spotting fraudulent transactions. The study demonstrates the potential of machine learning-based technologies to enhance the security and integrity of mobile payment systems and offers insightful information on the significance of various indicators in predicting fraudulent transactions.