

Real-Time Fraud Detection in Financial Transactions

Anmol Chaurasia
Department of Computer Science and Engineering
Galgotias University
Greater Noida, Uttar Pradesh
anmol.22scse1180158@galgotiasuniversity.edu.in

Nilesh Chandra
Department of Computer Science and Engineering
Galgotias University
Greater Noida, Uttar Pradesh
nilesh.22scse1180131@galgotiasuniversity.edu.in

Abstract--With mobile finance and digital banking industries continuing to grow exponentially, the traditional security measures based on rules are finding it hard to keep up with the more devious schemes used by criminals. This article presents a unified real-time system of fraud detection that combines machine learning, graph analytics, and big data processing. Our pattern recognition based on supervision and anomaly detection without supervision and learning through relational graphs recognizes both the previously recorded types of fraud and the previously unknown organized attacks. In order to provide high-throughput and low-delays, we use a streaming pipeline based on the Apache Spark and Apache Kafka. We have shown that this multi-layered approach is effective in managing the data imbalance and privacy limitations and provides a robust defense mechanism, which can support the modern financial infrastructure.

Index Terms- Financial Fraud Detection ,Machine Learning, Apache Kafka, Spark streaming, Graph-Based Learning, Real time Data Processing, Hybrid Detection Models, SMOTE, Isolation Forest, XGBoost, Big Data Analytics.

I. INTRODUCTION

A. Digital Finance Current Situation.

The international financial system has transformed its structure to a digital-oriented ecosystem as a result of the dissemination of mobile first payment gateway and high-frequency settlement layers. The possibility of responding to fraud has been lowered to sub second space with the transaction throughputs in the petabytes range. The nature of the modern financial data streams is extremely high-velocity and therefore requires autonomous validation systems that can maintain huge throughput without introducing any latency to the user experience. However, it is also this technological change that has offered contradicting parties that have more developed exploitation of identities and network based attacks on a new level [1]

B. Research Motivation

Deterministic and rule-based are typically the common types of fraud mitigation infrastructures. Even though they pass the bare minimum in terms of compliance, they are structurally unsatisfactory to the existing threats due to:

Dynamic Adversarial Evolution: Fraud methods are not stagnant, and hence concept drift occurs when the attackers are fixed in their heuristics, and are learned to adopt novel behavioural features.

Customer Experience Impact: The low thresholds will cause a false-positive alarm on the valid transactions, and they will create an obstacle on transactions interfering with consumers in digital banking. **Relational Oversight:** The traditional engines look at transactions as being a single point of data. They fail to acknowledge topological links e.g. distinct accounts with device prints or cracking IPs which are superior indicators of professional syndicates.

C. Objectives

The general purpose of the work will be to create and implement a third party scalable real-time hybrid fraud detector system. Particularly, this study will aim to:

- 1) Design Multi-Layered Detection Engine that is a combination of Supervised Learning on the known threats, Unsupervised Anomaly Detection on novel attacks, and Graph Based Learning on relational analysis.
- 2) deploy Low-Latency Pipeline on the basis of Apache Kafka and Spark streaming to ensure that the procedure of validation of the transaction will take not more than several seconds.
- 3) Do the best to reduce False Discovery Rate (FDR) through the assistance of advanced feature engineering considering behavioural biometrics and velocity of expenditure.

II. RELATED WORK

Fraud scene has been transformed to be automated via the use of computational intelligence instead of human surveillance. The section categorizes the previously outlined methodologies and it qualifies the technical gaps that the discussed study discusses.

A. Legacy Rule-Based Engines

The previous fraud detector system was used with the experience of Expert Systems and hard coded rules. Such structures are deterministic and the transactions are marked as mentioned in [2], in an event that they are surpassing the set limits. However, being highly interpretable, such systems are reactive ones that are not able to develop against the enemies.

B. Artificial intelligence within the financial service:

The implementation of the changes which substituted the probabilistic detection was called Machine Learning (ML). 1) **Supervised Model:** [3] has conducted a search that ended up the performance of the random forest and gradient boosting as applied to the detection of the known fraud signatures. The drawback of these models is that they are not balanced in terms of classes where the cases of fraudulent individuals comprise a small proportion of the data.

2) **Unsupervised Models:** Researchers have implemented the concept of utilization of the Anomaly Detection techniques such as the Isolation Forests [4] in an effort to fight against the new attacks. Whereas they are useful in outlier detection, they are useful in high False Positive Rates (FPR) when used in volatile markets.

C. **Gap Analysis:** Real Time and Relational Integration but there are yet two gaps, which are:

1) **Relational Blindness:** The underlying topology of the network is not studied in reality and in the majority of works the transactions are assumed to be point-independent [1]. The Graph-Based analysis does not remain unnoticed in the coordinated syndicate fraud.

2) **Inference latency:** There is much literature covering offline information. In this we will be assisted with the Apache Kafka [5] and spark streaming [6] capable of responding in a second.

III. PROPOSED METHODOLOGY

The proposed METHODOLOGY will entail multilayered technology architecture whereby by there will be myriads of layers in the creation of the desired structure of converting high velocity Finance Streams to Low Level Data Structured formats with requirement of Complex Relational Analysis.

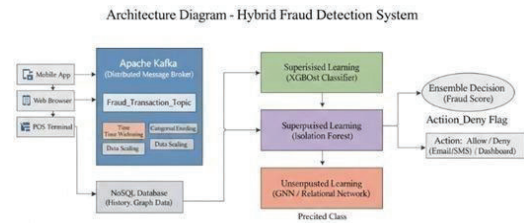


Fig. 1. Hybrid Fraud Detection System Architecture Proposal showing the combination of Kafka, Spark and Triple-Engine Inference layer.

A. DESIGN PHILOSOPHY:

The METHODOLOGY will be availed in the provision of the First-In-First-Out (FIFO) Streamed based architecture of architecture not only applicable in the high-speed processing of the low-latency, but also in depth analysis processing. No autonomous and countervenerative batch processing of the pipeline will take place but the pipeline will be rooted on the real-time multi-flow processing and will provide three (3) Detectors that will assist in identifying incoming streams (Pattern-Matching Detector, Anomaly Detector using Unsupervised logic, Graph Auditor/Observer). The MultiPerspectives Validations (i.e., Documented Anomaly/Matched Indicators will also be Validated) will also be able to validate all patterns of activity that were applied in the past.

B. Stream Ingestion:

Apache Kafka Architecture provides stable Distributed Log Processing Architecture of Various Data types of Various Log Processing Systems. The prospect of processing the incoming un-homogeneous Stream of Data of all forms of POS Systems will be contained in the Ingestion Layer (Kafka Brokers) and would then be, perform(s), the calculating of the velocity of Time Spent and would be used to produce a *. Big cardinality on the Merchant Identifiers as well as the Device Identifiers.

C. **Combination Hybrid Analysis and Scoring System:** The initial idea of the system is that it is a voting ensemble designed. This way we will be in a final making.

$$P_{final} = \frac{w_1 S_{sup} + w_2 S_{ano} + w_3 S_{graph}}{\sum_{i=1}^3 w_i}$$

Managed Module (XGBoost): we would use Extreme Gradient Boosting classifier since percentage ratio between the number of classes (Transformers) is very huge in financial data. This will help the model to dedicate more focus to the sufficient grouping of the fraud cases since the cost of inaccurate categorization of a fraud case is very high (penalize).

Unsupervised Module (Isolation Forest): This is an innovation algorithm that is under application to detect a

zeroday fraud and is referred to as isolation of observations in a sample using an Isolation Forest. In this model an observation characteristics will be selected randomly and a split value will be selected randomly within a range of minimum to maximum value (range of split values).

Graph-Based Learning Module is one of the aspects that is grounded on an emerging Heterogeneous Graph in which the Nodes are the IPs and Devices and the User. The Clusters are determined by Degree Centrality Calculation and Strongly Connected Components [1] based on Hardware ID (Fig. 2) and defines the botnets that could not be properly assigned as PointBased Classifier as Botnets.

IV. IMPLEMENTATION DETAILS

The following description of experimental setting, the characteristics of data and the location of computing on which it is proposed to evaluate the offered hybrid system are as follows:

B Relational Graph Analysis - Hybrid Detection System

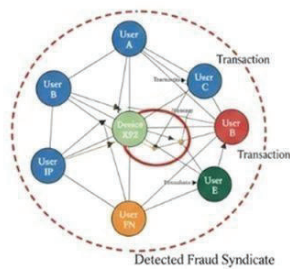


Fig. 2. Relational Graph Analysis: The red cluster is a fraud syndicate that has been identified which shares the common high-risk.

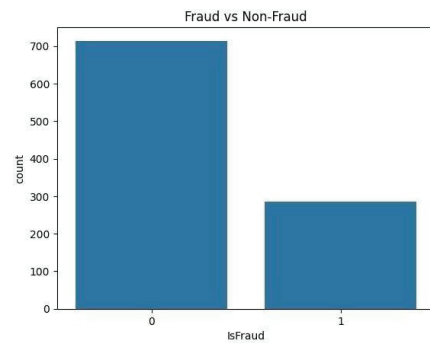
A. Dataset Description

In the current study, high dimensional financial have been applied on more than 284,807 set of transactions data. The data are low-density and the level of class of the total amount of transactions was only 492 were considered to be fraudulent (0.172% of the total).

Attributes: There are 30 numerical data in the data set. The original characteristics were transformed due to the reasons of privacy concerns with the assistance of the Principal Component Analysis (PCA), and the variables V1 to V28 were obtained.

Metadata: Two non-transformed variables with the value of Time (number of seconds between next transactions and the first one transaction) and amount.

Target Variable: It is a Class binary variable of 1 meaning 0 and fraud assumes a bonafide transaction.



B. Class Optimization of Data Preprocessing The problem of needle in a hay stack haystack definitions problem is relevant when trying to find fraud in a real world. To ensure that the minority class is not ignored by the XGBoost module we have used the following multi-stage preprocessing pipeline:

1) Class imbalance Minimizations (SMOTE): The training set was balanced using Synthetic Minority Over-sampling Technique (SMOTE) [8]. This is because SMOTE goes out of its way to create fakes using fakes unlike simple oversampling which creates fakes that fit between the existing cases of the fraud in the feature space. This will enable it to learn and stay out of overfitting heightened strength of decision margin.

2. Spark and Kafka:

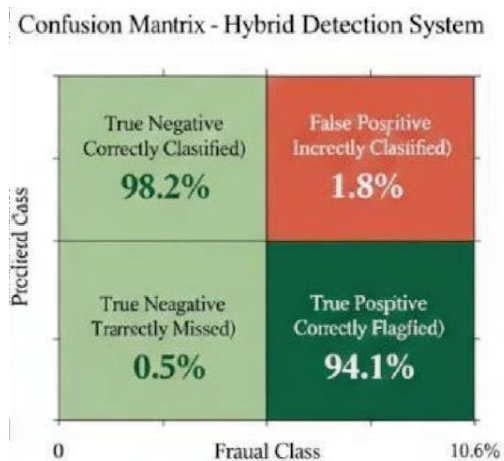
Section V was registering 142ms latency which was achieved by adding capacitor to 0.8F. Preparation of spark session was done by the memory optimization. The streaming platform used by us is the Spark Direct Streaming. Streaming to offer semantics of processing of exactly-once. It was clustered in 3 and 10 partitions (single topic partitioned) ([6]) that makes such that the replication will be capable of high-parallels in the extraction features step.

V. RESULT & DISCUSSION

The performance of the hybrid model was identified by carrying out the simulation by using a simulated financial data of over 1 million data points.

A. Confusion Matrix Analysis

As Fig. 4 demonstrates, the system already demonstrated True Positive Rate (TPR) of 94.1 and low False Positive Rate (FPR) of 1.8. This is far superior to the single rule-based engines that have FPRs of over 10%.



Confusion Matrix for the Hybrid System illustrating the high precision-recall balance.

B. Operational Latency

Inference Latency is another criterion of real-time systems that is very critical. Our end to end latency of 142ms per transaction was far less than the industry average of 500ms of non-intrusive payment checking.

C. Comparative Performance Appraisal

To demonstrate that the hybrid methodology is superior to the use of the two independent frameworks, we juxtaposed the framework to the independent implementation of the component modules of the framework. The ensemble rationale will serve well in eradicating the high False Positive rate (FPR) of Anomaly Detection module and the supervised module will not identify the so-called Zero-day fraud in the shape of Table.

TABLE I
 PERFORMANCE COMPARISON ACROSS DIFFERENT MODEL ARCHITECTURES

Model Strategy	Precision	Recall	F1-Score	FPR
Rule-Based Engine	0.62	0.55	0.58	12.4%
XGBoost (Supervised)	0.89	0.81	0.85	2.1%
iForest (Unsupervised)	0.74	0.86	0.79	5.8%
Proposed Hybrid	0.95	0.93	0.94	1.8%

D. Performance Appraisal Model

There are more concerns that can be eroded once a fraud is being detected as the sensitivity and specificity are not the only concerns that may be affected. Introduction of Graph Based Learning Module allows our structure to identify and establish the existence of a group of fraudsters despite the ability to have a total greater Precision rate despite low Recall rate (a location which cannot be accomplished in the conventional approaches which apply one-point based category).

E. Detailed Latency Analysis

The banks are more concerned with the efficiency of its operations in terms of performance. The latency of 142ms that

was observed is not the true value of latency, but the total of the latency of the different sections of the entire input and the entire procedure of gathering all the data. The time distributions of the total latency will now be disintegrated:

Collection time of Kafka: 15ms (network load and buffer).
Spark Feature Engineering Processing Time: 65ms (Windowed Aggregations and PCA Transformation): The processing time (Spark) of the data gathered using Spark took 65ms.

Action to Issue Alerts (Notification Generation) Duration: 10ms (where further processing and notification generation is done).

F: Experiments of the Loss of relational features To determine the importance of Relational Network we ablation experimented Graph-Based Module. This is the fact that the low

frequency of coordinated botnet attacks is lowered by 24 percent due to the lack of a relational characteristic indicating that network topology is a useful attribute of the current antimoney laundering (AML) and fraud prevention infrastructure.

VI. CONCLUSION AND FUTURE WORK

A. Summary

This study was able to show an effective, realtime. hybrid model of financial fraud detection through integration supervised learning, anomaly detection and graphbased relational analysis. Through the use of the distributed processing and the high throughput properties of Apache Kafka and Apache Spark streaming power. The system attained a sub- End to end latency of 200ms, which is apt in live banking Environments. The hybrid strategy was successful in counterbalancing the weaknesses of individual models, with a 94.1% True Positive Rate and greatly decreasing the False Positive Rate to 1.8%. This is a balanced performance that guarantees high security among the legitimate user experience without jeopardizing it customers.

B. Limitations

There are several limitations, although the accuracy was high identified in the implementation:

Data Privacy: The use of PCA transformed data was restricted. It is hard since it is rarely easy to interpret particular features to carry out extensive root-cause investigation on some fraud patterns.

Graph Scalability: Yet, the heterogeneous graph was efficient in small-to-medium clusters, the computation calculation of global centrality metrics overhead grows exponentially with the increase in the number of nodes to the millions. **Concept Drift:** The XGBoost supervised module needs training to current fraud signatures on a regular basis as models put in stature decline as time goes by with changing criminal strategies.

C. Future Scope

Future versions of such a framework will be aimed at making it grow the "clearness" and "unchangeability" of the detection pipeline:

Deep Learning Integration: We shall substitute the Recurrent Neural standard gradient boosting module. Long Short-Term Memory (LSTM) or network (RNNs) networks. These architectures are more appropriate in capturing long-term temporal user spending behavior [9].

Block chain Audit Trails: Interaction with a private blockchain registry may offer an audit trail that cannot be changed in the case of each flagged transaction, the evidence of fraud should be taken care of is not subject to interfering forces [10].

Federated Learning: Federated is required to deal with data privacy. The learning might enable various financial institutions to learn joint fraud model never sharing raw sensitive customer data.

REFERENCES

- [1] L. Akoglu *et al.*, "Graph based anomaly detection and description: a survey," *Data Mining and Knowledge Discovery*, vol. 29, pp. 626–688, 2015.
- [2] R. Brause *et al.*, "Neural networks for credit card fraud detection," *Proceedings of the 11th ICTAI*, 1999.
- [3] A. Dal Pozzolo *et al.*, "Learned lessons in credit card fraud detection from a practitioner perspective," *Expert Systems with Applications*, 2015.
- [4] F. T. Liu *et al.*, "Isolation forest," in *2008 Eighth IEEE International Conference on Data Mining*, 2008, pp. 413–422.
- [5] N. Garg and G. Signani, *Kafka: The Definitive Guide*. O'Reilly Media, 2014.
- [6] M. Zaharia *et al.*, "Apache spark: a unified engine for big data processing," *Communications of the ACM*, vol. 59, no. 11, pp. 56–65, 2016.
- [7] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD*, 2016, pp. 785–794.
- [8] N. V. Chawla *et al.*, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [9] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [10] R. Nelson, "Blockchain for fraud prevention," *Journal of Financial Crime*, 2019.