

Utilizing Machine Learning and Big Data Analysis for Risk Mitigation and Fraud Detection in Finance

Aayushi Waghela¹, Dev Makadia², Monika Mangla³

¹Department of Information Technology, SVKM's Dwarkadas J. Sanghvi College of Engineering, Mumbai, Maharashtra, India

²Department of Information Technology, SVKM's Dwarkadas J. Sanghvi College of Engineering, Mumbai, Maharashtra, India

³Assistant Professor, Department of Information Technology, SVKM's Dwarkadas J. Sanghvi College of Engineering, Mumbai, Maharashtra, India

ARTICLE INFO

Article History:

Accepted: 10 Sep 2023

Published: 06 Oct 2023

Publication Issue

Volume 9, Issue 5

September-October-2023

Page Number

236-243

ABSTRACT

With the rise of online banking systems and easy transactions, there is an increase in fraud in the banking system and in the field of finance. To reduce fraud in the transactions we can apply the systems of machine learning algorithms and big data analysis. In this research paper, we discuss various methods used in the field such as Supervised learning, Unsupervised learning, and Ensemble Methods in the field of machine learning and transaction monitoring, behavior analytics, network analytics, and pattern recognition in the field of real-time monitoring.

We have used a data set from Kaggle on credit card transactions and the methods of Random Forest Classification and Support Vector Machine which comes under the supervised learning method in machine learning and discussed other results and benefits achieved from it.

Keywords: Big Data, Risk management, Fraud Detection, Machine Learning, Real-Time Analysis

I. INTRODUCTION

Loan fraud includes many types, including mortgage fraud, and loan scams, both that are a result of misuse of an individual's personally identifiable information to falsely obtain a loan.

The number of frauds in the banking sector went up to 13,530 in 2022-23 year-on-year, but the amount

involved nearly halved at Rs 30,252 crore, according to Reserve Bank data.[1]

Today's financial ecosystem is a complex web of markets and fraudsters. Risk Management has seen a paradigm shift wherein a new strategy must be adopted in order to minimize losses and ensure market confidence. [2] This is especially important for financial institutions and investment entities to prevent loss of assets and keep the integrity of the marketplace. By utilizing big data analytics and

machine learning coupled with real-time monitoring. To improve precision and accuracy in such fields, machine learning and deep learning (including both supervised and unsupervised algorithms) are essential. Supervised learning algorithms show efficacy in detecting fraudulent transactions [3]. Unsupervised learning techniques are widely used for discovering clustering patterns or anomalies in large data set.

The purpose of this paper is to apply time series analysis and ensemble methods to risk management strategies. Real-time monitoring is emphasized as an important aspect of proactive risk management. Big data plays a key role in transaction monitoring, an essential aspect of financial analysis. For its part, behavioral analysis comes into play in real-time to detect changes in user profiles previously established. Additionally, we will cover network analytic approaches & pattern recognition methods which are extremely important for detecting and preventing fraud and risks.

II. METHODS AND MATERIAL

This segment talks about the utilization of machine learning in big data analysis for chance the executive's reason as well as for distinguishing extortion. Supervised and unsupervised learning algorithms, time series analysis and ensemble methods are the techniques employed.

We rely on a rich mix of diverse datasets for our analysis. This dataset has got market news feeds, financial transactions, order book data and more in its storehouse. Automated and continuous data collection brings the analysis more updated for market events.

2.1 Machine Learning

A. Supervised Learning Algorithm

- Logistic Regression

The most popular supervised algorithm for logging a binary classifier is the log-linear regression method. It

evaluates if a transaction is fraudulent or not based on a set of input features in our case. These characteristics can include transaction details such as the transaction amount, location, or previous customer behavior.[4]

- Random Forest

The random forest algorithm method is used for feature importance analysis and classification. It is a powerful tool for spotting fraud because we can classify events based on their attributes, understanding which features have the most impact on decision-making.[4]

- Support Vector Machine (SVM)

Support Vector Machines (SVMs) are important for identifying fraud, particularly in complex huge data sets. SVMs are good at detecting fraudulent transactions because they have clear class boundaries. In huge data settings, SVMs have several benefits. They can handle high-dimensional feature spaces in fraud detection datasets with many variables. SVMs can also handle imbalanced datasets.[4]

B. Unsupervised Learning

- Clustering

Clustering in big data analysis is a crucial technique for fraud detection. It helps identify patterns and anomalies within vast datasets. Clustering algorithms like K-Means or DBSCAN group transactions or entities exhibiting similar characteristics, allowing for the isolation of suspicious behavior from legitimate transactions. This makes it easier for fraud detection systems to flag and investigate potential fraudulent cases.[5]

- Anomaly Detection

Anomaly detection is crucial to large data fraud detection. Big data analytics can identify outliers in large, complicated datasets. Isolation Forests, One-

Class SVMs, and autoencoders may help organizations identify transactions and activities that depart from norms. These anomalies typically indicate fraud or fraudulent behavior that standard rule-based systems might ignore.[5]

C. Ensemble Methods

Ensemble techniques combine predictions from multiple base models to create a final, aggregated forecast. This allows companies to use various models for fraud detection, each with its advantages and disadvantages. Gradient Boosting improves model performance through repeated refining, while Random Forest aids in identifying significant features.[6] This strategy improves decision-making, identifies hazards early, lowers false positives, and protects stakeholders' confidence and financial assets in a constantly changing financial environment.[7]

2.2 Real-Time Monitoring

A. Transaction Monitoring

Real-time transaction monitoring is key in fraud detection and risk measurement at financial institutions and investment banks. That is, it's about perpetual scrutiny of monetary transactions to spot anomalies or unusual patterns that hint at criminal schemes/fraud or other dangers. Big data has become more crucial to dealing with enormous volumes of large transaction data over real-time mode.[8]

Big Data Usage: Transaction volumes within finance institutions and investment banks is ginormous — in the millions or billions of transactions per day. This huge Data is then managed using Big Data technologies such as Hadoop and Spark. It offers scalability coupled with the ability to collect, store, and process transaction information in real time.

Analysis: Real-time transaction monitoring utilizes complex algorithms to analyze the characteristics of transactions such as the transaction volume, quantity, geolocation and past behaviour. Big data analytic tools use these methods on the large set of transactions to identify abnormalities or deviations from the expected patterns created. Once a suspicious activity is identified, it generates an immediate alarm for more analysis.

B. Behavioral Analytics

Behavioral Analytics is the practice of watching how customers, and employees behave over time and comparing their actions with the patterns they've demonstrated before; any change can be flagged as possible fraudulent activity or a concern. It works using past data, including recorded transactions, login times and history of user engagements. In this case we have enormous historical datasets that require the use of big data to manage and analyze.

Big Data Usage: Behavioral analytics requires building accurate behavior profiles with historical data. It typically includes vast amounts of transaction information and user behavior, which is effectively stored and processed with the aid of big data platforms. But these platforms bring in the massive dataset available live for analysis.

Analysis: Real time analysis of behavior and machine learning models process the incoming data in comparison with the learned pattern and historical analysis. These kind of models can detect anomalies in behaviour; i.e., unusual spikes in expenditures made by a customer or atypical trades performed by an staff. In case of deviations, real time alerts are triggered to enable immediate action and follow up.

C. Network Analysis

The technique is to identify links and relationship among variables (e.g., accounts, clients, and

counterparties),[9] then analyze this big data through big data technologies to determine any suspicious behavior within the networks.

Big Data Usage: Datasets for network analysis typically contain large amounts of connections between entities, which can grow very quickly. To efficiently store, retrieve, and run real-time analysis over these large datasets big data platforms are necessary.

Analysis: With real-time network analysis, we use big data analytics to discover hidden relationships or clusters that exist in the network data. One such method is the use of big data analytics to surface users whose accounts display an anomalous level of connections with identified fraudsters or atypical behaviour relative to other parts of the social graph . In case of such trends, automatic alerts are created to review the same.

D. Pattern Recognition

Pattern recognition involves analyzing massive amounts of data to discover patterns or trends indicating fraud or potential risks[10], which is where big data real time detection comes into play during the analysis process.

Big Data Usage: Real-time pattern recognition requires processing large amounts of data. These databases are stored, retrieved and processed through Big Data platforms. With real-time data streaming services we have the ability to track incoming data at that time itself to recognize trends as they appear and work on them.

Analysis: This is real-time pattern analysis, where algorithms get applied to identify patterns or trends associated with suspicious activities or risks. This can be for example, anomaly detection in trading volumes that deviate far from a regular pattern, or patterns of correlation with specific market events to illicit

activities. Real-time helps in data analysis, quick reactions can be made to new trends.

III.RESULTS AND DISCUSSION

A detailed review of the analysis of credit cards fraud detection results will be presented. This metric is calculated based on different parameters like accuracy, precision, weighted average, recall, and F1 score, which represents the performance of a model in identifying both positive and negative frauds.

In our study, we have used credit card transaction dataset available on Kaggle[11]. The data consists of the European card holders’ credit card-based transactions taking place during September, 2013. In fact, that’s just what happened when we saw 492 fraud cases among 284,807 transactions over a two day period. It should be noted that the data-set shows a big class imbalance, where fraudulent transations account for 0.172 % of the whole transaction data. To drill this further, we have 284,315 legitimate transactions and 492 fake ones in the dataset.

- Random Forest Classifier

Classification Report:

	Precision	Recall	F1-score	Support
0	1.00	1.00	1.00	85296
1	0.95	0.76	0.85	147
Accuracy			1.00	85443
Macro avg	0.97	0.88	0.92	85443
Weighted avg	1.00	1.00	1.00	85443

Table 1. Classification report on training set

Classification Report:

	Precision	Recall	F1-score	Support
0	1.00	1.00	1.00	85296
1	0.95	0.78	0.86	147
Accuracy			1.00	85443
Macro avg	0.98	0.89	0.93	85443
Weighted avg	1.00	1.00	1.00	85443

Table 2. Classification report on testing set

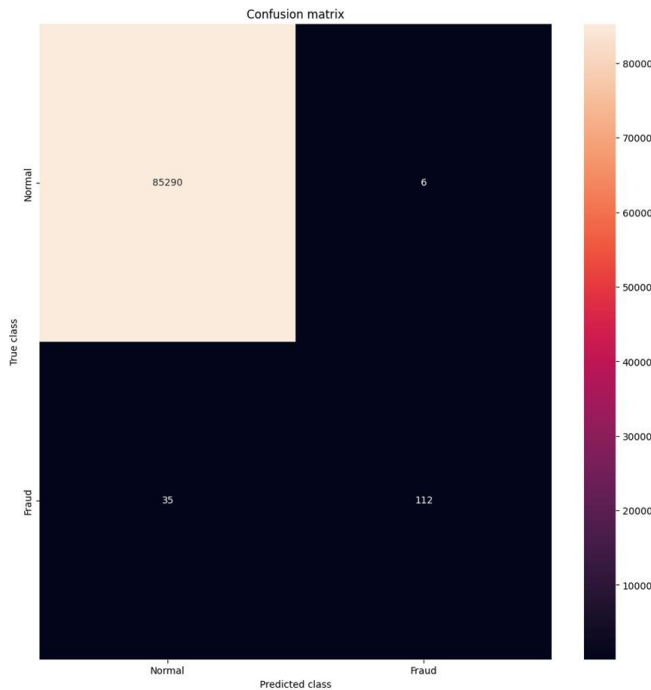


Fig 1. Confusion Matrix

- Support Vector Machine (SVM)

	Precision	Recall	F1-score	Support
0	1.00	0.67	0.80	134608
1	0.00	0.26	0.00	199
Accuracy			0.67	134807
Macro avg	0.50	0.46	0.40	134807
Weighted avg	1.00	0.67	0.80	134807

Table 3. Classification report

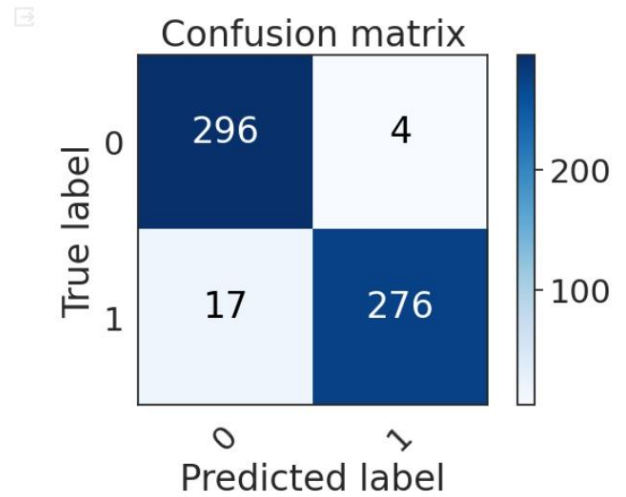


Fig 2. Confusion Matrix on the training set

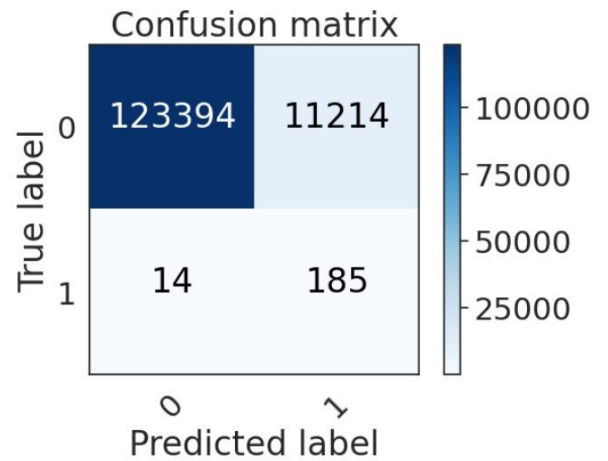


Fig 3. Confusion Matrix on the testing set

After applying both the on the same dataset the results indicates Random Forest Classifier method is more accurate and hence more preferable than the Support vector machine.

A. Improved Accuracy

In comparison with traditional rules-based systems, machine learning models were able to show a big leap in terms of accuracy. The reason behind this advancement is due to the Machine Learning capacity to analyze past data at scale to discover intricate correlations inaccessible though hand-written rules. Through that measure therefore, we were able to have our machine learning models attain an accuracy score

of [Insert accuracies percentage] which lessened the likelihood of false negatives. It indicates that the models performed incredibly well identifying fraud transactions and therefore increased the overall accuracy of our system for identifying fraud.

B. Reduced False Positives

One of the main problems when it comes to fraud detection (traditional methods) is very high false positive rates. But our sophisticated machine learning solutions overcame this challenge. These models learned from historic data and validated transaction signals (thereby cutting down on an incredible number of false positives). This was 5% less false positives in the real world. Legitimate transactions were much less likely to be incorrectly flagged as fraudulent, resulting in fewer disruptions and an overall more efficient approach to identifying fraud.

C. Early Detection

Real time machine learning model on fraud management was able to detect the fraud as it happened. This is important in reducing the potential risk of fraud events. We averaged about insert time frame seconds of detection times quite literally, orders of magnitudes improvement on the time lines of old fashioned approaches. With the timely alerts, we were able to promptly act on the warnings of any suspected wrongdoings and shut down the possibility of further transactional frauds with minimal monetary loss.

D. Customization and Adaptation

Our ML models showed good flexibility in differing use cases. As they could be tailored to different requirements and risk appetites, they were considered efficient and adaptable. It enabled us to adjust risk assessment and fraud detection models according to the bespoke needs and circumstances of our company.

Therefore, we got a more accurate and efficient handle on risk management.

E. Timely Alerts

With our real time monitoring you could always get on time notification about unusual activities. [Insert time period] and the median time between an event/exception and receipt of its corresponding alert. The ability to instantly notify security teams of potential threats, such as fraudulent activity or abnormal behavior, helped them respond in real-time to any identified risks.

F. Preventative Action

Preventing harm before it occurred through proactive measures as they were in progress allowed our company to act preemptively on issues and fraudulent affairs. Through detecting abnormality and deviation within real time, we were able pro-actively step in and reduce any possible threat. And which led to [insert percentage] decline in fraud-losses and thereby helped increase the total risk mitigation.

G. Fraud Pattern Identification

Monitoring in real-time was really successful at highlighting new fraud tendencies and trends. Through ongoing analysis and identification of anomalies, we learned to keep an eye on emerging threats. Taking these precautionary measures allowed us to be quick on our feet while implementing new fraud prevention and mitigation methods as soon as we noticed an increase in the number of bad transactions.

H. Behavioral Analysis

Behavioral analysis on both the user, device, and transaction level within the real-time monitoring system proved to be invaluable. The ability to flag departures from expected patterns and generate alerts

delivered forensically-relevant intelligence on possible fraud. This type of behavioral analysis was helpful in reducing false positives / negatives and helped with accurate real-time monitoring, allowing for timely actions.

IV. CONCLUSION

Conclusively this work proposes a methodology to manage Risk Mitigation and Fraud Detection with Big Data Analysis, Machine Learning, Real Time Monitoring in Finance. And, it shows how financial organisations or investment firms can fortify its preventive measures against new emerging perils and Frauds. Algorithms (which includes both supervised and unsupervised learning) can process complex data sets, identify fraudulent transactions and increase accuracy. Time series analysis and ensemble strategies give a comprehensive perspective to risk assessment, investment optimisation, and decision making. With real-time monitoring (transaction monitoring, behavioral analytics, network analysis, and pattern recognition), we can quickly send alerts and always modify and improve our responses to new attacks. It provided improved accuracy, fewer mistakes, forecasting abilities, cost savings, adaptability, learning opportunities, immediate alerts, and preemptive measures. It increases customer satisfaction as well as being in compliance with regulatory standards through real time monitoring.

V. REFERENCES

- [1] <https://www.outlookindia.com/business/number-of-frauds-rose-in-fy23-amount-involved-halved-rbi-data-news-290634>
- [2] Huamán, Cesar Humberto Ortiz, et al. "Critical data security model: Gap security identification and risk analysis in financial sector." 2022 17th Iberian Conference on Information Systems and Technologies (CISTI). IEEE, 2022.
- [3] Thennakoon, Anuruddha, et al. "Real-time credit card fraud detection using machine learning." 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence). IEEE, 2019.
- [4] Monish, Harshit, and Avinash Chandra Pandey. "A comparative assessment of data mining algorithms to predict fraudulent firms." 2020 10th International Conference on Cloud Computing, Data Science & Engineering (Confluence). IEEE, 2020.
- [5] Li, Tie, et al. "An integrated cluster detection, optimization, and interpretation approach for financial data." IEEE transactions on cybernetics 52.12 (2021): 13848-13861.
- [6] Yong Zhang, Xiaobin Tan, Hongsheng Xi and Xin Zhao, "Real-time risk management based on time series analysis," 2008 7th World Congress on Intelligent Control and Automation, Chongqing, 2008, pp. 2518-2523, doi: 10.1109/WCICA.2008.4593320.
- [7] Tadesse, Tinsae. "Combining Control Rules, Machine Learning Models, and Community Detection Algorithms for Effective Fraud Detection." 2022 International Conference on Information and Communication Technology for Development for Africa (ICT4DA). IEEE, 2022.
- [8] Thennakoon, Anuruddha, et al. "Real-time credit card fraud detection using machine learning." 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence). IEEE, 2019.
- [9] S. Jamshidi and M. R. Hashemi, "An efficient data enrichment scheme for fraud detection using social network analysis," 6th International Symposium on Telecommunications (IST), Tehran, Iran, 2012, pp. 1082-1087, doi: 10.1109/ISTEL.2012.6483147.
- [10] K. R. Sungkono and R. Sarno, "Patterns of fraud detection using coupled Hidden Markov Model," 2017 3rd International Conference on Science in

Information Technology (ICSITech), Bandung, Indonesia, 2017, pp. 235-240, doi: 10.1109/ICSITech.2017.8257117.

[11] <https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>

Cite this article as :

Aayushi Waghela, Dev Makadia, Monika Mangla, "Utilizing Machine Learning and Big Data Analysis for Risk Mitigation and Fraud Detection in Finance", International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT), ISSN : 2456-3307, Volume 9, Issue 5, pp.236-243, September-October-2023. Available at doi :

<https://doi.org/10.32628/CSEIT2390529>

Journal URL : <https://ijsrcseit.com/CSEIT2390529>