

ISSN(O): 2320-9801 ISSN(P): 2320-9798



International Journal of Innovative Research in Computer and Communication Engineering

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.771

Volume 13, Issue 4, April 2025

⊕ www.ijircce.com 🖂 ijircce@gmail.com 🖄 +91-9940572462 🕓 +91 63819 07438

www.ijircce.com



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Big Data Driven Credit Card Fraud Detection in Banking

Shaik Fara¹, V Keerthi², K Vishwesh³, J Nikhitha⁴ B.M.Rajasekhar⁵

UG Students, Dept. of CSE-DS, SRK Institute of Technology, Enikepadu, Vijayawada, Andhra Pradesh, India¹⁻⁴

Assistant Professor, Dept. of CSE-DS, SRK Institute of Technology, Enikepadu, Vijayawada, Andhra Pradesh, India⁵

ABSTRACT: Banks face challenges in detecting fraud due to the high volume and complexity of transactions. In an era where digital transactions dominate the global economy, fraud detection has become a cornerstone of financial security. Traditional fraud detection systems, which rely heavily on rule-based methodologies, are increasingly being outpaced by the sophisticated techniques employed by modern fraudsters. These legacy systems struggle with adapting to the fast evolving landscape of digital fraud, often producing a high number of false positives and suffering from delayed detection. As financial transactions increase in both volume and complexity, the demand for more agile, accurate, and real-time fraud detection systems is paramount.

KEYWORDS: Credit Card Fraud Detection, Big Data, Fraud Detection, Machine Learning, Real-Time Analytics, Anomaly Detection, Supervised Learning, Unsupervised Learning, Banking Security

I. INTRODUCTION

The rise of digital transactions, mobile banking, and e-commerce has led to a significant increase in fraud-related activities. Financial institutions are under pressure to detect fraudulent transactions in real time to protect customers and prevent financial losses. Traditional fraud detection systems, which rely on rule-based approaches, are increasingly inadequate for handling complex and adaptive fraud patterns. These systems are often too rigid to adapt to new types of fraud, resulting in delayed detection and high rates of false positives, which can disrupt legitimate customer transactions.

This project aims to develop a big data-driven fraud detection system using advanced analytics and machine learning. The system will analyze transaction data, detect anomalies, and flag suspicious activities in real-time . By leveraging big data, it can process vast amounts of information efficiently, identifying fraud patterns that traditional systems might miss. Both supervised and unsupervised learning methods will be used—supervised models will learn from past fraud cases, while unsupervised methods will detect new fraud patterns. Real-time alerts will notify the bank's fraud team, allowing quick intervention to minimize losses. The goal is to improve fraud detection accuracy, reduce false positives, and enhance banking security.

The paper is structured as follows: Section II discusses related work, highlighting previous studies in ML-based fraud transaction. Section III provides a detailed background on the algorithms used in the project. Section IV introduces the proposed system, detailing the methodology and model architecture. Section V presents comparative results using graphical visualizations, and Section VI concludes the study with insights and future research directions.

II. RELATED WORKS

Big Data-driven credit card fraud detection has seen significant advancements through machine learning and deep learning techniques. Researchers utilize algorithms like Random Forest, SVM, and XGBoost to identify fraudulent patterns within large transaction datasets, while deep learning models such as RNNs and CNNs analyze time-series and user behavior data. Feature selection focuses on variables like transaction amount, location, frequency, and user spending habits. Clustering methods (K-Means, DBSCAN) help group similar fraudulent transactions, and NLP techniques analyze transaction descriptions and user reviews for anomaly detection. Explainable AI (XAI) enhances model transparency, aiding fraud analysts in understanding and validating predictions. Recent improvements include real-time monitoring of transaction streams with distributed processing frameworks and adaptive learning models that

www.ijircce.com



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

dynamically adjust to evolving fraud patterns. However, challenges persist in handling imbalanced datasets and ensuring model generalizability across diverse user populations and banking systems.

Paper Information Description		Limitations/Inference		
Li, X., Chen, Y., & Zhang, Z. (2022) [1]	Developed a real-time fraud detection system using a Random Forest model on a large dataset of transaction logs and user profiles.	Achieved high accuracy (95%), but scalability for extremely high transaction volumes needs further optimization.		
Kumar, S., Patel, A., & Sharma, R. (2023) [2]	Proposed a Deep Learning approach utilizing Recurrent Neural Networks (RNNs) to capture sequential patterns in transaction data for fraud detection	Showed improved performance in detecting complex fraud patterns compared to traditional methods. However, interpretability of the model remains a challenge.		
Nguyen, T., Tran, H., & Le, D. (2021) [3]	Implemented a Hybrid Model combining Isolation Forest for anomaly detection and XGBoost for classification, leveraging features extracted from social media data.	Demonstrated robust performance with high precision, but data privacy concerns related to social media integration require careful consideration		
Garcia, M., Lopez, J., & Rodriguez, P. (2024) [4]	Explored the use of Graph Neural Networks (GNNs) to model relationships between accounts and transactions for fraud detection.	Achieved superior detection of collusive fraud rings. However, computational complexity and trainin		
Kim, J., Park, S., & Choi, H. (2023) [5]	Applied a Federated Learning approach to train a fraud detection model across multiple banks without sharing raw transaction data.	Addressed data privacy concerns and showed promising results in terms of model accuracy. Further research is needed to improve communication efficiency and handle data heterogeneity.		
Wang, L., Liu, Q., & Wu, J. (2022) [6]	Introduced a Stream Processing framework using Apache Flink to enable real-time fraud detection on streaming transaction data	Demonstrated low latency and high throughput for real-time detection. However, model drift and concept drift in dynamic fraud patterns need to be addressed.		

III. BACKGROUND

- 1. Machine Learning Models
- 1.1Random Forest



Fig1: Random Forest

Figure 1This architecture illustrates a standard machine learning workflow for a predictive task, starting with dataset selection. The raw data undergoes pre-processing, including feature selection, data filtering, and cleaning to enhance its quality. Subsequently, the data is split into training and testing samples, often employing techniques like SMOTE for



imbalanced datasets. Machine learning models are then instantiated and trained using the training data. The trained model's performance is evaluated on both the training and testing samples to assess its generalization ability. Finally, the results are analyzed and compared to draw conclusions about the model's effectiveness and identify areas for improvement. This iterative process aims to build a robust and accurate predictive model for the given task.

1.2 Decision Tree



Fig2: Decision Tree

Figure 2 demonstrates a decision tree, a supervised learning algorithm used for classification. The tree is structured with nodes representing decision points based on specific features, and branches indicating the outcomes of those decisions. Each internal node tests an attribute (feature), and each branch represents an outcome of the test. The leaves (green nodes) represent class labels or decisions. The blue lines indicate the "True" branch, meaning the condition is met, while the red dashed lines represent the "False" branch, where the condition is not met. The conditions are typically comparisons, such as "Median ≥ 0.07 ?" or "Third person plural pronouns ≥ 2.0 ?". The tree progresses from the root (top) to the leaves (bottom), classifying data by following the path based on the data's attributes. The "v:1" and "v:0" at the leaves likely represent the final classification results, such as "1" for positive and "0" for negative. This tree illustrates how decisions are made sequentially based on features to arrive at a final classification.

2. DEEP LEARNING MODELS

2.1 LSTM (Long Short Term-Memory)





Fig 3 This model illustrates a machine learning pipeline specifically designed for credit card fraud detection. The process begins with a Credit Card Fraud (CCF) Training Dataset, which is then subjected to Data Pre-processing to clean and prepare the data for analysis. Notably, GA Feature Selection (Genetic Algorithm Feature Selection) is employed to identify the most relevant features for fraud detection, enhancing model performance by reducing dimensionality and noise. The pre-processed data is split into Training Data and a separate Test Subset. The Training Data is used to train various machine learning models, including Decision Trees (DT), Random Forests (RF), Logistic Regression (LR), Naive Bayes (NB), and Artificial Neural Networks (ANN), resulting in a Trained Model. The Test Subset also undergoes Data Pre-processing and is used as Test Data to evaluate the model's performance. Finally, the Trained Model is applied to detect fraud, and the results are assessed to determine its effectiveness in identifying fraudulent transactions. The "Train with vn" and "Train the models" annotations indicate the training phase and the use of a variable "vn" during the training process.

2.2 CNN



Fig-4: CNN

Fig 4 This diagram illustrates a Convolutional Neural Network (CNN) architecture designed for image classification, specifically identifying a "car" in the input image. The process begins with the input image being fed into a convolutional layer, which extracts features using filters. This is followed by a pooling layer that reduces the spatial dimensions of the feature maps, decreasing computational complexity while retaining important information. Another convolutional layer further extracts higher-level features, followed by another pooling layer for dimension reduction. Finally, a fully-connected layer combines all the learned features to make a classification decision, resulting in the predicted image label, "car." The CNN's hierarchical structure, with alternating convolutional and pooling layers, allows it to learn increasingly complex features, enabling accurate image recognition. This architecture is fundamental in various computer vision applications, showcasing the power of CNNs in image analysis

3. Dataset

4	А	В	С	D	E	F
1 ind	lex	trans_date_trans_time	merchant	category	amt	city
2	(0 2019-01-01 00:00:44	Heller, Gutmann and Zieme	grocery_pos	107.23	Orient
3		1 2019-01-01 00:00:51	Lind-Buckridge	entertainment	220.11	Malad City
4	1	2 2019-01-01 00:07:27	Kiehn Inc	grocery_pos	96.29	Grenada
5		3 2019-01-01 00:09:03	Beier-Hyatt	shopping_pos	7.77	High Rolls Mountain Park
6	4	4 2019-01-01 00:21:32	Bruen-Yost	misc_pos	6.85	Freedom
7	5	5 2019-01-01 00:22:06	Kunze Inc	grocery_pos	90.22	Honokaa
8	(5 2019-01-01 00:22:18	Nitzsche, Kessler and Wol	shopping_pos	4.02	Valentine
9		7 2019-01-01 00:22:36	Kihn, Abernathy and Douglas	shopping_net	3.66	Westfir
10	8	8 2019-01-01 00:31:51	Ledner-Pfannerstill	gas_transport	102.13	Thompson
11	9	9 2019-01-01 00:34:10	Stracke-Lemke	grocery_pos	83.07	Conway
12	10	0 2019-01-01 00:40:50	Cummerata-Jones	gas_transport	70.53	Athena
13	1:	1 2019-01-01 00:41:45	Huel-Langworth	misc_net	177.57	Thompson
14	12	2 2019-01-01 00:46:18	Ferry, Lynch and Kautze	misc_net	2.76	San Jose
15	13	3 2019-01-01 00:49:25	Little, Gutmann and Lynch	shopping_net	83.52	Ravenna
16	14	4 2019-01-01 00:56:12	Swaniawski, Lowe and Robel	shopping_pos	317.14	Parks
17	15	5 2019-01-01 00:56:59	Reichert, Huels and Hoppe	shopping_net	113.4	Fort Washakie
18	16	5 2019-01-01 01:00:48	Howe Lt	misc_pos	218.71	Littleton
19	17	7 2019-01-01 01:02:16	Wolf Inc	grocery_pos	89.11	Meadville
20	18	8 2019-01-01 01:04:48	Vandervort-Funk	grocery pos	50.68	Moab

Table 1: presents a dataset



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

This dataset appears to represent credit card transaction data, likely used for fraud detection or customer behavior analysis. It's structured with six columns: "index," "trans_date_time," "merchant," "category," "amt," and "city." The "index" column serves as a unique identifier for each transaction. The "trans_date_time" column records the precise timestamp of each transaction, down to the second, indicating a high level of temporal granularity. The "merchant" column lists the names of various businesses where the transactions occurred, suggesting a diverse range of purchasing activities. The "category" column classifies the transactions into categories like "grocery_pos," "entertainment," "shopping_pos," "misc_net," and "gas_transport," providing insight into the nature of the purchases. The "amt" column specifies the transaction amount, showing the monetary value of each transaction. Finally, the "city" column indicates the geographical location where the transaction took place. This dataset could be used to analyze spending patterns, identify anomalies, or train machine learning models to detect fraudulent transactions based on the time, location, merchant, and amount of purchases. The precise timestamps and diverse categories suggest a focus on detailed transactional behavior for potential fraud analysis or customer profiling.

IV. PROPOSED SYSTEM

1. Architechture

Fig. 5 This architecture outlines a multi-phased approach to credit card fraud detection, starting with Phase 1, the User Interface. In this phase, transaction data such as "Txion Date Time," "Payment Amount (\$)," and "Credit Card Details" are captured. This initial data collection forms the basis for subsequent analysis. The data flows into the ML Match Data processing block, where various feature engineering and matching operations are performed. These include extracting the "P_R Email domain," counting IP address masses ("CL_Cl4"), analyzing time differences between previous transactions ("DL_Dl5"), matching card name and address ("Ml_Ma"), identifying entity relations ("Vxx"), and categorizing transactions ("Category"). These processed features are then fed into an Ensemble ML Alg composed of Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Long Short-Term Memory networks (LSTM).

The outputs from the ensemble of machine learning algorithms are used to make a final determination: whether the transaction is "Fraud detected" or "Not detected." The "Accuracy" block indicates that the system's performance is evaluated to ensure the reliability of its fraud detection capabilities. The use of an ensemble model, combining the strengths of CNNs, RNNs, and LSTMs, suggests a sophisticated approach aimed at capturing both spatial and temporal patterns within the transaction data. This architecture highlights a comprehensive, multi-layered strategy for enhancing the precision of credit card fraud detection by leveraging detailed transaction features and advanced machine learning techniques.



ML Match Data

Fig5 Overall Architecture

This raw data is then processed to extract and match critical features, including email domain analysis, IP address counting, time-based transaction comparisons, address matching, entity relationship identification, and category checks.



These engineered features are subsequently fed into an ensemble of machine learning models, combining CNNs, RNNs, and LSTMs, to leverage their respective strengths in pattern recognition. The system's output is a binary classification, indicating whether a transaction is fraudulent or legitimate, with an accuracy assessment to ensure reliability. This architecture demonstrates a comprehensive strategy to enhance fraud detection precision by integrating detailed feature engineering and advanced machine learning techniques.







2. Workflow

Fig. 6 illustrates the workflow This workflow diagram outlines a typical machine learning process for fraud detection, likely in a financial or transactional context. It begins with two distinct datasets: Training Data and Test Data. The Training Data, representing historical transactions, undergoes Pre-processing, which involves cleaning, transforming, and feature engineering to prepare it for model training. This prepared data is then fed into various Machine Learning Algorithms, including Convolutional Neural Networks (CNN), Long Short-Term Memory networks (LSTM), and Recurrent Neural Networks (RNN), among others. These algorithms are trained on the pre-processed data to learn patterns indicative of fraudulent activity. The Test Data, representing unseen transactions, is used to evaluate the model's performance. The Result from the trained model, which is a prediction of whether a transaction is fraudulent or not, is then passed through a decision point: IS FRAUD? If the result indicates fraud (YES), an ALERT!! is triggered, prompting further investigation or action. If the result to the model training stage suggests an iterative process where the model is continuously refined based on its performance. This workflow highlights a standard machine learning approach to fraud detection, emphasizing data preparation, model training and evaluation, and decision-making based on model predictions.

V RESULTS AND DISCUSSIONS'

1. Comparitive Analytic table

	Model Name	Precesion	F1_Score
1	ML Algorithm	0.70	0.77
2	Random Forest	0.82	0.94
3	Decision Tree	0.88	0.85
4	Neural Network	0.85	0.92
5	LSTM	0.90	0.96

Table2: Different model evaluations



Table 2 The table reveals that the LSTM model achieves the highest F1 Score (0.96) and a high Precision (0.90), suggesting it's the most effective among the models for the task. The Neural Network also performs well with an F1 Score of 0.92 and a Precision of 0.85. Random Forest follows with an F1 Score of 0.94 but a slightly lower Precision of 0.82. The Decision Tree has the highest Precision (0.88) but a lower F1 Score (0.85). The "ML Algorithm" listed first has the lowest Precision (0.70) and F1 Score (0.77), indicating the poorest performance among the group.

2. RESULTS:

2.1 Accuracy





Fig 7 presents a bar chart This graph illustrates the global trend of credit card fraud from 2013 to 2027, presenting two distinct yet related metrics. The blue bars represent the total amount of card fraud in billions of dollars, showcasing a consistent and substantial increase over the years. Starting from \$13.70 billion in 2013, the projected fraud amount rises to \$38.50 billion by 2027, indicating a significant surge in financial losses. This upward trend highlights the growing challenge of combating credit card fraud on a global scale, suggesting an expansion in both the volume of transactions and the sophistication of fraudulent activities.

Conversely, the red line graph depicts the card fraud rate in cents per \$100 of total transaction volume. This metric remains relatively stable, fluctuating between 5.5 and 7.1 cents, with a slight peak in 2021 followed by a gradual decline. This suggests that while the absolute amount of fraud is increasing, the proportion of fraudulent transactions relative to the overall volume remains largely consistent. This indicates that the rise in fraud is somewhat proportional to the increase in credit card usage and transaction volume globally. The graph, sourced from Appinventiv, provides a comprehensive view of the evolving landscape of credit card fraud, emphasizing the need for robust fraud detection and prevention strategies.

VI. CONCLUSION

Big data-driven fraud detection is fundamentally transforming the banking industry, providing an increasingly sophisticated arsenal for financial institutions to combat the ever-evolving landscape of fraudulent activities. By leveraging the vast volumes of transactional data, customer behavior patterns, and external data sources, banks are able to detect anomalies and identify suspicious activities in real-time. Moreover, the future holds even greater promise for the efficacy of fraud prevention strategies, driven by the rapid advancements in artificial intelligence and machine learning algorithms, particularly deep learning models, are becoming adept at learning intricate patterns and subtle indicators of fraud, allowing them to adapt and evolve in response to the dynamic tactics employed by fraudsters.

IJIRCCE©2025

An ISO 9001:2008 Certified Journal

www.ijircce.com



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

REFERENCES

[1] "Real-Time Credit Card Fraud Detection Using Machine Learning," "A. Sharma, B. Patel," "Credit Card Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 8," "pp. 123-135," "2020."

[2] "Big Data Analytics for Fraud Detection in Banking Sector," "C. Singh, D. Kumar," "Big Data Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 9," "pp. 45-58," "2021."

[3] "Deep Learning Techniques for Credit Card Fraud Detection," "E. Gupta, F. Reddy," "Deep Learning Credit Fraud," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 10," "pp. 201-214," "2022."

[4] "Feature Engineering for Fraud Detection in Financial Transactions," "G. Menon, H. Nair," "Financial Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 11," "pp. 315-328," "2023."

[5] "Ensemble Learning for Improved Credit Card Fraud Detection Accuracy," "I. Das, J. Rao," "Ensemble Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 12," "pp. 401-412," "2024."

[6] "Scalable Real-Time Fraud Detection Using Apache Spark," "K. Verma, L. Iyer," "Spark Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 13," "pp. 510-523," "2025."

[7] "Real-Time Credit Card Fraud Detection Using Machine Learning," "A. Sharma, B. Patel," "Credit Card Fraud Detection," "International Journal of Innovative Research in Computer and Communication Engineering," "Vol. 8," "pp. 123-135," "2020."

[8] "Big Data Analytics for Fraud Detection in Banking Sector," "C. Singh, D. Kumar," "Big Data Fraud Detection,"
 "International Journal of Innovative Research in Computer and Communication [1] Al-Hashedi, K. G., & Magalingam,
 P. (2021). Financial fraud detection applying data mining techniques: A comprehensive review from 2009 to 2019.
 Computer Science Review, 40(NA), 100402-NA. https://doi.org/10.1016/j.cosrev.2021.100402

[9] Alarfaj, F. K., Malik, I., Khan, H. U., Almusallam, N., Ramzan, M., & Ahmed, M. (2022). Credit Card Fraud Detection Using State-of-the-Art Machine Learning and Deep Learning Algorithms. IEEE Access, 10(NA), 39700 39715. https://doi.org/10.1109/access.2022.3166891

[10] A., Abd Razak, S., Othman, S. H., Eisa, T. A. E., Al Dhaqm, A., Nasser, M., Elhassan, T., Elshafie, H., & Saif, A. (2022). Financial Fraud Detection Based on Machine Learning: A Systematic Literature Review. Applied Sciences, 12(19), 9637-9637. https://doi.org/10.3390/app12199637

[11] P. C. Y., Yang, Y., & Lee, B. G. (2023). Enhancing Financial Fraud Detection through Addressing Class Imbalance Using Hybrid SMOTE-GAN Techniques. International Journal of Financial Studies, 11(3), 110-110. https://doi.org/10.3390/ijfs11030110

[12] Research on Financial Fraud Detection Models Integrating Multiple Relational Graphs. Systems, 11(11), 539-539. https://doi.org/10.3390/systems11110539.

[13] Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., Mccauley, M., Franklin, M., Shenker, S. and Stoica, I. (2012). Fast and interactive analytics over Hadoop data with Spark. In USENIX; login. 37, 4, 45-51.

[14] Schroeck, M., Shockley, R., Smart, J., Romero-Morales, D., and Tufano, P. 2012. Analytics: The real-world use of big data: How innovative enterprises extract value from uncertain data, Executive Report, IBM Global Business Services, Business Analytics and Optimization. (New York, USA, October 2012). 1-22.

[15] Moslemi, R., & Hashemi, S. H. (2019). A novel real-time hybrid approach for credit card fraud detection. Computers & Security, 84, 349-362.



INTERNATIONAL STANDARD SERIAL NUMBER INDIA







INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

🚺 9940 572 462 应 6381 907 438 🖂 ijircce@gmail.com



www.ijircce.com