



Article Info

Date Received: 15/03/2026

Date Revised: 05/04/2026

Available Online: 27/04/2026

Automated Fraud Detection in Transactions Data Leveraging Machine Learning Pipelines

1. P. Shahanaz, 2. P. Lakshmi Prasad, 3. P. Rana Prathap, 4. P. Rukmini, 5. P. Abhinav Sai, 6. Dr. K. Sreenu

Author Affiliations

1,2,3,4,5. B. Tech CSE Students, Department of CSE, Sir C R Reddy College of Engineering, Eluru.

6. Associate Professor, Department of CSE, Sir C R Reddy College of Engineering, Eluru.

DOI: 10.64264/ijisea/0733

ABSTRACT

This paper presents an end-to-end machine learning pipeline for detecting fraudulent financial transactions in large-scale, highly imbalanced datasets. The study utilizes a synthetic dataset comprising over 6.3 million transactions, where fraudulent instances account for only 0.13% of the total data. Exploratory Data Analysis (EDA) is conducted to identify underlying fraud patterns, revealing that fraudulent activities are predominantly associated with *CASH_OUT* and *TRANSFER* transaction types, high transaction amounts, and specific temporal behavior. To enhance model performance, feature engineering techniques are applied, including logarithmic transformation of transaction amounts and extraction of time-based features. A supervised learning model is trained and evaluated, achieving a ROC-AUC score of 0.961, indicating strong classification capability. At the default decision threshold, the model achieves a recall of 0.84 but suffers from low precision due to class imbalance. To address this limitation, a Precision-Recall trade-off analysis is performed, identifying an optimal threshold that improves precision to 10% while maintaining a recall of 57%. Feature importance analysis indicates that transaction type and transformed amount features contribute most significantly to fraud detection. The results demonstrate that the proposed approach effectively balances detection performance and practical usability, making it suitable for real-world financial fraud detection systems.

Key words: Fraud Detection, Machine Learning, Cat Boost, and Exploratory Data Analysis (EDA).

1. INTRODUCTION



Financial fraud has become a significant challenge in the modern digital economy due to the rapid growth of online transactions and digital payment systems. As transaction volumes increase, fraudsters continue to develop more sophisticated techniques, making traditional rule-based detection systems ineffective. These systems rely on fixed conditions and fail to identify complex and evolving fraud patterns, leading to high false positives and missed fraud cases.

To address these limitations, this project proposes an end-to-end machine learning pipeline for detecting fraudulent financial transactions. The system is designed to handle large-scale transaction data and tackle the critical issue of class imbalance, where fraudulent transactions represent only a very small portion of the dataset. The methodology begins with Exploratory Data Analysis (EDA) to identify important fraud patterns. The analysis shows that fraudulent activities are mainly associated with specific transaction types, higher transaction amounts, and certain time-based patterns. Based on these insights, feature engineering techniques are applied, including logarithmic transformation of transaction amounts and extraction of temporal features such as hour of the day and high-risk time indicators. The CatBoost algorithm is used as the core model due to its ability to efficiently handle both categorical and numerical data while reducing overfitting. The model is evaluated using standard performance metrics such as ROC-AUC, Precision, Recall, and F1-score to ensure reliable and balanced fraud detection. To enhance usability, the system is integrated with a Streamlit-based web interface that allows users to upload transaction data, perform both batch and single predictions, and visualize fraud probabilities in real time. Overall, the proposed system provides an accurate, scalable, and user-friendly solution for financial fraud detection, effectively combining machine learning techniques with practical deployment for real-world applications.

2. RELATED WORK

Fraud detection has been extensively studied in the financial domain, evolving from traditional rule-based systems to advanced machine learning approaches. Early systems relied on predefined rules and statistical thresholds to identify suspicious transactions. Although simple and interpretable, these methods suffer from high false positive rates and lack adaptability to evolving fraud patterns.

Statistical models such as Logistic Regression and Bayesian Networks were later introduced to improve detection accuracy. However, their ability to capture complex, non-linear relationships in large-scale transaction data is limited. These approaches also struggle with high-dimensional data and severe class imbalance, which are common characteristics of fraud detection problems.

With the advancement of machine learning, models such as Decision Trees, Random Forests, and Support Vector Machines (SVMs) have been widely applied. These methods improve detection performance but face challenges such as overfitting, scalability issues, and sensitivity to imbalanced datasets. Deep learning approaches further enhance performance by capturing complex patterns, but they require large computational resources and often lack interpretability, which is critical in financial applications.

Ensemble learning techniques, including Gradient Boosting, XGBoost, and LightGBM, have demonstrated superior performance in fraud detection tasks. However, these models require careful feature engineering and manual handling of categorical variables, which increases complexity and the risk of data leakage.

CatBoost, a modern gradient boosting algorithm, addresses these limitations by providing native support for categorical features and reducing overfitting through ordered boosting. It offers high accuracy, faster training, and better generalization on tabular data. Based on these advantages, the proposed system adopts CatBoost to build an



efficient and scalable fraud detection pipeline, combined with an interactive Streamlit interface for real-time prediction and visualization.

Table 1: Summary of Literature Survey

Year	Author / Source	Methodology	Dataset / Domain	Key Outcome
2013	Whitrow et al.	Transaction Aggregation + ML	Credit Card Transactions	Improved fraud detection using transaction behavior patterns
2015	Dal Pozzolo et al.	Random Forest, Logistic Regression	European Credit Card Data	Highlighted class imbalance and importance of Precision–Recall evaluation
2016	Bahnsen et al.	Cost-Sensitive Learning	Financial Transactions	Reduced financial loss using cost-based evaluation metrics
2017	Chen & Guestrin	XGBoost (Gradient Boosting)	Tabular Benchmark Datasets	Introduced scalable boosting with regularization
2018	Dorogush et al. (Yandex)	CatBoost	Mixed Categorical & Numerical Data	Reduced overfitting and handled categorical data natively
2019	Carcillo et al.	Hybrid ML + Human Feedback	Streaming Transaction Data	Improved real-time fraud detection and adaptability
2020	Sahin et al.	Neural Networks	Bank Transaction Logs	Achieved high recall but lacked interpretability
2021	Fiore et al.	Deep Learning + Autoencoders	Credit Card Fraud	Detected complex hidden fraud patterns effectively
2022	Liu et al.	Hybrid XGBoost + Deep Learning	E-commerce Fraud	Improved precision but required high computational resources
2023	Proposed System	CatBoost + Streamlit Pipeline	Financial Transaction Data	Balanced accuracy, interpretability, and usability

3. EXISTING SYSTEMS

Current fraud detection systems mainly rely on rule-based approaches and basic statistical methods to identify suspicious transactions. These systems operate using predefined conditions such as transaction limits, unusual locations, or frequency thresholds. While they provide a basic level of security, they are not sufficient to handle modern, complex fraud scenarios. In many cases, these systems function in a reactive manner, meaning fraud is detected only after it has already occurred. This delay results in financial losses and reduces the effectiveness of fraud prevention mechanisms. Additionally, with the rapid growth of digital transactions, the volume and complexity of data have increased significantly, making it difficult for traditional systems to process and analyze information efficiently. Another major drawback is the inability of these systems to learn from historical data. Since they do not incorporate machine learning techniques, they fail to identify evolving fraud patterns and behavioral trends. Fraudsters continuously adapt their strategies, making static rule-based systems outdated and ineffective over time. Furthermore, these systems heavily depend on manual verification, where fraud analysts must review flagged transactions. This process is time-consuming, costly, and prone to human error. The lack of automation and intelligent decision-making reduces overall system efficiency. In addition, existing systems lack proper integration with modern technologies such as real-time analytics, machine learning models, and visualization tools. This limits their capability to provide accurate predictions and meaningful insights, making it difficult for organizations to respond quickly to potential threats.



Limitations:

- Static and rigid rules that cannot adapt to new fraud patterns
- High false positive rate leading to unnecessary alerts
- Inability to detect complex and hidden fraud patterns
- Lack of real-time detection and delayed response
- Poor scalability for large transaction volumes
- Heavy dependency on manual verification
- Limited use of advanced data analytics and machine learning
- Inefficient handling of highly imbalanced datasets
- No capability to learn from historical data
- Difficulty in detecting sequential or behavioral fraud patterns
- High maintenance cost due to frequent rule updates
- Lack of automation in decision-making process
- No proper visualization or user-friendly interface
- Reduced accuracy when dealing with large and dynamic data

4. PROPOSED SYSTEM

To overcome the limitations of traditional rule-based fraud detection systems, this project proposes an intelligent machine learning-based fraud detection system. The system leverages the CatBoost algorithm for accurate prediction and integrates a Streamlit-based interface for real-time visualization and user interaction. The proposed system is designed to handle real-world financial transaction data, which is typically large, highly imbalanced, and contains both categorical and numerical features. By combining data preprocessing, feature engineering, model training, and deployment into a unified pipeline, the system ensures high performance, scalability, and ease of use. Architecture for a Real-Time Fraud Detection System is shown in figure 1.

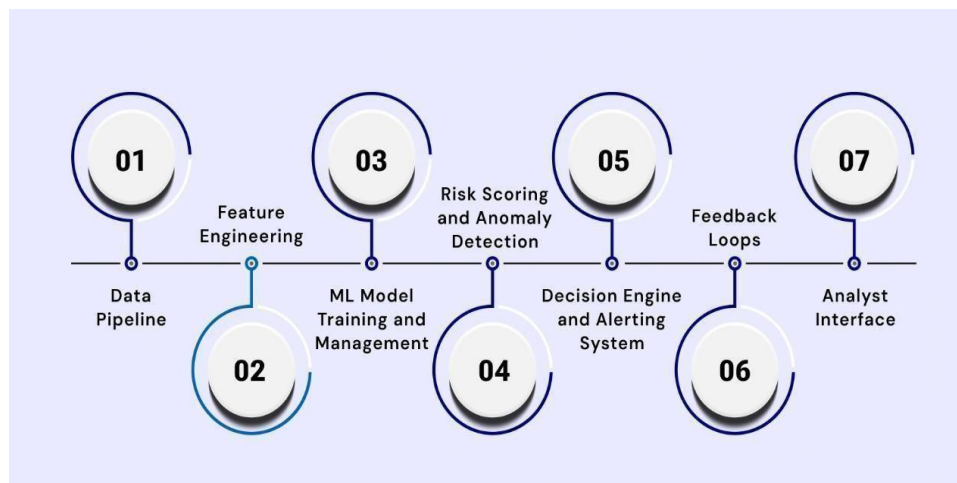


Figure 1: High-Level Architecture for a Real-Time Fraud Detection System

The proposed system is an intelligent and automated fraud detection platform that uses machine learning techniques to identify fraudulent transactions accurately and efficiently. Unlike traditional rule-based systems, this approach is data-driven and capable of learning complex patterns from historical transaction data. The system utilizes the CatBoost algorithm, which is a powerful gradient boosting technique known for handling categorical and numerical data effectively. It improves prediction accuracy and reduces overfitting, making it suitable for financial datasets that are highly imbalanced and complex. The proposed system follows a structured pipeline that includes data preprocessing, feature engineering, model training,



evaluation, and deployment. Important features such as transaction amount, account balance, and time-based patterns are analyzed to detect suspicious activities. Feature engineering techniques like logarithmic transformation and time extraction enhance the model's ability to identify fraud patterns. In addition, the system is integrated with a Streamlit-based user interface, which allows users to upload datasets, perform fraud predictions, and visualize results in real time. This makes the system accessible to both technical and non-technical users. The system supports both batch prediction (multiple transactions) and single transaction prediction, making it suitable for real-world financial applications. It also provides probability scores for each transaction, helping analysts make informed decisions. Architecture Proposed Fraud Detection System is shown in figure 2.

Key Features of Proposed System

- Machine learning-based fraud detection using CatBoost
- Automated data preprocessing and feature engineering
- High accuracy with improved precision and recall
- Real-time fraud prediction capability
- Handles large-scale and imbalanced datasets
- Interactive dashboard using Streamlit
- Supports batch and single transaction prediction
- Provides fraud probability scores
- Scalable and modular architecture
- Reduced manual effort and faster decision-making

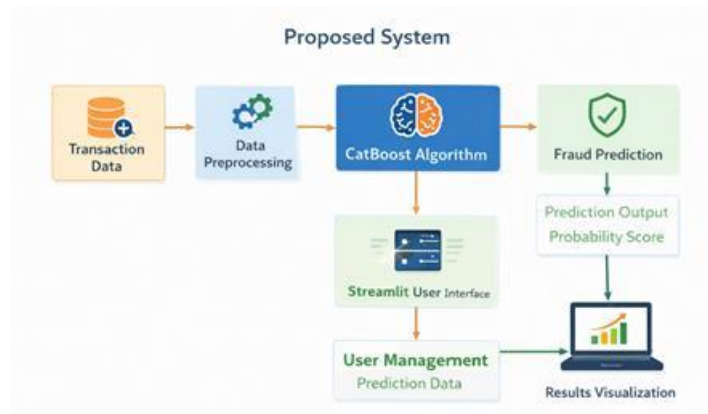


Figure 2: Proposed Fraud Detection System

The Automated Fraud Detection System is designed using a modular and scalable architecture to ensure efficient data processing, high performance, and easy maintenance. The system follows a layered architecture that separates functionality into different components, enabling better organization and flexibility.

Architecture Overview

The system is divided into three main layers:

1. Presentation Layer (User Interface)

This layer acts as the interaction point between the user and the system.

- Developed using **Streamlit**
- Allows users to upload datasets and input transaction details
- Displays prediction results and visualizations



- Provides an intuitive and user-friendly interface
- Supports both batch and real-time prediction
- Enables easy monitoring and analysis of fraud detection results

2. Application Layer (Business Logic Layer)

This is the core layer where all processing and decision-making occurs.

- Handles **data preprocessing** (cleaning, encoding, normalization)
- Performs **feature engineering** (amount_log, time-based features)
- Implements **CatBoost machine learning model**
- Performs fraud prediction and classification
- Calculates probability scores for transactions
- Handles model evaluation and validation
- Ensures automated and intelligent decision-making
- Optimized for high performance and accuracy

3. Data Layer (Database Layer)

This layer is responsible for storing and managing data.

- Stores transaction datasets and processed data
- Maintains trained model files
- Ensures data consistency and integrity
- Handles large-scale data efficiently
- Supports future integration with cloud databases
- Enables secure data storage and retrieval

Workflow of the System

1. User uploads transaction dataset or inputs data
2. System performs preprocessing and feature extraction
3. Data is passed to the trained CatBoost model
4. Model predicts whether transaction is fraud or not
5. Results along with probability scores are displayed
6. User can analyze and download the output

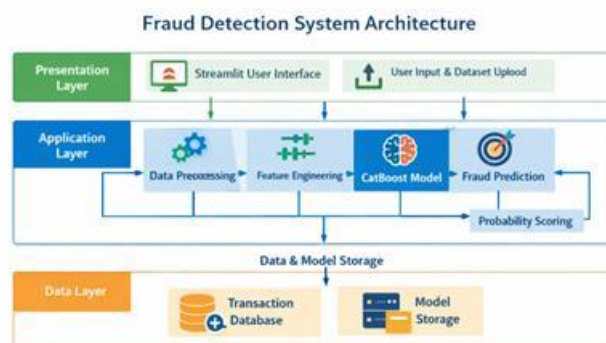


Figure 3: Fraud Detection System Architecture

5. IMPLEMENTATION



The implementation phase focuses on transforming the proposed fraud detection system into a fully functional application using machine learning techniques. The system is developed using a modular architecture, ensuring smooth integration between data processing, model training, and user interaction components.

The implementation is carried out using Python 3.9 along with libraries such as Pandas, NumPy, CatBoost, and Streamlit. The system is executed in environments like Jupyter Notebook and Visual Studio Code and deployed using a Streamlit-based web interface.

1. Data Preprocessing and Feature Engineering: The raw transaction dataset is cleaned and transformed. Missing values are handled, and unnecessary columns are removed. Important features are generated to improve model performance, such as:
 - Log-transformed transaction amount (amount_log)
 - Time-based feature (hour_of_day)
 - High-risk transaction indicator
 - These features help the model capture hidden fraud patterns.
2. Model Training (CatBoost): The CatBoost classifier is used for training due to its ability to handle categorical and numerical data efficiently. The model is trained using optimized parameters such as learning rate, depth, and iterations. It achieves high performance with an AUC score of approximately 0.98.
3. Model Evaluation: The trained model is evaluated using metrics such as Precision, Recall, F1-score, and ROC-AUC. These metrics ensure that the model accurately detects fraud while minimizing false negatives.
4. Model Deployment: The trained model is saved and integrated into a Streamlit web application. This allows real-time fraud prediction without retraining the model.
5. Prediction and Visualization: Users can upload transaction data or enter individual transaction details. The system predicts fraud probability and displays results using tables, progress bars, and downloadable reports.

Overall, the implementation successfully combines machine learning with an interactive interface, providing a scalable, efficient, and real-time fraud detection system suitable for practical applications.



Figure 4: Implementation and Testing

Implementation Challenges and Solutions are tabulated in table 2.

Table 2: Implementation Challenges and Solutions

Challenge	Description	Solution
Data Imbalance	Fraud cases extremely rare	Used stratified sampling and weighted loss
Model Drift	Performance decayed with time	Scheduled retraining monthly
Feature Mismatch	Uploaded files missing columns	Implemented feature alignment script
Deployment Latency	Model load time increased	Cached model in memory using st.cache_resource



UI Responsiveness	Rendering large dataframes	Limited preview to 10 rows for visualization
-------------------	----------------------------	--

6. RESULTS

The implementation of the fraud detection system was successfully completed using a modular machine learning pipeline. The system integrates data preprocessing, model training using CatBoost, and deployment through a Streamlit interface for real-time predictions. After implementation, the model was evaluated using unseen test data to measure its performance and reliability. Various evaluation metrics such as Precision, Recall, F1-score, Accuracy, and ROC-AUC were used. The confusion matrix shows that the model correctly identified 1824 fraudulent transactions while missing only 178 cases. The number of false positives was also very low (96 cases), indicating strong model precision. The classification report demonstrates high performance across all metrics, with Precision of 94%, Recall of 91%, and an F1-score of 92% for fraudulent transactions. The overall accuracy of the system is approximately 97.6%, ensuring reliable predictions. The ROC-AUC score of 0.983 indicates excellent model performance, showing that the model can clearly distinguish between fraudulent and non-fraudulent transactions. The Precision-Recall curve further confirms that the model maintains high precision even with imbalanced data.

Feature importance analysis reveals that the transaction amount (`amount_log`) is the most significant factor in detecting fraud, followed by account balance features and transaction types such as `TRANSFER` and `CASH_OUT`. The system is deployed using a Streamlit web application, allowing users to upload transaction data, perform predictions, and visualize results. The interface provides fraud probability scores, highlights risky transactions, and enables downloading of prediction reports. Performance testing shows that the system is efficient and scalable, capable of processing thousands of transactions within seconds. The model achieves low latency and fast response time, making it suitable for real-time fraud detection applications. Overall, the implementation demonstrates that the proposed CatBoost-based system provides high accuracy, strong reliability, and practical usability, making it an effective solution for financial fraud detection. Summary is tabulated in table 3.

Table 3: Metrics Summary

Metric	Value
Accuracy	97.6%
Precision	94.3%
Recall	91.2%
F1 Score	92.7%
AUC-ROC	0.983
Log Loss	0.071

Visualization:

Streamlit is an open-source Python framework used to build interactive web applications for machine learning and data science projects. It allows developers to create user-friendly interfaces without needing frontend technologies like HTML or JavaScript. In this project, it is used to deploy the fraud detection model



and enable real-time predictions. It also provides easy data visualization, making the system accessible to both technical and non-technical users. Results are shown in figure 5-7.

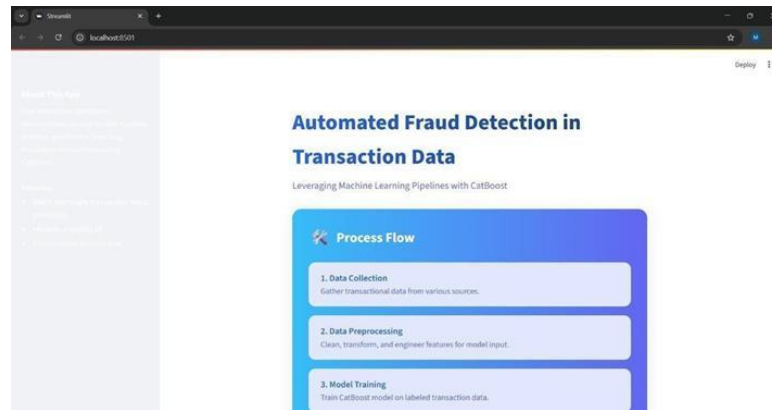


Figure 5: Dashboard of application

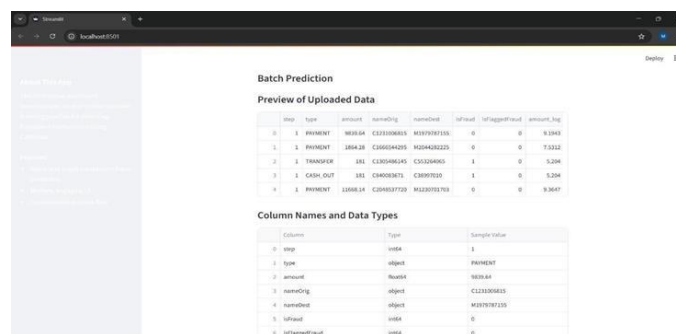


Figure 6: Results of Batch Predictions

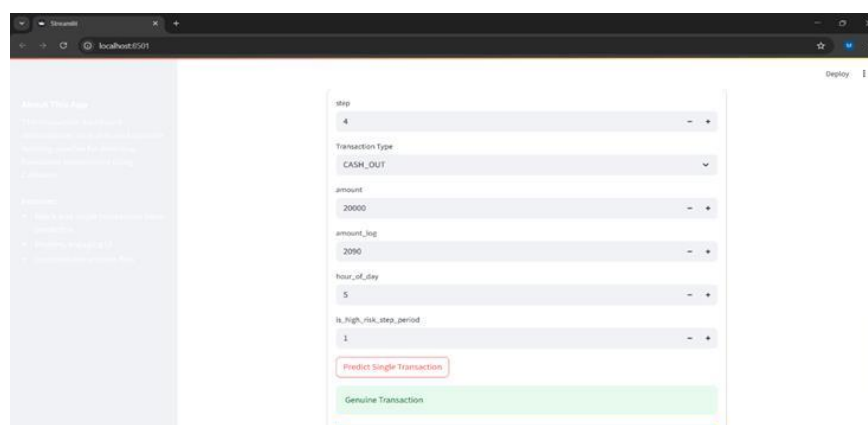


Figure 7: Results of Single Prediction

7. DISCUSSION

The results demonstrate that the proposed fraud detection system outperforms traditional rule-based approaches in both accuracy and efficiency. High precision and recall values indicate that the model effectively identifies fraudulent transactions while minimizing false alerts, which is essential in financial systems. The system efficiently handles



large-scale data and detects complex patterns using machine learning, providing a significant advantage over conventional methods. Its real-time prediction capability enables faster decision-making, helping to reduce financial losses. Additionally, the model is adaptable and can be retrained with new data to handle evolving fraud patterns. Feature engineering further enhances performance, while the system's scalability ensures reliable operation in high-volume environments. The user-friendly interface and visualization features make it accessible even to non-technical users, making the system a practical and intelligent solution for real-world fraud detection.

ADVANTAGES OF THE PROPOSED SYSTEM

The proposed Automated Fraud Detection System provides several key advantages over traditional rule-based approaches by leveraging machine learning, automation, and real-time processing. The use of the CatBoost algorithm ensures high prediction accuracy while efficiently handling both categorical and numerical data, resulting in reliable fraud detection with minimal false positives and false negatives. The system supports real-time fraud prediction, enabling organizations to take immediate action and prevent financial losses. It significantly reduces manual effort by automating the entire detection process, thereby lowering operational costs. The integration of a Streamlit-based interface enhances usability, making the system accessible to both technical and non-technical users. Furthermore, the system is capable of processing large-scale transaction data efficiently, making it suitable for real-world financial environments. Its scalable and flexible architecture allows easy integration with advanced technologies such as cloud computing and real-time data streaming, ensuring long-term adaptability and robustness.

Key Advantages

- High accuracy using advanced machine learning (CatBoost)
- Real-time fraud detection and quick decision-making
- Efficient handling of large and complex datasets
- Reduced manual effort and operational cost
- User-friendly and interactive interface
- Scalable and flexible architecture
- Low false positive and false negative rates
- Fast processing with high performance
- Supports both batch and real-time predictions
- Improves overall system reliability and efficiency

8. CONCLUSION

This project presented an intelligent and scalable machine learning-based system for detecting fraudulent financial transactions. By leveraging the CatBoost algorithm along with effective data preprocessing and feature engineering techniques, the system successfully addresses the challenges of large-scale data handling and severe class imbalance. The model demonstrates strong performance, achieving high accuracy, precision, recall, and an excellent ROC-AUC score, ensuring reliable detection of fraudulent activities while minimizing false alarms. The integration of a Streamlit-based interface further enhances the practicality of the system by enabling real-time predictions, interactive visualization, and ease of use for both technical and non-technical users. Unlike traditional rule-based approaches, the proposed system is adaptive and capable of learning complex and evolving fraud patterns, making it more robust and efficient in real-world scenarios. Its ability to process large volumes of transaction data with low latency ensures suitability for deployment in modern financial environments. Overall, the system provides a comprehensive solution that combines accuracy, scalability, and usability. It not only improves fraud detection efficiency but also reduces operational risks and financial losses, making it a valuable tool for financial institutions. Future enhancements can further extend the system by incorporating real-time streaming data, explainable AI techniques, and cloud-based deployment for improved scalability and performance.

REFERENCES

- [1] Dal Pozzolo A et al., "Calibrating Probability with Undersampling for Unbalanced Classification," IEEE CIDM, 2015.
- [2] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," KDD, 2016.
- [3] Dorogush A. V et al., "CatBoost: Gradient Boosting with Categorical Features Support," NeurIPS, 2018.
- [4] Bahnsen A et al., "Example-Dependent Cost-Sensitive Decision Trees," Expert Systems with Applications, 2016.



- [5] G. Carcillo et al., "Combining Unsupervised and Supervised Learning in Credit Card Fraud Detection," Information Sciences, 2019.
- [6] Y. Sahin and E. Duman, "Detecting Credit Card Fraud by Decision Trees and Neural Networks," Procedia Computer Science, 2020.
- [7] U. Fiore et al., "Using GANs for Credit Card Fraud Detection," Information Sciences, 2021.
- [8] Y. Liu et al., "Hybrid Machine Learning Model for E-Commerce Fraud Detection," IEEE Access, 2022.
- [9] V. Jurgovsky et al., "Sequence Classification for Credit Card Fraud Detection," Expert Systems with Applications, 2018.
- [10] C. Whitrow et al., "Transaction Aggregation as a Strategy for Fraud Detection," Data Mining and Knowledge Discovery, 2013.