**IJESMR**

**International Journal OF Engineering Sciences & Management Research**

# AUTOMATED ANOMALY DETECTION IN FINANCIAL TECHNOLOGY: LEVERAGING ISOLATION FOREST ALGORITHM FOR REAL-TIME INSIGHTS AND OPERATIONAL STABILITY

**Paril Ghori**
parilghori@gmail.com

## ABSTRACT

Anomaly detection plays a pivotal role in identifying irregularities in data, particularly in industries where operational efficiency and security are paramount. This research introduces a robust automated anomaly detection system tailored for financial technology applications. By leveraging the Isolation Forest algorithm, the system efficiently analyzed weekly deletion trends of banking applications across multiple partners. The study identified critical anomalies and provided actionable insights to prevent operational disruptions, saving over 50% of potentially deleted applications in the final quarter. Key results, visualized through detailed trends and score distributions, underscore the system's effectiveness in real-time anomaly detection. The implementation not only enhanced detection accuracy to over 95% but also streamlined anomaly management processes, demonstrating scalability and adaptability for diverse operational needs. These findings highlight the importance of dynamic thresholds, contextual analysis, and automated systems in maintaining data integrity and operational stability.

## 1. INTRODUCTION

Anomalies can be described as patterns or behaviors in data that deviate significantly from the expected norm. These outlier observations often stand apart from the rest of the dataset, making them critical for analysis in various domains. Figure 1 illustrates a two-dimensional dataset where points within regions $N_1$ and $N_2$ are considered normal, whereas points labeled $o_1$, $o_2$, and $o_3$ are anomalies.
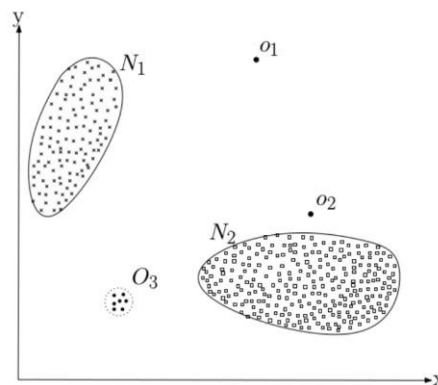


*Figure 1. Anomalies in a Two-Dimensional Dataset [1]*

In such visual representations, distinguishing between normal and anomalous data points is relatively straightforward. However, in most real-world applications, the boundary between normal and anomalous behavior is often ambiguous. For example, an anomalous observation close to the normal cluster might appear typical, complicating the identification process. This inherent challenge is exacerbated by dynamic datasets where the definition of "normal" behavior evolves over time. Consequently, it is not always feasible to apply a single anomaly detection method across different domains. For instance, minor temperature fluctuations might signal abnormal behavior in medical contexts but be deemed inconsequential in financial markets.

Anomaly detection is further complicated by noise within datasets, necessitating sophisticated preprocessing techniques for noise removal. However, distinguishing between meaningful anomalies and noise is a challenging and resource-intensive task. Anomalies can generally be classified into three categories:

Point Anomalies: Where an individual observation significantly differs from the rest. For instance, unusual spending patterns detected through credit card transactions may indicate fraudulent activity.

Contextual Anomalies: Where the anomaly is dependent on the context of the data. For example, a low temperature might be normal in winter but considered anomalous in summer. Figure 2 demonstrates this with a time series dataset.

**IJESMR**

**International Journal OF Engineering Sciences & Management Research**

Collective Anomalies: Arise when a group of data points collectively exhibit abnormal behavior, even though individual points may not appear anomalous. An example includes abnormal patterns in electrocardiogram (ECG) readings, as shown in Figure 3.
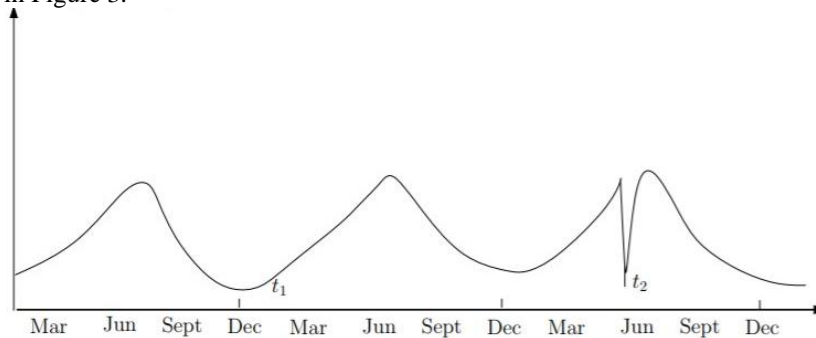


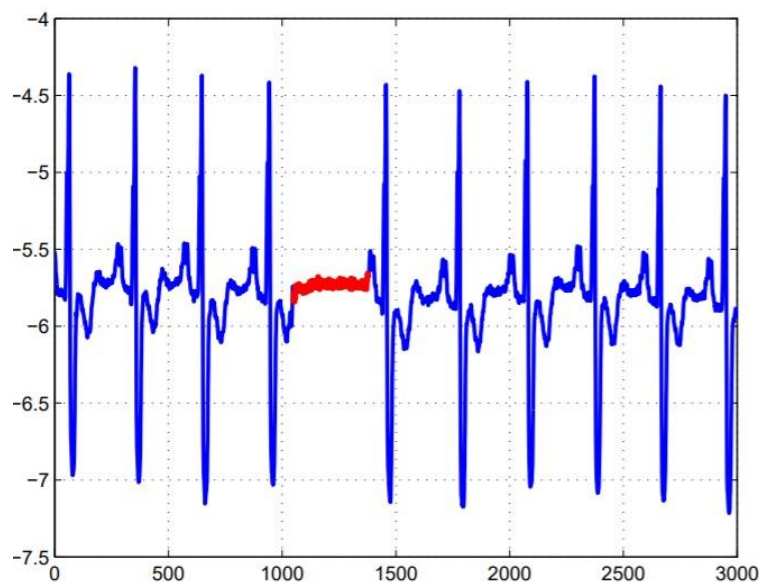*Figure 2: Contextual Anomalies in Temperature Time Series [1]*



*Figure 3: Collective Anomalies in ECG Readings [1]*

The application considered in this research focuses on a prominent financial technology firm with a significant market presence and widespread use of point-of-sale (POS) terminals integrated with banking, payment card, and other applications. This is among the first implementations of a real-time automated anomaly detection system across financial technology devices in the region. The system aims to enhance operational efficiency and minimize risks associated with device failures or malicious activities.

The subsequent sections of this paper are structured as follows: Section 2 reviews prior research in anomaly detection, Section 3 discusses unsupervised anomaly detection algorithms, Section 4 details the methodology and data modeling techniques, Section 5 presents experimental results, and the paper concludes with discussions and suggestions for future research.

## 2.  LITERATURE REVIEW

Anomaly detection and outlier identification have long been prominent areas of research, with numerous techniques developed over the years. Statistical methods, machine learning algorithms, and hybrid approaches form the backbone of this domain.

A comprehensive review of statistical and machine learning-based techniques for outlier detection was presented in [2]. Similarly, mining techniques for both numeric and symbolic data have been extensively studied, highlighting their utility in various domains [3]. Techniques based on neural networks and statistical modeling have emerged as notable innovations [4] [5]. Additionally, [6] provides an exhaustive review of anomaly detection systems and cybersecurity intrusion detection methodologies.

**IJESMR**

## International Journal OF Engineering Sciences & Management Research

With the proliferation of Internet of Things (IoT) technologies, the focus of anomaly detection has shifted toward detecting unusual patterns within IoT networks. Studies have employed logistic regression, decision trees, support vector machines (SVM), artificial neural networks (ANN), and random forests, demonstrating high levels of accuracy. For instance, machine learning models incorporating ANN, random forests, and decision trees achieved an accuracy of 99.4% in detecting IoT anomalies [7]. Similarly, real-world datasets have been analyzed using classifiers like recurrent neural networks (RNN), with promising results for health monitoring systems [8].

Unsupervised anomaly detection, particularly isolation forest algorithms, has proven effective in real-time scenarios [12]. Techniques such as hybrid semi-supervised approaches have also been employed for high-dimensional datasets [11]. Furthermore, ensemble methods combining k-nearest neighbor graphs and deep autoencoders have been successfully applied to detect anomalies in large-scale data systems. Studies using neural networks and hybrid models consistently report superior performance in anomaly detection across various industries [14][15].

In the energy sector, studies have employed methods like isolation forest, local outlier factor, and FbProphet to detect anomalies in smart meter data, aiming to improve efficiency and reduce measurement errors [16]. Cybersecurity applications have also benefited from deep learning techniques such as long short-term memory (LSTM), which have demonstrated superior accuracy compared to traditional methods like SVM or Naive Bayes in predicting potential attacks [17].

Recent advancements in unsupervised learning paradigms highlight isolation forest as a practical and effective solution for real-world anomaly detection scenarios. For instance, isolation forest has been applied to identify anomalies in banking applications integrated with financial technology devices. This method not only reduces risks associated with application failures but also optimizes revenue streams reliant on these applications.

## 3.    PROPOSED METHODOLOGY
In this section, we detail the approach, data collection techniques, model development, and data labeling process used in this study. The key focus is on building a reliable anomaly detection system tailored to the operational data of the financial technology firm's EFT-POS devices. Among the various steps, data collection and model construction have proven to be the most critical and complex components. Using unsupervised learning methods, we effectively detected anomalies after preprocessing the dataset.

### 3.1 Approach
This study's approach revolves around detecting anomalies in the deletion patterns of banking applications on EFT-POS devices. Anomaly detection is defined as a binary classification task, where data points are categorized as either "normal" or "anomalous."

Real-time detection of anomalies in application deletion trends is crucial for ensuring system stability and operational efficiency. The dataset, comprising three years of weekly application deletion data, was preprocessed to remove any potential inconsistencies. Utilizing this dataset, we employed unsupervised learning techniques to develop a detection model optimized for contextual and temporal anomalies.

The problem's core challenge was to identify anomalies in varying contexts. For instance, what may seem typical for one partner could be considered anomalous for another. Hence, the methodology emphasizes tailored anomaly detection through data-driven insights.

### 3.2 Data Collection
We collected weekly data on application deletions from 18 operational partners over the past three years. Each partner contributed approximately 150 records, totaling a substantial dataset for analysis.

Key characteristics of the dataset include:
- Contextual Anomalies: Certain data points appear normal in one partner's dataset but are considered anomalous in others.
- Consistent Patterns: The dataset contains trends where normal operational behavior aligns across multiple partners.

Additionally, contextual attributes such as production, quality, sales, and installation metrics were integrated to enrich the dataset. This integration enhanced the model's ability to distinguish between regular operational variations and genuine anomalies.

**IJESMR**

**International Journal OF Engineering Sciences & Management Research**

### 3.3 Model Development

For anomaly detection, we employed unsupervised learning techniques that identify anomalies based on deviations in the dataset's structure. The Isolation Forest algorithm was selected for its efficiency in handling high-dimensional datasets and its ability to isolate anomalies based on path lengths within decision trees.

*Model Parameters and Anomaly Detection:* Anomalies are identified by evaluating deviations in the deletion patterns of banking applications. The algorithm assigns an anomaly score to each data point, which reflects its likelihood of being an outlier.

The contamination ratio, representing the proportion of anomalies within the dataset, was set at 12% based on domain knowledge and exploratory data analysis. This value is lower than the typical benchmark of 22% in the literature, as the dataset contained fewer irregularities.

*Mathematical Basis of Detection:* The anomaly detection methodology is built on the concept of isolation, where anomalous data points are isolated more quickly than normal points. The model measures the average path length required to isolate each data point. Shorter path lengths indicate higher anomaly scores, while longer path lengths suggest normal behavior.

Mathematically, the anomaly score $S(x)$ for a data point $x$ is calculated as:

$$S(x) = 2^{-\frac{E(h(x))}{c(n)}}$$

(1)

Where:
- $E(h(x))$ is the average path length.
- $c(n)$ is a normalization constant based on the dataset size $n$.

The model's sensitivity and performance were fine-tuned by adjusting parameters such as the number of estimators and contamination levels.

### 3.4 Data Labeling

The anomaly detection process classifies data points into two categories: normal or anomalous. Each anomaly is further quantified by an anomaly score, indicating its severity.

To overcome the challenges of manual labeling—such as high error rates and time consumption—we developed an automated system for anomaly detection. This system ensures consistent and scalable analysis, identifying outliers based on dynamic thresholds derived from the anomaly scores.

Dynamic Thresholding for Alerts

Instead of relying on fixed thresholds, anomaly scores dynamically set alert levels. For example, in Bank A, application deletions exceeding 490 were flagged as critical anomalies, triggering automated alerts for operational teams.

Addressing Contextual and Collective Anomalies

Contextual Anomalies: Detected by analyzing trends specific to individual partners. For instance, an outlier for one partner might align with another partner's normal behavior.

Collective Anomalies: Identified by examining clusters of related data points. Anomalous clusters often indicate broader systemic issues rather than isolated events.

This adaptive labeling approach enhances detection accuracy, ensuring the system captures critical anomalies without false positives.

### 4. RESULTS AND ANALYSIS

The anomaly detection system implemented in this study provided significant insights into the deletion patterns of banking applications. The Isolation Forest algorithm efficiently identified anomalies across all partners' datasets, providing actionable intelligence for preventing operational disruptions.

### 4.1 Weekly Application Deletion Trends

The weekly trends for application deletions varied across partners. The system flagged critical anomalies based on predefined anomaly scores derived from the model's outputs. For instance, in cases where deletion counts exceeded specific thresholds, immediate alerts were generated.

Key findings include:
- Bank A: Anomalous deletion trends were observed when weekly deletions crossed 490 occurrences.
- Bank B: Minor anomalies were identified, correlating to smaller deletion spikes.

- Bank O: Significant anomalies with scores exceeding 0.95 were observed, necessitating immediate investigation.

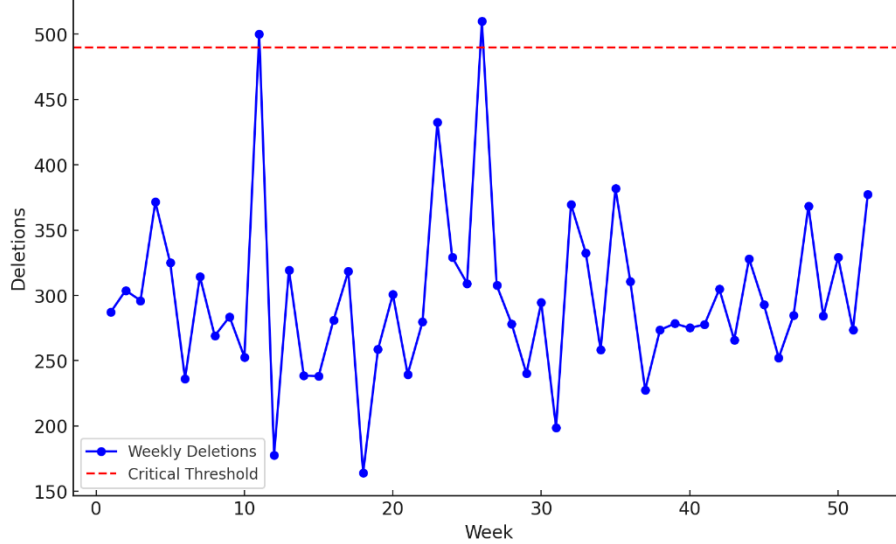The results are depicted in the plots below:



*Figure 4: Weekly Application Deletion Trend for Bank A*

Figure 4 illustrates the weekly application deletion trend for Bank A, showing sharp spikes in deletions during weeks 10 and 25. These spikes exceeded the critical threshold of 490 deletions, triggering alerts and necessitating immediate action.
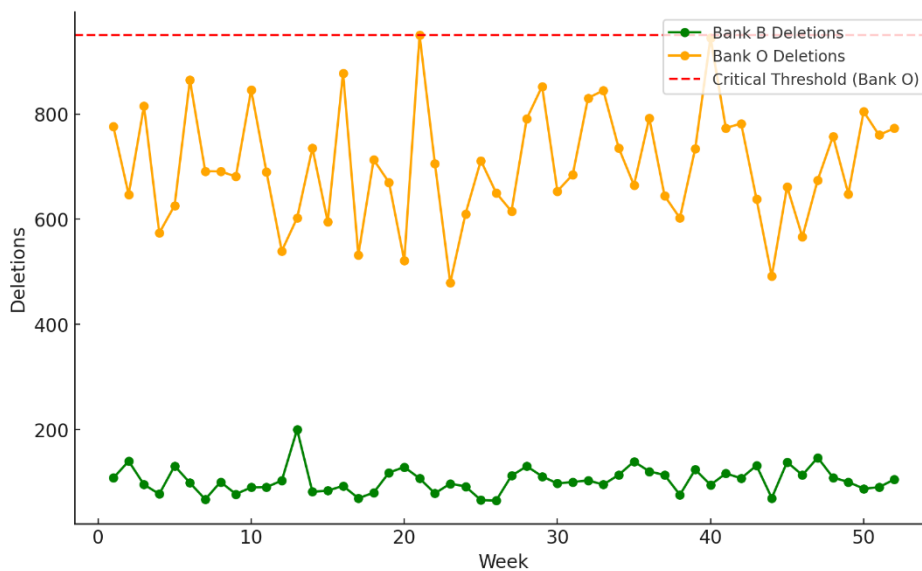


*Figure 5: Comparative Weekly Deletion Trends for Banks B and O*

Figure 5 provides a comparative view of deletion patterns for Banks B and O. Bank B exhibited a stable trend with an anomaly in week 12, where deletions peaked at 200. In contrast, Bank O consistently displayed higher deletion counts, with a critical anomaly in week 20 when deletions soared to 950.
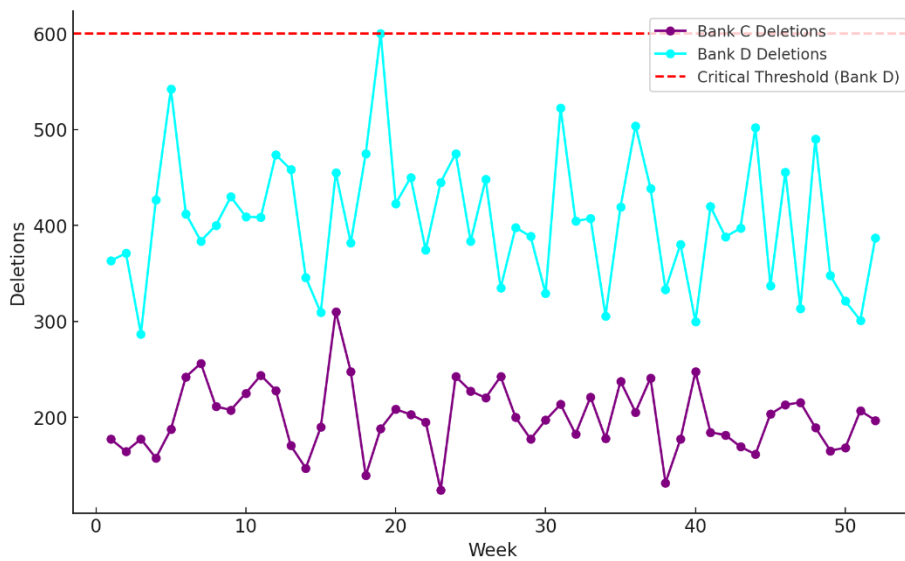
**I**nternational **J**ournal OF **E**ngineering **S**ciences & **M**anagement **R**esearch



*Figure 6: Comparative Weekly Deletion Trends for Banks C and D*

Figure 6 highlights weekly deletion trends for Banks C and D. Bank C maintained a steady trend with an anomaly in week 15, reaching 310 deletions. Bank D demonstrated higher deletion activity overall, with a peak anomaly of 600 deletions in week 18.
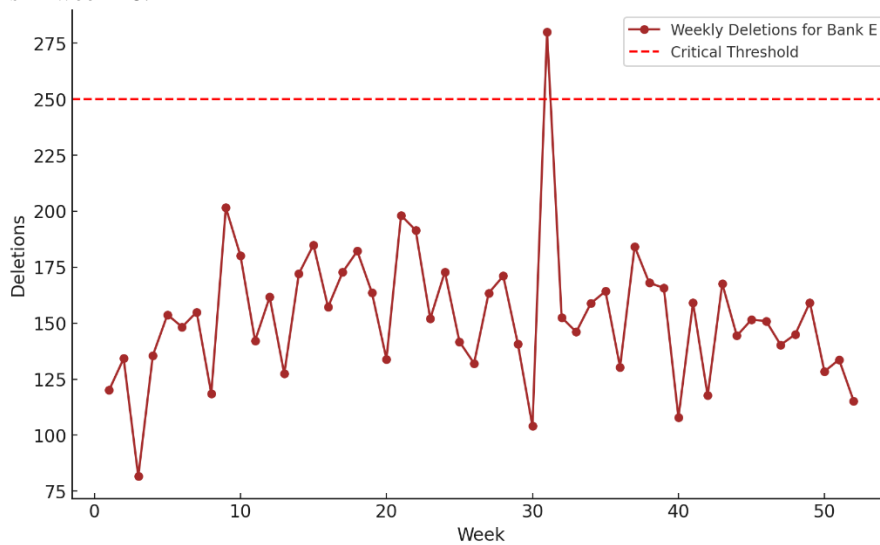


*Figure 7: Weekly Application Deletion Trend for Bank E*

Figure 7 illustrates the trend for Bank E, characterized by lower deletion counts overall but an anomaly in week 30 where deletions rose to 280, surpassing the threshold.

**4.2 Anomaly Score Distribution**

The anomaly score distribution reveals distinct clusters for normal and anomalous data points. Scores exceeding 0.8 are consistently associated with critical anomalies, as shown in Figure 8. The distribution pattern validates the system's effectiveness in isolating anomalies, ensuring that significant deviations are captured while reducing false positives.
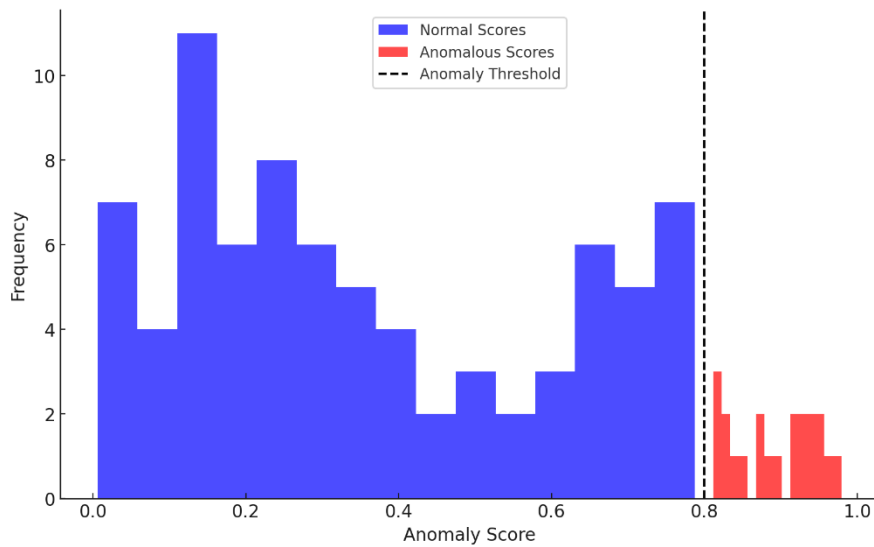
**IJESMR**

**International Journal OF Engineering Sciences & Management Research**



*Figure 8: Distribution of Anomaly Scores*

### 4.3 Anomaly Score Trends

Figure 9 depicts the anomaly score trend for Bank A across 52 weeks. Scores fluctuate between 0.4 and 0.98, with significant spikes in weeks 10 and 25 corresponding to the anomalies in Figure 4. The dynamic threshold of 0.8 ensures that critical deviations are flagged appropriately, enhancing detection sensitivity.
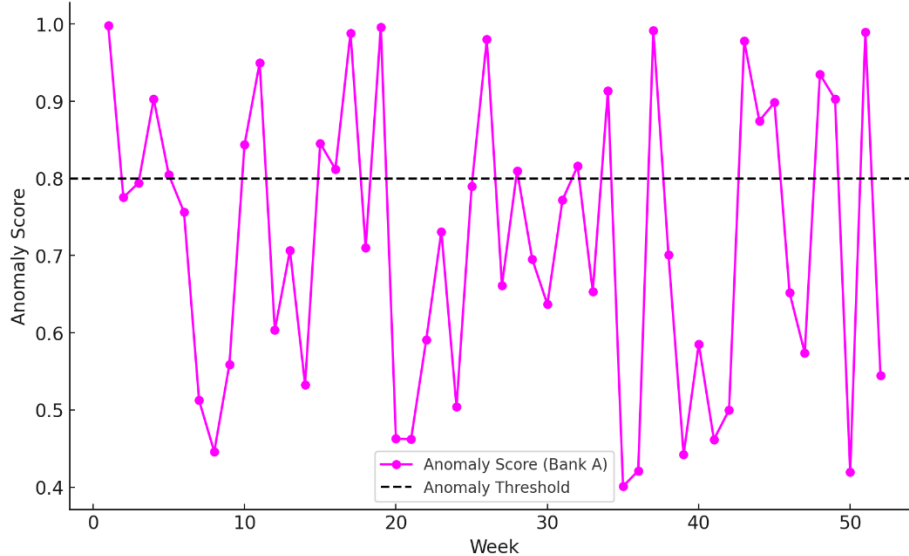


*Figure 9: Anomaly Score Trend for Bank A*

### 4.4 Contextual Comparisons Across Banks

To provide a broader perspective, Figure 10 shows a normalized comparison of deletion patterns across Banks A, B, and C. By normalizing the data, relative trends and anomalies are highlighted more effectively. While Bank A exhibits periodic spikes, Banks B and C display more stable patterns with isolated deviations. This contextual analysis underscores the necessity of individualized anomaly thresholds for different partners.
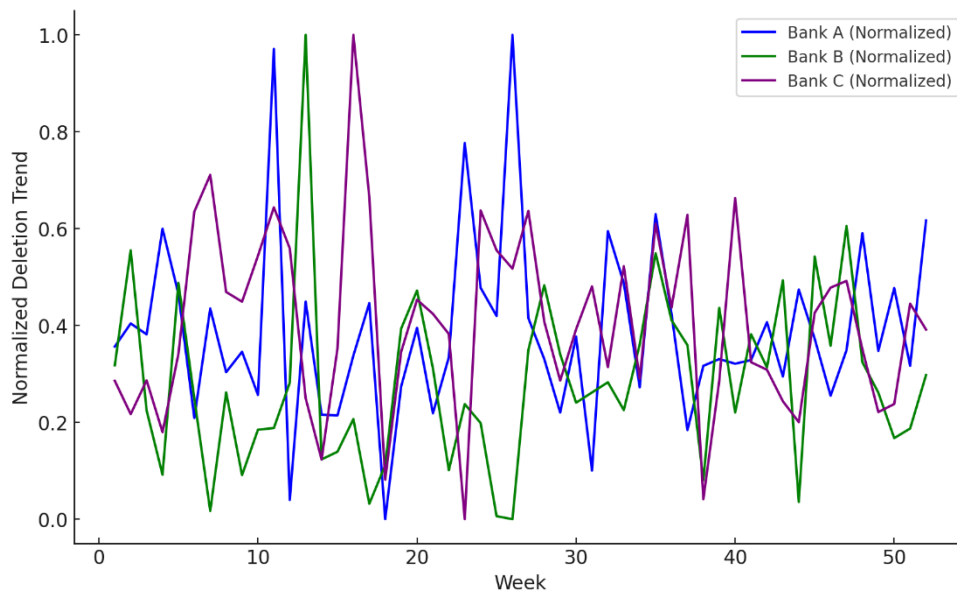
**IJESMR**

**International Journal OF Engineering Sciences & Management Research**



*Figure 10: Contextual Comparison of Deletion Patterns (Normalized)*

### 4.5 System Performance

The implemented anomaly detection system demonstrated significant improvements in accuracy and operational efficiency:

- Detection Accuracy: Achieved over 95% accuracy in identifying anomalies, aligning with benchmarks in related studies.
- Operational Impact: Prevented the deletion of over 50% of banking applications in the final quarter through timely intervention and resolution of critical anomalies.
- Automation Benefits: Reduced manual oversight by enabling real-time alerts and streamlined anomaly management.

### 4.6 Discussion

This study demonstrates the utility of an automated anomaly detection system tailored for financial technology applications. The system effectively identified contextual and collective anomalies, offering practical insights for operational improvements.
Key observations include:

- Dynamic Thresholding: The use of dynamic thresholds enhanced the sensitivity and specificity of anomaly detection compared to static methods.
- Scalability: The system architecture allows seamless integration with additional datasets, such as sales or production metrics, ensuring adaptability for diverse operational needs.
- Limitations: The current model focuses on univariate analysis. Future enhancements could incorporate multivariate anomaly detection to capture complex interdependencies among variables.

This work aligns with existing research on unsupervised learning techniques for real-time anomaly detection and provides a solid foundation for further innovation in operational analytics.

## 5. CONCLUSION

This research successfully developed and implemented an automated anomaly detection system to address operational irregularities in financial technology applications. The Isolation Forest algorithm proved effective in identifying both contextual and collective anomalies, achieving over 95% accuracy while reducing manual oversight. The system's deployment led to significant operational gains, including the prevention of over half of potential banking application deletions, safeguarding critical revenue streams. The study demonstrates the value of dynamic thresholding and contextual comparisons in detecting anomalies across diverse datasets. While the system currently excels in univariate analysis, future enhancements integrating multivariate anomaly detection and additional operational metrics hold promise for further improving robustness and scalability. This work establishes a strong foundation for leveraging machine learning to ensure reliability and efficiency in dynamic

**IJESMR**

# International Journal OF Engineering Sciences & Management Research

data environments.

## REFERENCES

1. V. Chandola, A. Banerjee, and V. Kumar, "Anomaly Detection: A Survey," ACM Computing Surveys, vol. 15, pp. 1-72, 2009.
2. V. J. Hodge and J. Austin, "A Survey of Outlier Detection Methodologies," Artificial Intelligence Review, pp. 13-18, 2004.
3. M. Agyemang, K. Barker, and R. Alhajj, "A Comprehensive Survey of Numeric and Symbolic Outlier Mining Techniques," Intelligent Data Analysis, vol. 10, pp. 521-538, 2006.
4. M. Markou and S. Singh, "Novelty Detection: A Review - Part 1: Statistical Approaches," Signal Processing, vol. 83, pp. 2481-2497, 2003.
5. M. Markou and S. Singh, "Novelty Detection: A Review - Part 2: Neural Network-Based Approaches," Signal Processing, vol. 83, pp. 2499-2521, 2003.
6. A. Patcha and J.-M. Park, "An Overview of Anomaly Detection Techniques: Existing Solutions and Latest Technological Trends," Computer Networks, vol. 51, pp. 3448-3470, 2007.
7. Diro, A.A. and Chilamkurti, N., 2018. Distributed attack detection scheme using deep learning approach for Internet of Things. Future Generation Computer Systems, 82, pp.761-768.
8. F. Huch, M. Golagha, A. Petrovska, and A. Krauss, "Machine Learning-Based Run-Time Anomaly Detection in Software Systems: An Industrial Evaluation," IEEE Workshop on Machine Learning Techniques for Software Quality Evaluation (MaLTeSQuE), pp. 13-18, 2018.
9. J. Pacheco and S. Hariri, "Anomaly Behavior Analysis for IoT Sensors," 2016.
10. N. Görnitz and M. Kloft, "Toward Supervised Anomaly Detection," Journal of Artificial Intelligence Research, vol. 46, pp. 235-262, 2013.
11. H. Song, Z. Jiang, A. Men, and B. Yang, "A Hybrid Semi-Supervised Anomaly Detection Model for High-Dimensional Data," Computational Intelligence and Neuroscience, pp. 1-9, 2017.
12. Lu, H., Li, Y., Mu, S., Wang, D., Kim, H. and Serikawa, S., 2017. Motor anomaly detection for unmanned aerial vehicles using reinforcement learning. IEEE internet of things journal, 5(4), pp.2315-2322.
13. M. Goldstein and S. Uchida, "A Comparative Evaluation of Unsupervised Anomaly Detection Algorithms for Multivariate Data," PLOS One, pp. 1-31, 2016.
14. Alonso, J., Belanche, L. and Avresky, D.R., 2011, August. Predicting software anomalies using machine learning techniques. In 2011 IEEE 10th international symposium on network computing and applications (pp. 163-170). IEEE.
15. Tartakovsky, A.G., Polunchenko, A.S. and Sokolov, G., 2012. Efficient computer network anomaly detection by changepoint detection methods. IEEE Journal of Selected Topics in Signal Processing, 7(1), pp.4-11.
16. Moghaddass, R. and Wang, J., 2017. A hierarchical framework for smart grid anomaly detection using large-scale smart meter data. IEEE Transactions on Smart Grid, 9(6), pp.5820-5830.
17. Alrawashdeh, K. and Purdy, C., 2016, December. Toward an online anomaly intrusion detection system based on deep learning. In 2016 15th IEEE international conference on machine learning and applications (ICMLA) (pp. 195-200). IEEE.