**ANOMALY DETECTION IN CLOUD SERVER LOGS USING BASIC TEXT MINING**

VARSHITHA YANAMALA

Data Science, University of Maryland, Baltimore County

AZ70378

Prof. Najam Hassan

DATA 603 Platforms for Big Data Processing

## Table of Contents

## Introduction

The security and dependability of cloud computing systems depend on the capacity to identify irregularities in cloud server logs. Finding anomalous activity, security risks, and system failures is made easier by the abundance of unstructured data in these logs, which provide important insights into the operational health of cloud systems. The main argument of this paper is that basic text mining techniques are an efficient and practical way to detect anomalies in unstructured cloud server log data, ultimately enhancing system monitoring and reducing risks associated with security breaches and operational failures. Traditional rule-based and statistical methods are often inadequate for handling the complexity and scale of modern cloud environments. In contrast, text mining combined with machine learning offers a more adaptable and scalable solution.

Le & Zhang (2022) highlight the role of deep learning in capturing intricate patterns within unstructured log data, something traditional techniques struggle to achieve. Vervaet (2023) presented MoniLog, an automated system capable of real-time anomaly detection, demonstrating the increasing need for adaptive, sophisticated solutions. Meng et al. (2022) emphasize the value of unsupervised methods such as LogAnomaly, which can identify both sequential and quantitative anomalies in large log datasets, illustrating the potential of text mining approaches. Goldstein & Uchida (2016) conducted comparative evaluations of various anomaly detection algorithms, providing insights into the strengths and limitations of different models for managing multivariate data.

Yasarathna & Munasinghe (2020) studied anomaly detection in cloud network data, highlighting the difficulties in detecting rare events in large datasets. Luo et al. (2021) proposed a robust unsupervised anomaly detection framework, showcasing how machine learning advances can improve detection accuracy in cloud environments. Alzoubi et al. (2024) discuss recent trends in machine learning for cloud security, underscoring the importance of scalability and explainability in building trust in automated detection systems.

An extensive examination of the development of anomaly detection techniques, from early rule-based systems to contemporary machine learning-driven approaches, will be presented in this presentation. The primary techniques for text mining anomaly detection, the difficulties in putting these technologies into practice, and the possible effects of sophisticated anomaly detection techniques on cloud infrastructure security will all be covered. Through

an assessment of the shortcomings and possibilities presented by existing methods, this research seeks to advance knowledge of how text mining techniques might improve cloud system security and dependability.

Le & Zhang (2022) demonstrate the effectiveness of deep learning in capturing complex patterns in unstructured log data, a task at which traditional approaches often fall short. Vervaet (2023) introduced MoniLog, an automated system capable of detecting anomalies in cloud environments in real time, highlighting the growing need for adaptive and sophisticated solutions. Meng et al. (2022) also emphasize the importance of unsupervised methods like LogAnomaly for identifying sequential and quantitative anomalies in large log datasets, further showcasing the potential of text mining-based anomaly detection approaches. Scholars like Goldstein & Uchida (2016) have provided comparative evaluations of various anomaly detection algorithms, shedding light on the strengths and limitations of different approaches in handling multivariate data.

Furthermore, Yasarathna & Munasinghe (2020) explored anomaly detection in cloud network data, underscoring the challenges associated with identifying rare events in vast datasets. A strong unsupervised anomaly detection framework was presented by Luo et al. (2021), showing how machine learning developments can be applied to improve anomaly detection accuracy in cloud environments. Recent developments in machine learning for cloud security are covered by Alzoubi et al. (2024), who stress the importance of scalable and explicable models in building confidence in automated anomaly detection systems.

This paper provides a detailed analysis of the development of anomaly detection methods, from the first rule-based systems to the most recent machine learning-driven methods. We go over the main approaches for detecting anomalies in text mining, the difficulties and problems associated with putting these technologies into practice, and the possible effects of sophisticated anomaly detection methods on the security of cloud infrastructure.

The arguments presented in this paper are built upon recent advancements in anomaly detection research. Studies such as those by Li et al. (2022) demonstrate the effectiveness of deep learning in capturing complex patterns in unstructured log data, which traditional approaches fail to address. Luo et al. (2020) introduced MoniLog, an

unsupervised learning system capable of real-time anomaly detection in cloud environments, showcasing the growing need for more

sophisticated and adaptive solutions. Meng et al. (2019) further highlight the effectiveness of unsupervised methods like LogAnomaly, which can detect both sequential and quantitative anomalies in large volumes of log data.

From early rule-based systems to contemporary machine learning-driven techniques, the development of anomaly detection techniques is thoroughly examined in this work. The main approaches used in text mining for anomaly detection will be covered, along with the tech's drawbacks and difficulties and the possible effects of implementing sophisticated anomaly detection methods on cloud infrastructure security. Through an analysis of the shortcomings and potential of existing approaches, this work seeks to advance knowledge of how simple text mining techniques might improve cloud-based systems' security and dependability.

## **Literature Review**

Anomaly detection in cloud systems has evolved significantly over time. Initially, anomaly detection relied on rule-based systems and statistical models to flag deviations in system behaviour. Early research primarily focused on using pattern-matching techniques and statistical baselines, which were often inadequate for handling the growing complexity and volume of cloud log data.

Modern anomaly detection has its roots in machine learning developments, which moved the emphasis to automated, data-driven techniques. The benefits of employing machine learning to find intricate patterns in unstructured log data have been highlighted by recent research that have investigated deep learning techniques for log-based anomaly detection, such Li et al. (2022) (Li et al., 2022). MoniLog, an automated method for identifying irregularities in cloud infrastructures, was presented by Luo et al. (2020). It uses unsupervised learning to detect anomalies in real time. LogAnomaly, a method introduced by Meng et al. (2019), is effective in real-world log settings and employs unsupervised algorithms to detect both sequential and quantitative anomalies.

Additional research, including that by Alzoubi et al. (2024), emphasizes the need for scalable and explainable anomaly detection models to enhance cloud security. The evolution of techniques has also seen the adoption of advanced text mining methods, such as transformer-based models, which offer a more nuanced understanding of log data and improve anomaly detection capabilities (Alzoubi et al., 2024). Goldstein & Uchida (2016) provided a comparative evaluation of different unsupervised anomaly detection algorithms, highlighting the strengths and weaknesses of various approaches for handling multivariate cloud data (Goldstein & Uchida, 2016).

In more recent advancements, Luo et al. (2021) proposed a robust unsupervised anomaly detection framework (RUAD) using AutoEncoders and Gaussian Mixture Models (GMM) to adapt to different data types, which can enhance the accuracy of anomaly detection in cloud environments. This progression towards more adaptive and resilient models is crucial for the rapidly evolving landscape of cloud computing.

## 3. Technical Details

### 3.1 Protocols and Standards

Anomaly detection in cloud server logs adheres to specific protocols and standards to ensure uniform data processing. Logs are gathered using standardized formats, such as Common Event Format (CEF) and syslog standards, which help maintain consistency across different cloud platforms. These standardized formats are crucial for effective preprocessing and compatibility with various anomaly detection models (Le & Zhang, 2022).

### 3.2 Preprocessing Methodology

The raw, unstructured log data must be pre-processed as the first stage in anomaly detection. In order to facilitate successful analysis, this stage converts the data into an organized format. Preprocessing involves stop-word removal,

which eliminates unnecessary words, and tokenization, which converts log entries into tokens. The significance of particular phrases within the dataset is then captured by features extracted using techniques like Term Frequency-Inverse Document Frequency (TF-IDF) (Meng et al., 2022). The data is guaranteed to be ready for the feature extraction and anomaly detection stages that follow thanks to these pretreatment procedures.

### 3.3 Feature Extraction Techniques

To capture the semantic links in log data, feature extraction uses sophisticated techniques including transformer-based models like BERT and word embeddings (like Word2Vec and GloVe). A better contextual understanding of log messages is provided by transformers like BERT, which enhances the detection of anomalies, while word embeddings convert textual data into a numerical representation that machine learning models can interpret (Vervaet, 2023).

### 3.4 Algorithms for Anomaly Detection

After features are extracted, machine learning algorithms are used to detect anomalies. For datasets with few anomalies, One-Class Support Vector Machines (OCSVM) work well. Neural networks called autoencoders are used to detect deviations, which point to anomalies, and to restore typical patterns. TLong Short-Term Memory (LSTM) networks must be able to capture temporal dependencies in order to identify sequential anomalies in log data. Two examples of ensemble learning methods that integrate the benefits of many models to improve detection accuracy are Long Short-Term Memory (LSTM) and Convolutional Neural Networks (CNN) (Meng et al., 2022; Luo et al., 2021).

### 3.5 Explainable AI (XAI) Approaches

The decision-making process of the model is made more transparent by XAI techniques, which offer insights into the reasoning behind the flagging of particular log entries as anomalies. Building stakeholder trust in circumstances when misclassification could have serious repercussions depends on this transparency. IT teams may assess and respond appropriately to anomalies because XAI makes the anomaly detection process interpretable (Alzoubi et al., 2024).

### 3.6 Federated Learning for Privacy-Preserving Anomaly Detection

A possible method for improving anomaly detection while protecting data privacy is federated learning. Sensitive data is protected by federated learning, which trains models locally on various datasets without requiring raw data sharing. This approach is particularly valuable in cloud environments where privacy is a major concern. Federated

learning allows cloud service providers to collaboratively build robust anomaly detection models, improving their effectiveness while ensuring data privacy (Yasarathna & Munasinghe, 2020; Alzoubi et al., 2024).

## 4. Obstacles

While basic text mining techniques offer significant potential for anomaly detection, several risks, issues, and limitations hinder their effective deployment in cloud environments. One of the primary obstacles is the unstructured nature of log data, which requires substantial preprocessing to make it suitable for analysis. Logs are often generated in diverse formats, complicating data normalization and increasing the overall complexity of preprocessing (Le & Zhang, 2022).

Another significant issue is the imbalance in datasets, as anomalies are relatively rare compared to normal behaviour. This imbalance presents a major challenge for training machine learning models, as standard algorithms often struggle to effectively detect rare events amidst a high volume of normal data points (Goldstein & Uchida, 2016). Techniques such as oversampling, under sampling, and specialized anomaly-focused model adjustments are frequently employed to address this issue, but these methods may not always provide effective results in all scenarios.

Privacy concerns are another critical risk, particularly given the sensitive information often contained in server logs. To implement anomaly detection solutions that respect data privacy while maintaining their effectiveness, innovative techniques are necessary. One potential remedy is federated learning, which permits decentralized model training while maintaining local data, protecting privacy (Alzoubi et al., 2024). Federated learning does have certain drawbacks, though, like node coordination and communication overhead.

The computational cost of processing and analyzing the vast amounts of log data generated by cloud environments is another significant limitation. Anomaly detection models often require substantial computational resources, which makes deploying these solutions in real-time settings challenging. To address these issues, lightweight machine learning algorithms and scalable data processing frameworks are being developed to ensure that detection is both efficient and accurate (Yasarathna & Munasinghe, 2020).

Explainability is another critical limitation, especially when using complex deep learning models. Explainable AI (XAI) approaches must be used to boost confidence and openness, guarantee that the model's results are reliable and useful, and help stakeholders comprehend the reasoning behind the model. Notwithstanding their advantages, these models are usually regarded as "black-box" systems, which makes it challenging for interested parties to understand why a certain anomaly was discovered.

In summary, although basic text mining techniques provide an effective approach for anomaly detection in cloud server logs, addressing the unstructured nature of log data, managing data imbalance, ensuring privacy, minimizing computational costs, and improving model explainability are all significant challenges that must be tackled for successful implementation in real-world scenarios. Another major challenge is the inherent imbalance in the dataset, as anomalies are rare compared to normal activities. This imbalance complicates the training of machine learning models, as most algorithms struggle to effectively identify rare events amidst a large number of normal instances (Goldstein & Uchida, 2016). Techniques such as oversampling, undersampling, and anomaly-specific model tuning are often necessary to address this issue but may not always yield satisfactory results in all scenarios.

Privacy concerns also pose a significant barrier, especially considering the sensitive information frequently stored in server logs. Ensuring that anomaly detection solutions respect data privacy without compromising their effectiveness requires innovative techniques. One possible remedy is federated learning, which permits decentralized model training while maintaining local data, protecting privacy (Alzoubi et al., 2024). Federated learning does, however, have certain drawbacks, including node coordination and communication complexity.

The computational cost involved in processing and analyzing large volumes of log data is another significant limitation. Cloud environments generate massive amounts of log data, and the real-time nature of anomaly detection makes it challenging to deploy resource-intensive algorithms. To mitigate these issues, lightweight machine learning algorithms and scalable data processing frameworks are being developed to enhance the efficiency of detection without compromising accuracy (Yasarathna & Munasinghe, 2020).

Another critical challenge is the explainability of deep learning-based anomaly detection models. Although these models are quite effective, they are frequently seen as "black-box" systems, which makes it hard for stakeholders to comprehend why an anomaly was reported (Luo et al., 2021). According to Alzoubi et al. (2024), explainable AI (XAI) techniques are therefore essential for enhancing trust and transparency in these systems, particularly in sectors like healthcare and finance where comprehension of the model's decision is critical for compliance and operational integrity.

## 5. The Promise

Basic text mining can uncover anomalies in cloud server logs, potentially improving the reliability, security, and overall resilience of cloud-based systems. By automating the detection of anomalous or suspicious behavior within logs, these strategies not only minimize the need for manual monitoring, but also allow for shorter response times to developing threats. This improvement in real-time detection can mitigate the impact of security breaches, thereby contributing to more robust data protection practices.

### 5.1 Transformative Potential in Public Health

The transformative potential of anomaly detection extends beyond IT infrastructure and into areas such as public health. For example, healthcare cloud systems that store patient records and other sensitive data can benefit from automated anomaly detection by ensuring unauthorized access attempts are quickly identified and mitigated. By safeguarding patient data, anomaly detection technologies help maintain trust in cloud-based healthcare solutions, supporting efforts to digitize health records and improve data accessibility for medical professionals.

### 5.2 Industrial Benefits and Societal Trust

From an industrial standpoint, anomaly detection systems help avoid costly disruptions to cloud services that many organizations rely on for everyday operations. Organizations can reduce downtime and improve productivity and financial outcomes by recognizing and responding to system faults in real-time. Furthermore, the application of explainable AI (XAI) in these systems can assist bridge the gap between technical anomaly detection approaches

and stakeholder comprehension, increasing trust in automated systems.

Anomaly detection also holds societal value in ensuring the security and privacy of user data, which is particularly important as more services and personal data are moved to the cloud. By incorporating privacy-preserving technologies like federated learning, cloud systems can collectively benefit from improved anomaly detection models while maintaining individual data privacy. This approach is especially valuable in environments like financial services, where data sensitivity is paramount, and breaches can have widespread societal repercussions.

Overall, the adoption of advanced anomaly detection techniques in cloud computing infrastructures promises to bolster security, enhance system reliability, and foster trust among both end-users and industry stakeholders, ultimately contributing to the stable growth of cloud-based services in society.

## 6. Suggested Course of Action

To maximize the potential of text mining in anomaly detection, several actions are recommended. First, cloud service providers should standardize log formats to reduce the complexity of preprocessing. Standardized formats facilitate easier integration of different anomaly detection systems and promote consistency in data handling across platforms. Second, to protect data privacy while enabling efficient anomaly detection, privacy-preserving strategies like federated learning should be used. Federated learning improves security while preserving user confidence by enabling collaboration between many cloud environments without disclosing private information (Alzoubi et al., 2024).

Third, algorithms that are scalable and lightweight must be created in order to effectively manage the massive amounts of log data. This involves investigating the application of hybrid models, which blend the advantages of many machine learning approaches. For example, ensemble learning techniques that incorporate CNNs, LSTMs, and attention mechanisms to enhance detection accuracy.

Training IT staff to use cutting-edge text mining and machine learning techniques for anomaly identification should also be funded. Working together, academia and industry can push this field's research forward and produce

creative solutions that overcome existing constraints. Enhancing the interpretability of anomaly detection systems through the creation of explainable AI tools should be a top priority in order to promote increased adoption and trust.

## 7. Conclusion

Simple text mining methods provide a workable way to find anomalies in cloud server logs, which aids in efficiently identifying operational problems and security risks. Even if there are still issues with data imbalance, computing expense, and privacy, new developments in machine learning, transformer models, federated learning, and explainable

AI offer encouraging answers. The cloud industry may improve its anomaly detection skills and guarantee the dependability and security of cloud infrastructures by implementing standardized log formats, privacy-preserving techniques, and lightweight algorithms.

Incorporating cutting-edge methods like explainable AI and federated learning improves anomaly detection accuracy while also fostering stakeholder trust, which is essential for wider adoption. These solutions will be crucial in protecting the dependability and integrity of cloud-based systems as cloud environments continue to get more complex.

## 8.References

1. **Le, V.-H., & Zhang, H. (2022). Log-Based Anomaly Detection with Deep Learning: How Far Are We?** *International Conference on Software Engineering (ICSE)*. https://doi.org/10.1145/3510003.3510155

2. **Vervaet, A. (2023). MoniLog: An Automated Log-Based Anomaly Detection System for Cloud Computing Infrastructures.** *Institut Supérieur d'Électronique de Paris*.

3. **Meng, W., Liu, Y., Zhu, Y., Zhang, S., Pei, D., & Liu, Y. (2022). LogAnomaly: Unsupervised Detection of Sequential and Quantitative Anomalies in Unstructured Logs.** *Proceedings of the International Joint Conference on Artificial Intelligence*.

4. **Goldstein, M., & Uchida, S. (2016). A Comparative Evaluation of Unsupervised Anomaly Detection Algorithms for Multivariate Data.** *PLOS ONE, 11*(4), e0152173. **https://doi.org/10.1371/journal.pone.0152173**

5. **Yasarathna, T. L., & Munasinghe, L. (2020). Anomaly Detection in Cloud Network Data. Smart Computing and Systems Engineering. University of Kelaniya, Sri Lanka.**

6. **Luo, Z., He, K., & Yu, Z. (2021). A Robust Unsupervised Anomaly Detection Framework.** *Applied Intelligence, 52*, 6022–6036.

   **https://doi.org/10.1007/s10489-021-02736-1**

7. **Alzoubi, Y. I., Mishra, A., & Topcu, A. E. (2024). Research Trends in Deep Learning and Machine Learning for Cloud Computing Security.** *Artificial Intelligence Review, 57*, 132.

   **https://doi.org/10.1007/s10462-024-10776-5**

8. **Zhang, B., Zhang, H., & Moscato, P. (2020). Anomaly Detection via Mining Numerical Workflow Relations from Logs. IEEE Symposium on Reliable Distributed Systems (SRDS).**

9. **Jayaweera, M. P. G. K., Kithulwatta, W. M. C. J. T., & Rathnayaka, R. M. K. T. (2023). Detect Anomalies in Cloud Platforms by Using Network Data: A Review.** *Cluster Computing, 26*, 3279–3289. **https://doi.org/10.1007/s10586-023-04055-1**

10. **Shahzad, F., Mannan, A., Javed, A. R., Almadhor, A. S., Baker, T., & Al-Jumeily, D. (2022). Cloud-Based Multiclass Anomaly Detection and Categorization Using Ensemble Learning.** *Journal of Cloud Computing, 11*(74).

    **https://doi.org/10.1186/s13677-022-00329-y**