

Developing an Automated Anomaly Detection System for Network Logs

LITERATURE SURVEY

submitted by

Dachepalli Leela Kesava Trinadh

21BAI1579

Nidumolu Sairam Gopal

21BRS1459

Rishi Patri

21BPS1396

of

B. Tech. Computer Science and Engineering

School of Computer Science and Engineering



VIT[®]
Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

October 2024

S.No	Title	Methodology	Takeaways
1	Adanomaly: Adaptive Anomaly Detection for System Logs with Adversarial Learning	<p>BiGAN Model: Uses Bidirectional GANs for feature extraction, enhancing detection accuracy.</p> <p>Ensemble Learning: Combines multiple classifiers to address class imbalance and reduce hyperparameter reliance.</p> <p>Log Parsing: Applies the Drain method to convert raw log data into accurate templates.</p> <p>Feature Extraction: Extracts features from log sequences with BiGAN and balances data for classifier training.</p>	<p>Adanomaly Framework: Introduces a novel log-based anomaly detection method using BiGAN and ensemble learning to improve accuracy and manage class imbalance.</p> <p>Feature Extraction: BiGAN derives features through reconstruction and discriminative losses, enhancing detection precision.</p> <p>Experimental Results: Shows superior recall and accuracy compared to six baseline methods across three public datasets.</p>
2	Anomaly Detection on Servers Using Log Analysis	<p>Log Parsing: Uses the Drain algorithm to structure raw log data.</p> <p>Feature Extraction:</p> <ul style="list-style-type: none">Event Count/TF-IDF: Compiles counts and applies TF-IDF.Sliding Window Counts: Creates matrices from event sequences.Final Matrix: Merges sliding window matrices with TF-IDF vectors. <p>Anomaly Detection Model: Implements a CNN model trained on labeled log data.</p>	<p>Deep Learning Model: CNN achieved up to 99% accuracy in anomaly detection.</p> <p>Log Parsing: Used the Drain algorithm for structuring raw log data.</p> <p>Feature Extraction: Employed TF-IDF and Sliding Window Event Counts.</p> <p>Experimental Results: High performance with low error rates and high F1-scores.</p>
3	LogST: Log Semi-supervised Anomaly Detection Based on Sentence-BERT	<p>Log Parsing: Converts raw logs into structured templates using the Drain method.</p> <p>Semantic Embedding: Uses Sentence-BERT (SBERT) for semantic representations of log events.</p> <p>Clustering: Applies HDBSCAN for clustering log sequences based on semantics.</p> <p>Anomaly Detection: Implements a GRU neural network with semi-supervised learning for anomaly detection.</p>	<p>LogST: Semi-supervised log anomaly detection using SBERT and GRU.</p> <p>Improved Accuracy: Outperforms traditional methods through semantic relationships.</p> <p>Stability with Few Labels: Effective with limited labeled normal logs.</p> <p>Experimental Validation: Shows significant improvements on the HDFS dataset.</p>
4	Machine Learning to Detect Anomalies in Web Log Analysis	<p>Two-Level ML Algorithm: Decision tree for classification and HMMs for anomaly detection.</p> <p>Feature Extraction: Extracts HTTP status codes, URL length, and parameter counts from logs.</p> <p>Data Labeling: Automatically labels logs to identify attacks versus normal behavior.</p> <p>Performance Evaluation: Uses accuracy, precision, FPR, and TPR, comparing to Logistic Regression and SVM.</p>	<p>Anomaly Detection System: Uses a two-level algorithm with a Decision Tree and HMM for web log detection.</p> <p>High Accuracy: Achieved 93.54% accuracy and 4.09% false positive rate, outperforming Logistic Regression and SVM.</p> <p>Real-World Data: Tested on data from a real industrial environment.</p> <p>Future Improvements: Plans to add a retraining module for adapting to new attack patterns.</p>

5	Log-based Anomaly Detection Without Log Parsing	<p>Data Splitting: 80% training and 20% testing, with unseen log messages in the test set.</p> <p>Sliding Window: 20-message length with a step size of 1 for log sequences.</p> <p>Comparison of Methods: NeuralLog compared to SVM, LR, IM, LogRobust, and Log2Vec.</p> <p>Evaluation Metrics: Uses Precision, Recall, and F1-score across datasets (HDFS, BGL, Thunderbird, Spirit).</p>	<p>NeuralLog Approach: Transforms raw logs into semantic vectors using BERT, bypassing parsing.</p> <p>Results: Achieves F1-scores over 0.95 on four datasets, outperforming existing methods.</p> <p>Contributions: Highlights limitations of log parsing and NeuralLog's effectiveness with OOV words.</p>
6	A Study on Log Anomaly Detection using Deep Learning Techniques	<p>Feature Extraction: Uses TF-IDF and Word2vec to convert log data into dense vectors.</p> <p>Machine Learning: Employs algorithms like SVM, Decision Tree, and PCA for anomaly detection.</p> <p>Deep Learning: Utilizes LSTM, RNN, Autoencoder, and Bi-LSTM for advanced anomaly detection.</p>	<p>Importance of Anomaly Detection: Essential for maintaining reliability and performance in large-scale networked systems.</p> <p>Deep Learning Techniques: Utilizes models like RNN, LSTM & autoencoders for effective log anomaly detection.</p> <p>Challenges: Addresses complications from unstructured data, log instability, and log bursts.</p>
7	A Comprehensive Review of Anomaly Detection in Web Logs	<p>Rule-based Models: Detect known anomalies but rely on administrator expertise.</p> <p>Statistical Models: Use Regex for anomaly detection based on query parameters.</p> <p>Supervised & Unsupervised ML Models: Supervised methods use labeled data; unsupervised methods employ clustering for unknown anomalies.</p> <p>Deep Hybrid Models (DHM): Combine log sequence encoding and machine learning for semantic anomaly detection.</p>	<p>Focus on Web Logs: Reviews techniques for detecting anomalies in HTTP logs.</p> <p>Categorization: Classifies methods into rule-based, statistical, supervised, unsupervised, and deep hybrid models.</p> <p>Challenges: Discusses issues like high-dimensional data and adaptive thresholds.</p> <p>Applications: Highlights use in cybersecurity, including IDS and FDS.</p>
8			
9			
10			