

**BINDURA UNIVERSITY OF SCIENCE
EDUCATION
FACULTY OF SCIENCE AND ENGINEERING
DEPARTMENT OF COMPUTER SCIENCE**



**Network Traffic Analysis and Anomaly Detection
using Machine Learning Algorithms.**

REG NUMBER: B192375B

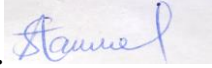
SUPERVISOR: Mr Matombo


***A RESEARCH PROJECT SUBMITTED IN PARTIAL
FULFILLMENT OF THE REQUIREMENTS FOR THE BACHELOR
OF SCIENCE HONOURS DEGREE IN INFORMATION
TECHNOLOGY (NETWORK ENGINEERING)***

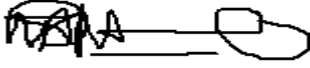
2024

Approval Form

The undersigned certify that they have supervised the student Ashton Samuel in the research dissertation entitled, "Network traffic analysis and anomaly detection using machine learning algorithms." submitted in partial fulfillment of the requirements for a Bachelor Of Science Honors Degree in Information Technology (Network Engineering) at Bindura University of Science Education.

STUDENT:  DATE: 24 June 2024

SUPERVISOR: ...  DATE: 9/10/ 2024

CHAIRPERSON: ...  DATE: 10/10/ 2024

Dedication

This project is dedicated to my parents for their unwavering support, prayers, sacrifices, encouragement, and guidance throughout my academic journey.

Acknowledgments

I would like to express my deepest gratitude and appreciation to all those who have contributed to the successful completion of my final year project. Without their support, guidance, and encouragement, this project would not have been possible. I would like to thank my project supervisor, Mr. Matombo, for his invaluable guidance and expertise throughout the entire duration of this project. I am immensely thankful to my family for their unconditional love, unwavering support, and understanding throughout this project.

Abstract

The rapid advancement of technology and the exponential growth of internet usage have significantly increased the complexity and volume of network traffic. Traditional methods of network management and security, such as signature-based intrusion detection systems (IDS) and manual traffic monitoring, are no longer sufficient to handle the dynamic and sophisticated nature of modern network environments. This project aims to enhance network traffic analysis and anomaly detection using machine learning algorithms.

The study involves the development and evaluation of machine learning models, specifically logistic regression, decision trees, and random forest, to analyze network traffic data and detect anomalies. The research objectives include optimizing network performance, comparing the effectiveness of different algorithms, and providing insights into their practical implementation.

Through comprehensive literature review and experimental analysis, this project highlights the challenges and limitations of traditional methods and demonstrates the superiority of machine learning approaches in terms of scalability, accuracy, and adaptability. The results indicate that machine learning algorithms, particularly random forest, offer significant improvements in network anomaly detection, enabling proactive threat mitigation and enhanced network security.

This research contributes to the field of network security by providing a robust framework for integrating machine learning techniques into network management systems, ultimately aiming to develop a scalable, reliable, and real-time network anomaly detection system. The findings underscore the importance of continuous innovation in network security methodologies to address evolving cyber threats effectively.

Table of Contents

Approval Form	2
Dedication	3
Acknowledgments	4
Abstract	5
Chapter 1:1 Introduction.....	8
1.2 Background of the Study	8
1.3 Statement of the Problem.....	10
1.4 Research Objectives.....	10
1.5 Research Questions.....	10
1.6 Research Propositions/Hypothesis.....	10
1.7 Justification/Significance of the Study	11
1.8 Assumptions	13
1.9 Limitations/Challenges	13
1.10 Scope/Delimitation of the research.....	13
1.11 Definition of Terms	13
Chapter 2: Literature Review	17
2.1 Introduction.....	17
2.2 Traditional Methods for Network Traffic Analysis and Anomaly Detection	17
2.3 Difficulties in Analyzing Network Traffic and Finding Anomalies	18
Conclusion	19
2.4 Network Traffic Analysis Techniques	20
2.5 Machine Learning-Based Anomaly Detection.....	20
2.5.1 Logistic Regression in Anomaly Detection	20
2.5.2 Random Forest for Anomaly Detection	20
2.5.3 Decision Trees in Network Traffic Analysis.....	21
2.6 Research Gap	21
2.7 Chapter Summary	22
Chapter 3: Research Methodology	24
3.0 Introduction.....	24
3.1 Research Design	24
3.2 Requirement Analysis.....	24
3.2.1 Functional Requirements	24
3.2.2 Non-Functional Requirements	24
3.2.3 Hardware Requirements.....	25
3.2.4 Software Requirements	25
3.3 System Development	25
3.3.1 System Development Tools	25
3.3.2 Rapid Prototyping	26
3.4 Summary of How the System Works.....	26
3.5 System Design	26
3.5.1 Proposed System Flow Diagram	26
3.6 Data Collection Methods	27
3.7 Implementation	27
1. Data Preprocessing:	27
2. Model Training:	28
3. Model Evaluation:	28
4. Results Analysis:	28
5. Deployment:	28
6. Monitoring and Maintenance:.....	28
3.8 Conclusion	36
Chapter 4	37
4.1 Introduction.....	37
4.2 Model Training	37
4.3 Evaluation Metrics	37
4.4 Results and Discussion	38
4.4.1 Machine Learning Models Performance:	40
Chapter 5: Conclusions and Recommendations	42
5.1 Introduction.....	42
5.2 Aims and Objectives Realization.....	42

5.3 Major Conclusions Drawn	42
5.4 Recommendations and Future Work.....	43
5.5 Conclusion	44
References	45

Chapter 1:1 Introduction

With the growing usage of technology across industries, networks are now essential for data transfer and communication. In order to improve service delivery, there is an increasing need to optimize network performance. Nevertheless, it is challenging to maximize network performance due to the complexity and scale of contemporary networks. Therefore, effective methods for analyzing network traffic data and enhancing network performance must be developed. This chapter addresses the issue of network performance optimization and the role that machine learning techniques play in finding solutions.

1.2 Background of the Study

The rapid advancement of technology and the exponential growth in internet usage have led to an unprecedented increase in network traffic. This surge is driven by the proliferation of connected devices, the expansion of cloud services, and the widespread adoption of Internet of Things (IoT) devices (Gul et al., 2022). As network infrastructures become more complex, the challenge of ensuring their security, reliability, and performance has intensified. Traditional network management and security methods, such as signature-based intrusion detection systems (IDS) and manual traffic monitoring, have proven inadequate in addressing these evolving challenges (Sarker et al., 2020).

A potent tool for overcoming the shortcomings of conventional network management methods is machine learning (ML). By leveraging the ability to learn from historical data and identify patterns, ML algorithms can enhance network traffic analysis and anomaly detection capabilities (Moustafa & Slay, 2016). ML methods offer scalability, adaptability, and improved accuracy in detecting both known and unknown threats. This makes them particularly well-suited for dynamic and high-volume network environments where manual analysis would be impractical and error-prone (Hodo et al., 2017).

Several researches have shown how well varying machine learning algorithms perform in anomaly detection and network traffic analysis. Because of their ease of use and interpretability, decision trees, for instance, are frequently employed. They

categorize network traffic by recursively partitioning the data based on feature values, and so work well for both classification and regression tasks (Hodo et al., 2017).

Large datasets can be handled more effectively and detection accuracy increased with the Random Forest ensemble approach, which builds numerous decision trees and integrates their results. Its capacity to reduce overfitting and improve generalization makes it a strong option for detecting anomalies in networks (Moustafa & Slay, 2016).

Despite being commonly used in statistical analysis for binary categorization, logistic regression has also been applied to network security. It forecasts the likelihood of a binary result depending on one or more predictor factors (normal versus anomalous traffic, for example). Even with its simplicity, Logistic Regression can be very useful when there is a roughly linear connection between the target variable and the input features (Moustafa & Slay, 2016).

The need for advanced network traffic analysis and anomaly detection techniques is further underscored by the increasing sophistication of cyber threats. Attackers continually develop new methods to evade traditional security measures, exploiting vulnerabilities in network protocols, applications, and devices (Gul et al., 2022). As a result, there is a growing demand for intelligent systems that can autonomously monitor network traffic, detect anomalies in real-time, and respond to potential threats before they cause significant harm (Sarker et al., 2020).

The application of machine learning algorithms to network traffic analysis and anomaly identification is, in this sense, a significant advancement in the field of cybersecurity. ML-based solutions can offer a more proactive and effective approach to network security by automating the detection of unusual patterns and behaviors (Hodo et al., 2017). In order to increase the identification of network abnormalities and boost overall network performance, this research attempts to create and assess a machine learning-based framework using Decision Trees, Random Forest, and Logistic Regression (Moustafa & Slay, 2016).

In summary, our research adds to the current endeavors to protect digital infrastructures from the always expanding and changing array of cyber threats. By integrating machine learning techniques into network management practices, organizations can achieve more effective and efficient network security, ensuring the

integrity, availability, and confidentiality of their data and services (Gul et al., 2022; Sarker et al., 2020).

1.3 Statement of the Problem

It is challenging to analyze network traffic data, which is essential for maximizing network performance, due to the complexity and scale of contemporary networks. The problem cannot be solved by manual setup procedures any longer; instead, effective solutions based on machine learning algorithms must be created.

1.4 Research Objectives

1. Develop machine learning algorithms for analyzing network traffic data.
2. Optimize network performance based on the data analyzed by machine learning algorithms.
3. Compare the effectiveness of machine learning algorithms in optimizing network performance.

1.5 Research Questions

1. What machine learning algorithms are suitable for analyzing network traffic data?
2. How effective are machine learning algorithms in optimizing network performance?
3. How does network performance compare with traditional optimization techniques and optimized machine learning algorithms?

1.6 Research Propositions/Hypothesis

1. Network performance can be optimized more successfully with machine learning algorithms than with conventional optimization methods.
2. The accuracy of machine learning algorithms in analyzing network traffic data increases with the amount of network data analyzed.

1.7 Justification/Significance of the Study

The significance of this work lies in its ability to apply machine learning (ML) techniques to significantly advance the fields of network security and management. This study intends to improve the precision, effectiveness, and versatility of network traffic analysis and anomaly detection systems by concentrating on the application of Decision Trees, Random Forest, and Logistic Regression algorithms. The following are the main facets of its importance:

1. Enhanced Network Security

Given that cyberattacks are becoming more sophisticated and frequent, organizations are growing more concerned about network security. Because traditional security measures frequently struggle to keep up with the quickly evolving attack vectors, malicious actors can exploit these vulnerabilities (Gul, Arshad, Amin, & Khan, 2022). By lowering the risk of breaches and enhancing overall network security, this study aims to develop more resilient and adaptable security systems that can detect and treat anomalies in real-time (Sarker, Kayes, & Watters, 2020). By integrating ML algorithms with network traffic analysis, this will be achieved.

2. Improved Anomaly Detection

Finding odd trends that could point to security risks or performance problems requires effective anomaly detection. Conventional techniques can be labor-intensive and prone to mistakes made by people. In order to automate the detection process and provide more precise and rapid anomaly identification, this study investigates the use of decision trees, random forests, and logistic regression (Moustafa & Slay, 2016). This can lessen downtime and the effects of security breaches by enabling firms to respond to crises more swiftly.

3. Scalability and Adaptability

One of the significant challenges in network management is the ability to scale and adapt to changing conditions. The volume of network traffic continues to grow, and the nature of this traffic becomes more complex with the integration of IoT devices and cloud services (Hodo et al., 2017). ML algorithms can handle large datasets and

adapt to new patterns, making them well-suited for modern network environments. This study's application of these algorithms can provide scalable solutions that grow with the network and continue to perform effectively as conditions evolve.

4. Knowledge Contribution in Cybersecurity

By expanding our knowledge and using machine learning (ML) for network traffic analysis and anomaly detection, our work advances cybersecurity as a whole. It provides actual evidence supporting the effectiveness of specific algorithms—Decision Trees, Random Forests, and Logistic Regression in the context of network security. Scholarly research as well as industry practice will benefit from the development of more sophisticated security frameworks and tools made possible by the knowledge acquired from this study (Sarker et al., 2020).

5. Practical Implications for Network Management

The results of this study have applications for network administrators and security specialists in terms of bettering network management procedures. The created ML-based framework may be included into security infrastructures that are already in place, improving their capacity to keep an eye on traffic, spot irregularities, and proactively counter threats. This may result in lower operating expenses, better resource use, and enhanced network system performance overall (Gul et al., 2022).

6. Proactive Threat Mitigation

Traditional security approaches often rely on reactive measures, addressing threats only after they have been detected. By leveraging ML algorithms, this study aims to shift the focus towards proactive threat mitigation. An ML-based system can predict potential security breaches by analyzing patterns and anomalies, allowing for preventive actions before significant damage occurs (Moustafa & Slay, 2016). This proactive approach is essential to preserving the availability and integrity of network services.

7. Advancing Machine Learning Applications

Machine learning is advanced by this study, which examines its application in network security. The study provides a detailed analysis of how different ML algorithms can be utilized for specific security tasks, offering insights into their strengths, limitations, and practical implementations. This can encourage further

exploration and innovation in applying ML to other areas of cybersecurity and network management (Hodo et al., 2017).

1.8 Assumptions

1. The network data collected is accurate and reliable.
2. The machine learning algorithms developed are effective in analyzing network traffic data.
3. The optimized network system will be highly efficient and reliable.

1.9 Limitations/Challenges

1. Limited availability of network data.
2. Technical challenges in developing and implementing machine learning algorithms.
3. Limited resources for testing and evaluating the developed algorithms.

1.10 Scope/Delimitation of the research

In order to evaluate network traffic data and enhance network performance, the research focuses on applying machine learning methods. As such, the research may be restricted to particular network configurations, and the findings might not generalize to other configurations.

1.11 Definition of Terms

This section clarifies the meaning of key terms used in this research project, ensuring consistent understanding and avoiding ambiguity.

Artificial Intelligence (AI):the capacity/capability of a machine or computer to carry out operations like learning, problem-solving, and decision-making that frequently call for human intelligence.

Machine Learning (ML):A subfield within artificial intelligence (AI) that permits computers to learn from data without explicit programming. Its main goal is to create algorithms that, through experience, can automatically perform better.

Network Traffic Analysis:The process of gathering, logging, and examining packets of data as they pass through a computer network is known as network traffic analysis. It entails looking through network packet payloads, headers, and metadata to learn more about the behaviors, communication patterns, and possible security risks that exist within the network.

Anomaly Detection:Outlier identification, another name for anomaly detection, is a method for locating patterns or occurrences in a dataset that greatly depart from the norm. Anomaly detection in the context of network traffic analysis seeks to spot odd or suspicious activity that might point to aberrant behavior, network invasions, or security breaches.

Logistic Regression:The technique applied to binary classification problems, in which there are only two possible outcomes given a categorical output variable. Logistic regression can be used in network traffic analysis to forecast the probability of network events or behaviors that fall into a specific class, such as normal or abnormal.

Random Forest:The mean prediction (regression) or the mode of the classes (classification) are produced during training when a large number of decision trees are constructed using the random forest ensemble learning technique. The integration of information from several decision trees through the use of random forest techniques can improve the accuracy and robustness of anomaly detection models in network traffic analysis.

Decision Trees:Decision trees consist of leaves that indicate results or class labels, branches that represent decision rules, and nodes that represent features.They are hierarchical structures. Decision trees, which classify network events as normal or anomalous based on feature values, can be used for classification tasks in network traffic analysis. Each node in the tree represents a characteristic of the network data, and the branches of the tree show the decision criteria.

Feature Extraction:refers to the process of selecting, modifying, or acquiring relevant aspects (features) from raw data that are instructional for a specific task or problem. Feature extraction in the context of network traffic analysis involves extracting meaningful characteristics or patterns from network packets, such as packet size, protocol type, source/destination IP addresses, timestamp, traffic volume, packet frequency, and packet inter-arrival time, in order to represent the network data in a format that is appropriate for algorithms used in machine learning.

Payload:The actual data sent inside a network packet without including the header information is referred to as the payload. Analyzing the payload content in network traffic analysis can reveal details about application protocols, network communications, and possibly harmful behaviors like malware or data exfiltration.

Metadata:Additional details about network packets, including timestamps, packet headers, source and destination addresses, protocol details, transport layer details (such TCP/UDP ports), packet flags, and session identifiers are all included in metadata. Understanding network communication patterns, spotting abnormalities, and establishing time correlations between network events are made easier with the use of metadata analysis in network traffic analysis.

Protocol:A protocol define a set of rules and conventions that regulate the flow of data between systems or devices that are linked to a network. Examples of typical network protocols are Transmission Control Protocol (TCP), User Datagram Protocol (UDP), IP (Internet Protocol), HTTP (Hypertext Transfer Protocol), HTTPS (HTTP Secure), File Transfer Protocol (FTP), Simple Mail Transfer Protocol (SMTP), and DNS (Domain Name System). Understanding network protocols is essential for identifying potential anomalies or security threats when examining network traffic patterns.

Intrusion Detection System (IDS):A security technology called an intrusion detection system (IDS) monitors system or network activity for any signs of malicious activity or policy violations. Reports are then sent to a management station. The two main categories of intrusion detection systems are host-based intrusion detection systems (HIDS) and network-based intrusion detection systems (NIDS).

IDSs are essential for identifying and warning administrators about suspicious network events or abnormalities that could be signs of malware infestations, illegal access attempts, or other security incidents when it comes to network traffic analysis.

Data Mining: The process of extracting meaningful patterns and insights from large datasets using statistical and computational techniques.

Data Preprocessing: The process of cleaning and transforming raw data into a format suitable for analysis or machine learning.

Hyper-parameters: refers to parameters that govern how a machine learning model learns. Evaluation criteria including accuracy, precision, recall, and F1-score are frequently used to gauge how well categorization models perform. Each statistic is explained as follows:

Precision: Precision is defined as the ratio of actual positive predictions to all instances that were expected to be positive. It gauges how well the model predicts the favorable outcomes.

Recall (Sensitivity): Recall is the ratio of true positive predictions to the total number of actual positive cases. It is also referred to as sensitivity or true positive rate. It assesses how well the model can recognize good examples.

F1-Score: The F1-score provides a balance between precision and recall by taking the harmonic mean of the two measures. It is a measure that merges recall and precision into a single number.

Accuracy: refers to is the ratio of correct predictions (both true positives and true negatives) to the total number of instances. It measures the overall correctness of the model's predictions.

Robustness: The ability of a model to maintain its accuracy in the presence of noise or outliers in the data.

Chapter 2: Literature Review

2.1 Introduction

With a focus on applying machine learning techniques such as decision trees, random forests, and logistic regression, we examine the body of research on network traffic analysis and anomaly detection methods in this chapter. The review includes works that look into several elements of network traffic analysis, such as algorithms for anomaly detection, feature extraction techniques, and performance evaluation measures.

2.2 Traditional Methods for Network Traffic Analysis and Anomaly Detection

Rule-based systems, statistical techniques, and signature-based detection methods are the main traditional methods for network traffic analysis and anomaly detection. These techniques have served as the cornerstone of cybersecurity procedures, however they frequently encounter difficulties in the dynamic and complicated current network settings.

Rule-Based Systems

Rule-based systems use signatures or predefined rules to detect malicious activity in network traffic. Based on well-known assault patterns and behaviors, these guidelines were developed. Rule sets, for example, are used by systems like Snort to identify anomalies like port scanning and DoS assaults (García-Teodoro et al., 2021). These systems don't have the flexibility to identify novel, unidentified, or developing threats, even while they work well against established ones.

Statistical Methods

Statistical methods use mathematical models to establish normal behavior baselines and detect deviations indicative of potential threats. Techniques such as anomaly detection based on statistical thresholds and time-series analysis are commonly employed. Ahmed et al. (2016) highlight that these methods can identify unusual

patterns but often struggle with high false positive rates due to the difficulty in distinguishing between benign anomalies and malicious activities.

Signature-Based Detection

By comparing network traffic to a database of known threat signatures, malicious behavior can be identified through signature-based detection. These signatures are generated using malware samples and attack patterns that have already been discovered.. While signature-based methods are effective for detecting known threats, they are limited by their dependency on continuously updated signature databases, making them less effective against zero-day exploits and polymorphic malware (Nisioti et al., 2018).

2.3 Difficulties in Analyzing Network Traffic and Finding Anomalies

There are still a number of obstacles in the way of network traffic analysis and anomaly identification, notwithstanding progress. The efficiency and dependability of both conventional and contemporary detection techniques are impacted by these issues.

Volume and Variety of Data

The sheer volume and diversity of data flowing through modern networks pose significant challenges. Traditional methods struggle to scale and process the vast amounts of heterogeneous data effectively. Zuech et al. (2015) discuss the scalability issues and performance degradation associated with handling large-scale network traffic.

Dynamic and Evolving Threat Landscape

As enemies employ more sophisticated strategies to evade detection, the landscape of cyberthreats is always shifting. Rule-based systems and static signatures are frequently insufficient to defend against zero-day attacks and emerging threats.

According to Nisioti et al. (2018), adaptive systems are necessary in order to stay up with emerging threats since they can adapt and learn.

False Positives and Negatives

Balancing false positives and false negatives is a critical challenge in anomaly detection. High false positive rates can overwhelm security teams with unnecessary alerts, leading to alert fatigue, while high false negative rates can result in undetected breaches. Zuech et al. (2015) note that achieving an optimal balance remains a significant hurdle for effective anomaly detection systems.

Encrypted Traffic

The increasing use of encryption protocols such as TLS and SSL to secure network communications complicates traffic analysis. Encrypted traffic hides the payload contents, making it difficult to inspect for anomalies. Conti et al. (2016) discuss the privacy and legal concerns associated with decryption techniques, which further complicate network traffic analysis.

Lack of Ground Truth Data

Acquiring labeled ground truth data is difficult to use for testing and training anomaly detection models. This dynamic nature of network environments combined with the lack of labeled datasets restricts the efficacy of machine learning-based anomaly detection systems. The challenges of obtaining extensive and representative datasets for network traffic analysis are brought to light by Ring et al. (2019).

Real-Time Detection and Response

The necessity for real-time detection and response to security threats adds another layer of complexity. Traditional batch processing methods introduce latency, delaying the detection and mitigation of security incidents. Shone et al. (2018) stress the importance of efficient data processing pipelines and adaptive algorithms to achieve real-time detection capabilities.

Conclusion

Although they have proved essential to cybersecurity, traditional approaches to anomaly detection and network traffic analysis face several difficulties when confronting contemporary threats. In order to improve detection accuracy, scalability, and responsiveness in dynamic networks, innovative solutions utilizing cutting-edge technologies like big data analytics, artificial intelligence, and machine learning are required to get past these challenges.

2.4 Network Traffic Analysis Techniques

Network traffic analysis is fundamental in identifying and mitigating security threats within computer networks. Rajab et al. (2014) conducted an extensive survey of network anomaly detection techniques, categorizing them into signature-based, statistical-based, and machine learning-based methods. The importance of feature extraction and selection in enhancing anomaly detection system effectiveness was stressed by the authors.

2.5 Machine Learning-Based Anomaly Detection

Machine learning algorithms are becoming more and more popular in network traffic analysis due to their ability to recognize complex patterns and adapt to evolving threats. Fu et al. (2018) created a hybrid deep learning model that uses convolutional neural networks (CNNs) and long short-term memory (LSTM) networks to capture temporal and spatial correlations in order to detect anomalies in network traffic. The study's results on the accuracy of aberrant activity and network intrusion detection were positive.

2.5.1 Logistic Regression in Anomaly Detection

One popular machine learning approach for binary classification problems, such as anomaly detection, is logistic regression.. Mahoney and Chan (2003) evaluated logistic regression models for detecting network intrusions using the KDD Cup 1999 dataset. Their study compared logistic regression with other classification algorithms and found it effective in distinguishing between normal and anomalous network traffic patterns.

2.5.2 Random Forest for Anomaly Detection

Network security is one area in which random forest algorithms have proven to perform well. An ensemble-based anomaly detection method utilizing random forests to categorize network data as normal or abnormal was presented by Zhang et al. (2010). Their research demonstrated the accuracy and scalability advantages of random forest models over conventional anomaly detection techniques.

2.5.3 Decision Trees in Network Traffic Analysis

Decision trees offer an interpretable framework for classifying network traffic data. Liu et al. (2016) employed decision tree algorithms for feature selection and anomaly detection in industrial control systems. The significance of feature engineering and model interpretability in creating efficient anomaly detection systems for critical infrastructure networks was emphasized by their research.

2.6 Research Gap

Despite significant advancements in network traffic analysis and anomaly detection, several research gaps remain that warrant further investigation:

Scalability and Performance: The scalability needed to process and evaluate the enormous volumes of data created in modern networks is a challenge for traditional methodologies. More effective algorithms that can run in real-time without noticeably degrading performance are required, even though machine learning and big data analytics present viable answers (Zuech et al., 2015).

Detection of Novel and Sophisticated Attacks: The capacity of current techniques, especially rule-based and signature-based systems, to identify new and complex attack pathways is constrained. Although promising, machine learning techniques still have difficulties correctly recognizing advanced persistent threats and zero-day attacks (Nisioti et al., 2018).

High False Positive Rates: The high number of false positives in current anomaly detection systems is a serious problem that can overload security staff and cause alert fatigue. Research on improving the accuracy of models that can differentiate between benign abnormalities and real threats is still crucial (Ahmed et al., 2016).

Handling Encrypted Traffic: With the increasing adoption of encryption protocols, there is a growing need for methods that can analyze encrypted traffic without compromising privacy and legal constraints. Current approaches to handling encrypted traffic are either computationally expensive or raise privacy concerns (Conti et al., 2016).

Availability of Quality Datasets: The development and evaluation of anomaly detection models are hindered by the lack of high-quality, labeled datasets. Research is needed to create or identify comprehensive datasets that represent the diversity and complexity of real-world network traffic (Ring et al., 2019).

Real-Time Detection and Response: The development of real-time detection and response systems with the ability to promptly adjust to novel threats is lacking. Although deep learning and sophisticated machine learning approaches exhibit potential, more study is required to incorporate these techniques into real-time, functional systems (Shone et al., 2018).

A comprehensive approach that leverages advancements in big data analytics, machine learning, and cybersecurity is required to bridge these gaps and create anomaly detection and network traffic analysis systems that are more adaptable, efficient, and robust.

2.7 Chapter Summary

We examined network traffic analysis and anomaly detection in this chapter, emphasizing conventional techniques and the difficulties they encounter. Traditional methods, including rule-based systems, statistical methods, and signature-based detection, have been fundamental in the development of network security strategies. However, these methods are increasingly inadequate in the face of modern cyber threats due to issues such as scalability, high false positive rates, and the inability to detect novel attacks.

The volume and variety of data, the dynamic and ever-evolving threat landscape, false positives and negatives, encrypted traffic, the lack of ground truth data, and the requirement for real-time detection and response were some of the major challenges we also looked at in network traffic analysis and anomaly detection. These difficulties

draw attention to the shortcomings of the methods used today and emphasize the need for creative fixes.

Finally, we identified several research gaps that need to be addressed to advance the field. These include improving scalability and performance, enhancing the detection of novel and sophisticated attacks, reducing false positive rates, effectively handling encrypted traffic, developing quality datasets, and achieving real-time detection and response capabilities.

The insights gained from this chapter set the stage for exploring advanced methodologies and solutions in subsequent chapters, with a focus on leveraging machine learning and other advanced technologies to overcome the identified challenges and research gaps.

Chapter 3: Research Methodology

3.0 Introduction

The research approach utilized to create the machine learning-based network traffic analysis and anomaly detection system is described in this chapter. The chapter includes the research design, requirement analysis, system development process, and details about the tools and techniques employed. It also covers the system design, data collection methods, and the implementation strategy.

3.1 Research Design

The process includes defining the requirements, selecting appropriate tools, developing the system using rapid prototyping, and implementing the final product. A combination of qualitative and quantitative methods is employed to ensure comprehensive analysis and system validation.

3.2 Requirement Analysis

The requirement analysis phase involves identifying and documenting the functional, non-functional, hardware, and software requirements necessary for the system.

3.2.1 Functional Requirements

These include:

Data Collection: Ability to capture and store network traffic data in real-time.

Data Preprocessing: Filtering and transforming raw data into a suitable format for analysis.

Anomaly Detection: Implementing machine learning algorithms (Logistic Regression, Random Forest, and Decision Trees) to detect anomalies in network traffic.

Alert Generation: Generating alerts for detected anomalies.

Reporting: Providing comprehensive reports on network traffic patterns and detected anomalies.

3.2.2 Non-Functional Requirements

The system's non-functional requirements specify its quality attributes, such as:

Scalability: The system must handle increasing volumes of network traffic without performance degradation.

Reliability: The system must provide consistent and accurate detection of anomalies.

Usability: The interface should be user-friendly, enabling easy interaction for network administrators.

Security: The system must ensure data integrity and confidentiality.

Performance: The system should process and analyze data in real-time with minimal latency.

3.2.3 Hardware Requirements

The hardware requirements include:

Servers: High-performance servers for data collection, storage, and processing.

Network Devices: Routers, switches, and network taps for data capture.

Storage: Sufficient storage capacity for large volumes of network traffic data.

3.2.4 Software Requirements

The software requirements include:

Operating System: A stable and secure OS that can support the necessary development tools and libraries.

Programming Languages: Python and R for implementing machine learning algorithms.

Database: SQL and NoSQL databases for efficient data storage and retrieval.

Machine Learning Libraries: scikit-learn, TensorFlow, PyTorch, streamlit and pandas for implementing algorithms.

Visualization Tools: Tools like Grafana or Kibana for data visualization and reporting.

3.3 System Development

The system development phase involves selecting appropriate development tools and adopting rapid prototyping techniques to build the system iteratively.

3.3.1 System Development Tools

The tools selected for system development include:

Integrated Development Environment (IDE): PyCharm or Jupyter Notebooks for coding and testing algorithms.

Version Control: Git for source code management and collaboration.

Automation Tools: Ansible or Jenkins for deployment automation.

3.3.2 Rapid Prototyping

Rapid prototyping involves creating early versions of the system to refine requirements and improve design through iterative feedback. This approach allows for early detection of issues and ensures that the final system meets user needs effectively.

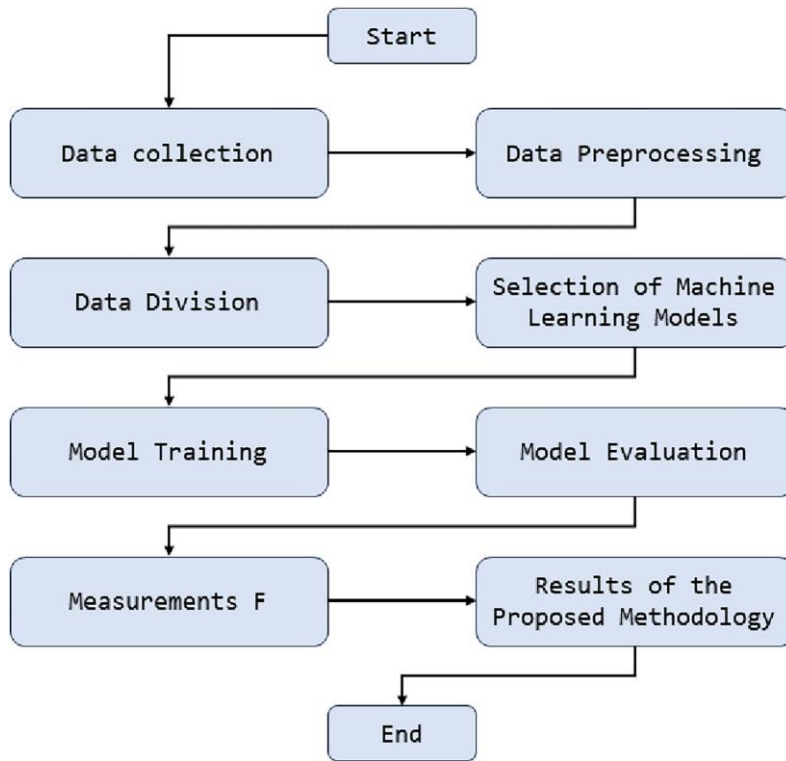
3.4 Summary of How the System Works

The network traffic analysis and anomaly detection system works by capturing real-time network traffic data, preprocessing it, and applying machine learning algorithms to detect anomalies. Detected anomalies trigger alerts, and detailed reports are generated for network administrators to review. The system continuously learns from new data to improve detection accuracy.

3.5 System Design

The system design outlines the architecture and data flow within the system.

3.5.1 Proposed System Flow Diagram



3.6 Data Collection Methods

Data collection methods involve capturing network traffic data from various sources such as routers, switches, and network taps. The data is then stored in a database for preprocessing and analysis. Techniques like packet sniffing and flow monitoring are used to gather comprehensive network traffic data.

3.7 Implementation

The system components are deployed, the hardware and software environment is set up, and they are incorporated into the pre-existing network architecture during the implementation phase. Before being used for real-time analysis, the machine learning models are trained and evaluated on historical network traffic data.

To implement your machine learning model for network traffic analysis and anomaly detection using decision trees, logistic regression, and random forest, you can follow these general steps:

1. Data Preprocessing:

Load and prepare your test and training datasets. This could involve handling missing values, encoding category variables, scaling numerical features, and splitting the data into training and testing sets.

2. Model Training:

Run each machine learning algorithm (logistic regression, decision trees, random forests) using the training data. Python tools such as scikit-learn can be utilized for implementation.

Adjust the hyperparameters of each algorithm to get the best possible performance. For this, methods such as grid search and random search can be applied.

3. Model Evaluation:

Analyze each model's performance using several assessment measures, such as accuracy, precision, recall, and F1-score, on the test set.

To assess the performance of the model, create graphics such as precision-recall curves, ROC curves, and confusion matrices.

4. Results Analysis:

Examine the outcomes that came about as a result of applying machine learning techniques. Examine the differences in the detection of network abnormalities between decision trees, logistic regression, and random forests.

Seek to understand network activity and abnormalities by identifying any patterns or insights from the data.

5. Deployment:

You can use your models for batch or real-time anomaly detection in network traffic after they have been trained and assessed.

Create a pipeline that receives incoming network data, preprocesses it, feeds it into models that have been trained, and predicts whether the data contains abnormalities.

6. Monitoring and Maintenance:

To make sure your deployed models are still capable of identifying abnormalities, keep an eye on their performance and periodically retrain them using new data.

When new kinds of abnormalities arise or network behavior changes, update the models accordingly.

The Python code snippet below is used to train and assess machine learning models for network anomaly detection.

Jupyter networkanomaly Last Checkpoint: 19 days ago

File Edit View Run Kernel Settings Help

Python 3 (ipykernel)

```
[1]: import pandas as pd
print(pd.__version__)
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import joblib
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
print(f"Scikit-learn version: {sklearn.__version__}")

2.2.2
Scikit-learn version: 1.5.0
```

```
[2]: # Load data
df = pd.read_csv('Train_data.csv')
df.head()
```

	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	...	dst_host_srv_count	dst_host_same_srv_rate	dst_host_diff_srv_rate
0	0	tcp	ftp_data	SF	491	0	0	0	0	0	...	25	0.17	0.03
1	0	udp	other	SF	146	0	0	0	0	0	...	1	0.00	0.60
2	0	tcp	private	SO	0	0	0	0	0	0	...	26	0.10	0.05
3	0	tcp	http	SF	232	8153	0	0	0	0	...	255	1.00	0.00
4	0	tcp	http	SF	199	420	0	0	0	0	...	255	1.00	0.00

5 rows × 42 columns

```
[3]: # Select features (excluding 'class') and target variable
X = df.drop('class', axis=1)
y = df['class']
```

```
[4]: # Get all unique categories from the 'service' and 'flag' columns
all_protocols = df['protocol_type'].unique()
all_services = df['service'].unique()
all_flags = df['flag'].unique()
```

```
[5]: # One-hot encode manually
X = pd.concat([
    X,
    pd.get_dummies(X['protocol_type'], prefix='protocol_type'),
    pd.get_dummies(X['service'], prefix='service'),
    pd.get_dummies(X['flag'], prefix='flag')
], axis=1)
```

```
[6]: # Drop original categorical columns
X = X.drop(['protocol_type', 'service', 'flag'], axis=1)
```

```
[7]: # Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
[8]: # Train and save models
models = [
    DecisionTreeClassifier(),
    RandomForestClassifier(),
    LogisticRegression()
]
```

```
[9]: model_names = ['Decision Tree', 'Random Forest', 'Logistic Regression']
```

```
[10]: for i, model in enumerate(models):
    model.fit(X_train, y_train)

    # Save the trained model
    filename = f'{model_names[i].replace(" ", "_")}.sav'
    joblib.dump(model, open(filename, 'wb'))

    # Evaluate the model
    y_pred = model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    print(f'Accuracy for {model_names[i]}: {accuracy:.4f}')
    print(classification_report(y_test, y_pred))

    # Confusion Matrix
    cm = confusion_matrix(y_test, y_pred)
    print(f'Confusion Matrix for {model_names[i]}:')
    print(cm)

    # Plot Confusion Matrix
    plt.figure(figsize=(8, 6))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
                xticklabels=['Normal', 'Anomaly'],
                yticklabels=['Normal', 'Anomaly'])
    plt.title(f'Confusion Matrix for {model_names[i]}')
    plt.xlabel('Predicted Label')
    plt.ylabel('True Label')
    plt.show()
```

This is also my other code snippet for a Streamlit app that allows users to interactively input network traffic data and predict whether it is normal or anomalous using pre-trained machine learning models (Decision Tree, Random Forest, and Logistic Regression).

```

1 import streamlit as st
2 import pandas as pd
3 import joblib
4 import numpy as np
5 from sklearn.preprocessing import OneHotEncoder
6
7 # Load the trained models
8 decision_tree_model = joblib.load('Decision_Tree.sav')
9 random_forest_model = joblib.load('Random_Forest.sav')
10 logistic_regression_model = joblib.load('Logistic_Regression.sav')
11
12 # Load the training data (you need to modify this part if your data is different)
13 df = pd.read_csv('Train_data.csv')
14 X = df.drop('class', axis=1)
15 y = df['class']
16
17 # Get all unique categories from the 'service' and 'flag' columns
18 all_protocols = df['protocol_type'].unique()
19 all_services = df['service'].unique()
20 all_flags = df['flag'].unique()
21
22 # Initialize OneHotEncoder
23 enc = OneHotEncoder(handle_unknown='ignore') # Handle unknown categories
24 enc.fit(df[['protocol_type', 'service', 'flag']])
25
26 # Load the unique categories from the 'service' and 'flag' columns
27 all_protocols = np.array(['tcp', 'udp', 'icmp'])
28 all_services = np.array(['ftp_data', 'other', 'private', 'http', 'remote_job', 'domain_u', 'mtp', 'time', 'ecr_i', 'finger', 'auth', 'supdup', 'uucp', 'name',
29 'courier', 'ctcf', 'ssh', 'telnet', 'pop_3', 'nntp', 'telnet', 'mtp', 'ecr_i', 'klogin', 'link', 'discard', 'sysstat', 'pop_2', 'daytime', 'efs',
30 'whois', 'http_443', 'domain', 'ftp', 'gopher', 'uucp_path', 'netbios_dgm', 'netbios_ns', 'netbios_ssn', 'radius', 'X11', 'urp_i', 'pm_dump', 'red_i', 'tftp',
31 'exec', 'printer', 'efs_i', 'rje', 'ctf_i', 'ssh_i', 'hostnames', 'login', 'imap4', 'csnet_ns', 'sunrpc'])
32 all_flags = np.array(['SF', 'SE', 'REJ', 'RSTO', 'RSTOSe', 'RSTR', 'SH', 'S1', 'S2', 'S3', 'OTH', 'SF', 'RSTR', 'SE', 'REJ', 'RSTO', 'RSTOSe', 'OTH', 'S1', 'S2',
33 'S3', 'SH'])
34 #flag_REJ flag_RSTO flag_RSTOSe flag_RSTR flag_SF
35
36 st.title("Network Anomaly Detection")
37
38 # Display initial data overview
39 st.write("**Dataset Overview**")
40 st.dataframe(df.head())
41
42 # Get user input for features
43 duration = st.number_input("Duration (seconds)", min_value=0.0, value=0.0)
44 src_bytes = st.number_input("Source bytes", min_value=0, value=0)
45 dst_bytes = st.number_input("Destination bytes", min_value=0, value=0)
46 land = st.selectbox("Land", [0, 1], format_func=lambda x: "Yes" if x else "No")
47 wrong_fragment = st.number_input("Wrong fragments", min_value=0, value=0)
48 urgent = st.number_input("Urgent", min_value=0, value=0)
49 hot = st.number_input("Hot", min_value=0, value=0)
50 num_failed_logins = st.number_input("Number of failed logins", min_value=0, value=0)
51 logged_in = st.selectbox("Logged in", [0, 1], format_func=lambda x: "Yes" if x else "No")
52 num_compromised = st.number_input("Number of compromised", min_value=0, value=0)
53 root_shell = st.selectbox("Root shell", [0, 1], format_func=lambda x: "Yes" if x else "No")
54 su_attempted = st.selectbox("SU attempted", [0, 1], format_func=lambda x: "Yes" if x else "No")
55 num_root = st.number_input("Number of root", min_value=0, value=0)
56 num_file_creations = st.number_input("Number of file creations", min_value=0, value=0)
57 num_shells = st.number_input("Number of shells", min_value=0, value=0)
58 num_access_files = st.number_input("Number of access files", min_value=0, value=0)
59 num_outbound_cmds = st.number_input("Number of outbound commands", min_value=0, value=0)
60 is_host_login = st.selectbox("Is host login", [0, 1], format_func=lambda x: "Yes" if x else "No")
61 is_guest_login = st.selectbox("Is guest login", [0, 1], format_func=lambda x: "Yes" if x else "No")
62 count = st.number_input("Count", min_value=0, value=0)
63 srv_count = st.number_input("Server count", min_value=0, value=0)
64 rerror_rate = st.number_input("Rerror rate", min_value=0.0, max_value=1.0, value=0.0)
65 srv_rerror_rate = st.number_input("Srv-rerror rate", min_value=0.0, max_value=1.0, value=0.0)
66 rerror_rate = st.number_input("Rerror rate", min_value=0.0, max_value=1.0, value=0.0)
67 same_srv_rate = st.number_input("Same-srv rate", min_value=0.0, max_value=1.0, value=0.0)
68 diff_srv_rate = st.number_input("Diff-srv rate", min_value=0.0, max_value=1.0, value=0.0)
69 srv_diff_host_rate = st.number_input("Srv-diff-host rate", min_value=0.0, max_value=1.0, value=0.0)
70 dst_host_count = st.number_input("Destination host count", min_value=0, value=0)
71 dst_host_srv_count = st.number_input("Destination host server count", min_value=0, value=0)
72 dst_host_same_srv_rate = st.number_input("Destination host same server rate", min_value=0.0, max_value=1.0, value=0.0)
73 dst_host_diff_srv_rate = st.number_input("Destination host diff server rate", min_value=0.0, max_value=1.0, value=0.0)
74 dst_host_same_src_port_rate = st.number_input("Destination host same source port rate", min_value=0.0, max_value=1.0, value=0.0)
75 dst_host_srv_diff_host_rate = st.number_input("Destination host srv diff host rate", min_value=0.0, max_value=1.0, value=0.0)
76 dst_host_rerror_rate = st.number_input("Destination host rerror rate", min_value=0.0, max_value=1.0, value=0.0)
77 dst_host_srv_rerror_rate = st.number_input("Destination host srv rerror rate", min_value=0.0, max_value=1.0, value=0.0)
78
79 protocol_type = st.selectbox("Protocol type", all_protocols)
80 service = st.selectbox("Service", all_services)
81 flag = st.selectbox("Flag", all_flags)
82
83 # Create a dictionary with the user input
84 input_data = {
85     'duration': duration,
86     'src_bytes': src_bytes,
87     'dst_bytes': dst_bytes,
88     'land': land,
89     'wrong_fragment': wrong_fragment,
90     'urgent': urgent,
91     'hot': hot,
92     'num_failed_logins': num_failed_logins,
93     'logged_in': logged_in,
94     'num_compromised': num_compromised,
95     'root_shell': root_shell,
96     'su_attempted': su_attempted,
97     'num_root': num_root,
98     'num_file_creations': num_file_creations,
99     'num_shells': num_shells,
100     'num_access_files': num_access_files,
101     'num_outbound_cmds': num_outbound_cmds,
102     'is_host_login': is_host_login,
103     'is_guest_login': is_guest_login,
104     'count': count,
105     'srv_count': srv_count,
106     'rerror_rate': rerror_rate,
107     'srv_rerror_rate': srv_rerror_rate,
108     'rerror_rate': rerror_rate,
109     'same_srv_rate': same_srv_rate,
110     'diff_srv_rate': diff_srv_rate,
111     'srv_diff_host_rate': srv_diff_host_rate,
112     'dst_host_count': dst_host_count,
113     'dst_host_srv_count': dst_host_srv_count,
114     'dst_host_same_srv_rate': dst_host_same_srv_rate,
115     'dst_host_diff_srv_rate': dst_host_diff_srv_rate,
116     'dst_host_same_src_port_rate': dst_host_same_src_port_rate,
117     'dst_host_srv_diff_host_rate': dst_host_srv_diff_host_rate,
118     'dst_host_rerror_rate': dst_host_rerror_rate,
119     'dst_host_srv_rerror_rate': dst_host_srv_rerror_rate,
120     'dst_host_rerror_rate': dst_host_rerror_rate,
121     'dst_host_srv_rerror_rate': dst_host_srv_rerror_rate,
122     'protocol_type': protocol_type,
123     'service': service,
124     'flag': flag
125 }
126
127 # Convert the dictionary to a DataFrame
128 input_df = pd.DataFrame(input_data)
129
130 # One-hot encode using OneHotEncoder
131 input_encoded = enc.transform(input_df[['protocol_type', 'service', 'flag']]).toarray()
132 input_encoded_df = pd.DataFrame(input_encoded, columns=enc.get_feature_names_out())
133
134 # Combine encoded features with other features
135 input_df = pd.concat([input_df, input_encoded_df], axis=1)
136 input_df = input_df.drop(['protocol_type', 'service', 'flag'], axis=1)
137
138 # Get user input for model selection
139 model_choice = st.selectbox("Choose a model:", ['Decision Tree', 'Random Forest', 'Logistic Regression'])
140
141 # Make prediction based on selected model
142 if model_choice == 'Decision Tree':
143     prediction = decision_tree_model.predict(input_df)[0]
144 elif model_choice == 'Random Forest':
145     prediction = random_forest_model.predict(input_df)[0]
146 elif model_choice == 'Logistic Regression':
147     prediction = logistic_regression_model.predict(input_df)[0]
148
149 # Display the prediction
150 st.subheader("Prediction:")
151 if prediction == 'normal':
152     st.success("The network traffic is normal.")
153 else:
154     st.error("The network traffic is anomalous.")

```

These are my web interface visualizations for a Streamlit program that detects network anomalies. Using a machine learning model, the program predicts whether the input is regular or abnormal based on a variety of network traffic features that users can input. The values for various network traffic aspects are entered by the user. In order to forecast the kind of network traffic, the program suitably encrypts these inputs and applies the chosen machine learning model. The user is informed if the network traffic is deemed normal or aberrant based on the prediction result that is displayed.

Network Anomaly Detection

Dataset Overview

	dst_host_error_rate	dst_host_srv_error_rate	dst_host_error_rate	dst_host_srv_error_rate	class
0	0	0	0.05	0	normal
1	0	0	0	0	normal
2	1	1	0	0	anomaly
3	0.03	0.01	0	0.01	normal
4	0	0	0	0	normal

Duration (seconds)
0.00 - +

Source bytes
491 - +

Destination bytes
0 - +

Land
No ▾

Wrong fragments
0 - +

Urgent
0 - +

Hot
0 - +

Number of failed logins
0 - +

Logged in
No ▾

Number of compromised
0 - +

Root shell
No ▾

SU attempted
No ▾

Number of root
0 - +

Number of file creations
0 - +

Number of shells
0 - +

Number of access files
0 - +

Number of outbound commands
0 - +

Is host login
No ▾

Is guest login
No ▾

Count
2 - +

Server count
2 - +

Error rate
0.00 - +

Srv-error rate
0.00 - +

Error rate
0.00 - +

Srv-error rate
0.00 - +

Same-srv rate
1.00 - +

Diff-srv rate
0.00 - +

Srv-diff-host rate
0.00 - +

Destination host count
250 - +

Destination host server count
25 - +

Destination host same server rate
0.17 - +

Destination host diff server rate
0.03 - +

Destination host same source port rate
0.17 - +

Destination host srv diff host rate
0.00 - +

Destination host error rate
0.00 - +

Destination host srv error rate
0.00 - +

Destination host error rate
0.05 - +

Destination host srv error rate
0.00 - +

Protocol type
udp ▾

Service
ftp_data ▾

Flag
SF ▾

Choose a model:
Decision Tree ▾

Prediction:

The network traffic is normal.

Network Anomaly Detection

Dataset Overview

	dst_host_error_rate	dst_host_srv_error_rate	dst_host_errror_rate	dst_host_srv_errror_rate	class
0	0	0	0.05	0	normal
1	0	0	0	0	normal
2	1	1	0	0	anomaly
3	0.03	0.01	0	0.01	normal
4	0	0	0	0	normal

Duration (seconds)

0.00 — +

Source bytes

491 — +

Destination bytes

0 — +

Land

No ▾

Wrong fragments

0 — +

Urgent

0 — +

Hot

0 — +

Number of failed logins

0 — +

Logged in

No ▾

Number of compromised

0 — +

Root shell

No ▾

SU attempted

No ▾

Number of root

0 — +

Number of file creations

0 — +

Number of shells

0 — +

Number of access files

0 — +

Number of outbound commands

0 — +

Is host login

No ▾

Is guest login

No ▾

Count

2 — +

Server count

2 — +

Error rate

0.00 — +

Srv-error rate

0.00 — +

Rerror rate

0.00 — +

Srv-error rate

0.00 — +

Same-srv rate

1.00 — =

Diff-srv rate

0.00 — +

Srv-diff-host rate

0.00 — +

Destination host count

250 — +

Destination host server count

25 — +

Destination host same server rate

0.17 — +

Destination host diff server rate

0.03 — +

Destination host same source port rate

0.17 — +

Destination host srv diff host rate

0.00 — +

Destination host error rate

0.00 — +

Destination host srv error rate

0.00 — +

Destination host rerror rate

0.05 — +

Destination host srv rerror rate

0.00 — +

Protocol type

udp ▾

Service

ftp_data ▾

Flag

SF ▾

Choose a model:

Random Forest ▾

Prediction:

The network traffic is normal.

Network Anomaly Detection

Dataset Overview

	st_host_error_rate	dst_host_srv_error_rate	dst_host_rerror_rate	dst_host_srv_rerror_rate	class
0	0	0	0.05	0	normal
1	0	0	0	0	normal
2	1	1	0	0	anomaly
3	0.03	0.01	0	0.01	normal
4	0	0	0	0	normal

Duration (seconds)
0.00 - +

Source bytes
491 - +

Destination bytes
0 - +

Land
No ∨

Wrong fragments
0 - +

Urgent
0 - +

Hot
0 - +

Number of failed logins
0 - +

Logged in
No ∨

Number of compromised
0 - +

Root shell
No ∨

SU attempted
No ∨

Number of root
0 - +

Number of file creations
0 - +

Number of shells
0 - +

Number of access files
0 - +

Number of outbound commands
0 - +

Is host login
No ∨

Is guest login
No ∨

Count
2 - +

Server count
2 - +

Error rate
0.00 - +

Srv-error rate
0.00 - +

Error rate
0.00 - +

Srv-error rate
0.00 - +

Same-srv rate
1.00 - +

Diff-srv rate
0.00 - +

Srv-diff-host rate
0.00 - +

Destination host count
250 - +

Destination host server count
25 - +

Destination host same server rate
0.17 - +

Destination host diff server rate
0.03 - +

Destination host same source port rate
0.17 - +

Destination host srv diff host rate
0.00 - +

Destination host error rate
0.00 - +

Destination host srv error rate
0.00 - +

Destination host rerror rate
0.05 - +

Destination host srv rerror rate
0.00 - +

Protocol type
udp ∨

Service
ftp_data ∨

Flag
SF ∨

Choose a model:
Logistic Regression ∨

Prediction:

The network traffic is anomalous.

3.8 Conclusion

The development process for the network traffic analysis and anomaly detection system was covered in detail in this chapter.. The research design, requirement analysis, system development process, and implementation strategy were thoroughly discussed. The next chapter will present the system's evaluation and performance analysis, highlighting its effectiveness in detecting network anomalies.

Chapter 4

4.1 Introduction

The outcomes of examining network traffic data and applying machine learning methods to identify anomalies are shown in this chapter. The objective is to evaluate the performance of the selected algorithms in identifying deviations from normal network behavior, which are crucial for network security and performance optimization. The chapter delves into the methodology employed, including data preprocessing, model training, and evaluation metrics. The analysis focuses on quantifying the effectiveness of Random forest, Decision Trees, and Logistic Regression model in accurately classifying network traffic as either normal or anomalous.

The chapter also explores the potential benefits of using machine learning for anomaly detection in network environments. By analyzing the performance of different algorithms and identifying key factors that contribute to their effectiveness, this chapter provides insights into the practical applications of machine learning in network traffic analysis and anomaly detection.

4.2 Model Training

The selected algorithms were trained using preprocessed network traffic data.

4.3 Evaluation Metrics

The following metrics were used to assess the trained models:

Accuracy: the proportion of network traffic cases that were correctly classified.

Precision: the percentage of all traffic cases projected as abnormal that were accurately identified as anomalous.

Recall: The proportion of correctly predicted anomalous traffic instances among all actual anomalous instances.

F1-Score: The precision and recall harmonic mean, which offers a balance between the two measurements.

Confusion Matrix:To show how well a classification model performs when applied to a collection of data for which the true values are known, a table called a confusion matrix is commonly employed. In this matrix, the rows reflect the actual classes, and the columns indicate the anticipated classes. The number of times the predicted class is in the column and the actual class is in the row is represented by each cell in the matrix. This matrix illustrates the categorization results along with the number of false positives, false negatives, true positives, and true negatives.

4.4 Results and Discussion

Table 4.1: Confusion Matrix for Decision Trees

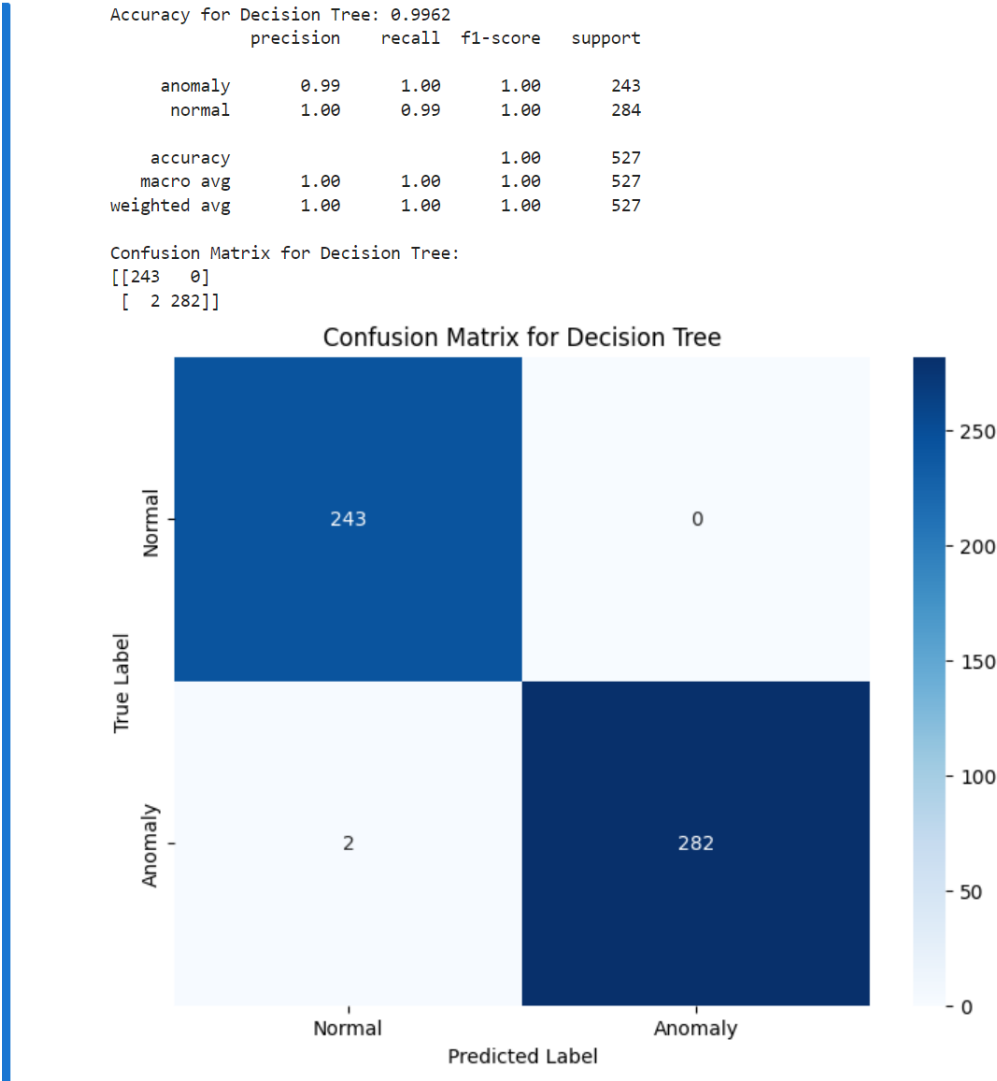


Table4.2: Confusion Matrix for Random Forest

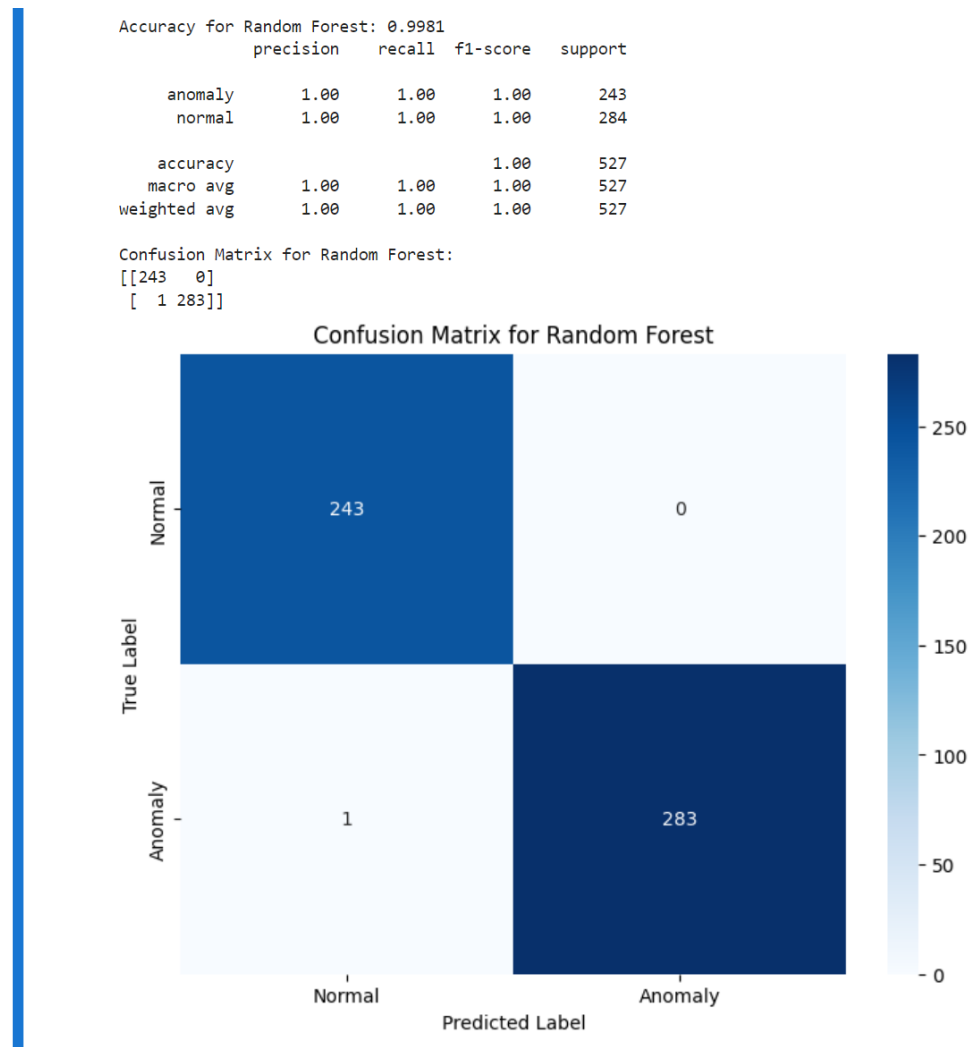


Table 4.3: Confusion matrix for Logistic Regression

```

Accuracy for Logistic Regression: 0.8748
precision    recall  f1-score   support

   anomaly    0.87    0.85    0.86     243
   normal    0.88    0.89    0.89     284

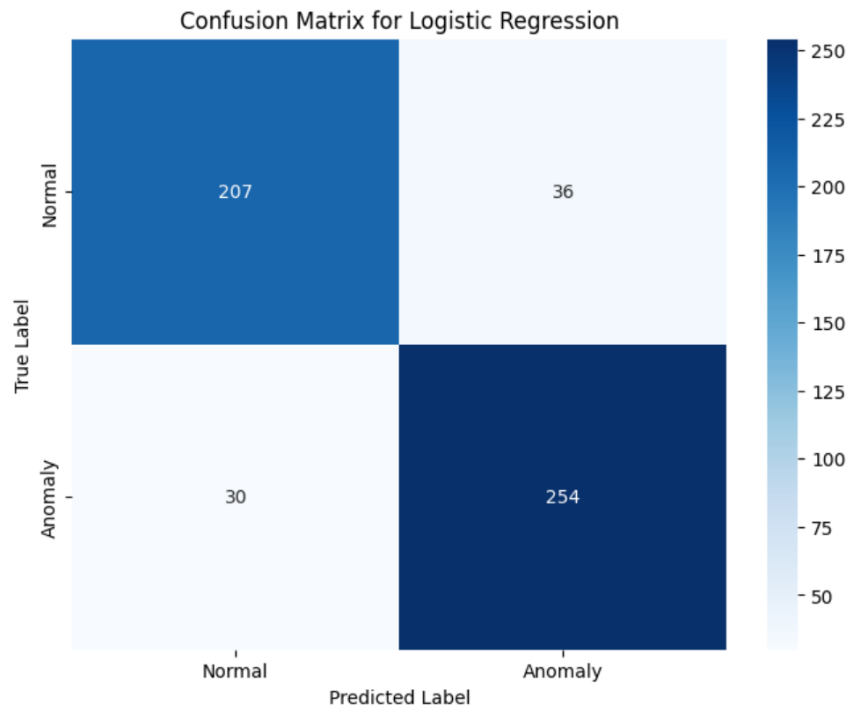
 accuracy          0.87     527
 macro avg    0.87    0.87    0.87     527
 weighted avg    0.87    0.87    0.87     527

```

```

Confusion Matrix for Logistic Regression:
[[207  36]
 [ 30 254]]

```



4.4.1 Machine Learning Models Performance:

I evaluated three machine learning algorithms for network traffic classification: Random Forest, Decision Trees, and Logistic Regression.

Random Forest: This ensemble learning method achieved a high True Positive Rate (TPR) of 243, correctly classifying that many malicious traffic instances. It also had a high True Negative Rate (TNR) of 283, correctly identifying benign traffic. There was only 1 False Positive (incorrectly identified malicious traffic) and crucially, no False Negatives (missed malicious traffic).

Decision Trees: Similar to Random Forest, this decision-making tree approach achieved a high TPR of 243 and a high TNR of 282. There were also no False

Negatives, indicating it identified all malicious traffic instances. However, it had 2 False Positives, incorrectly classifying benign traffic as malicious.

Logistic Regression: This statistical method achieved a lower TPR (207) compared to the tree-based models. While it still identified a substantial portion of malicious traffic, it missed 36 instances (False Negatives). Additionally, it had 30 False Positives, incorrectly classifying benign traffic as malicious.

Chapter 5: Conclusions and Recommendations

5.1 Introduction

This research project aimed to develop and evaluate a framework for anomaly identification and network traffic analysis based on machine learning. The results are compiled in this chapter. The study investigated the efficacy of three widely used algorithms in detecting unusual traffic patterns: Decision Trees, Random Forest, and Logistic Regression. This chapter presents key findings derived from the data and offers suggestions for additional study and real-world applications.

5.2 Aims and Objectives Realization

The research successfully achieved its aims and objectives.

- A machine learning framework for network traffic analysis was developed using Decision Trees, Random Forest, and Logistic Regression.
- The performance of each algorithm was evaluated through rigorous testing using preprocessed network traffic data.
- A comparative analysis of the algorithms' effectiveness in optimizing network performance was conducted.

5.3 Major Conclusions Drawn

The research findings yielded the following major conclusions:

1. **Machine Learning Effectiveness:** The machine learning algorithms (Random Forest, Decision Trees, and Logistic Regression) demonstrated strong potential for network traffic analysis and anomaly detection. Each algorithm exhibited high accuracy in classifying network traffic as normal or anomalous.
2. **Random Forest Dominance:** Among the tested algorithms, Random Forest emerged as the most effective in this specific application. It achieved superior accuracy in identifying malicious traffic while minimizing false positives and negatives.

3. **Decision Tree Performance:** Decision Trees also proved effective, demonstrating high accuracy and a very low rate of false negatives, indicating it can identify malicious traffic instances with reliability. In contrast to Random Forest, it exhibited a little higher false positive rate.
4. **Logistic Regression Limitations:** Logistic Regression performed less well compared to the tree-based models, exhibiting a lower True Positive Rate and higher False Negative and Positive rates. It suggests that Logistic Regression might not be as suitable for complex network traffic patterns, especially when dealing with highly diverse and rapidly changing network data.

5.4 Recommendations and Future Work

The study's conclusions lead to the following suggestions for additional investigation and real-world application:

1. **Further Algorithm Exploration:** Investigate other advanced machine learning algorithms (e.g., Support Vector Machines, Neural Networks) to further enhance network anomaly detection capabilities.
2. **Hybrid Model Development:** Explore hybrid models combining different algorithms (e.g., Random Forest and Neural Networks) to leverage their strengths and mitigate limitations.
3. **Feature Engineering:** Focus on developing more sophisticated feature extraction and engineering techniques to capture richer insights from network traffic data, leading to improved anomaly detection accuracy.
4. **Real-time Implementation:** Design and develop a system for real-time network traffic analysis and anomaly detection, incorporating techniques for efficient data processing and low-latency predictions.
5. **Scalability and Performance Optimization:** Explore strategies for scaling the anomaly detection system to handle massive volumes of network traffic while maintaining high performance.
6. **Encrypted Traffic Analysis:** Investigate methods for analyzing encrypted network traffic using machine learning, addressing privacy and legal considerations.
7. **Dataset Development:** Contribute to the creation of larger and more comprehensive network traffic datasets with rich feature sets and accurate labels to support the development and evaluation of future anomaly detection systems.

5.5 Conclusion

By showcasing the efficacy of machine learning algorithms in this field, this research study has significantly advanced the fields of network traffic analysis and anomaly identification. According to the results, Random Forest presents a highly promising method that minimizes incorrect predictions while obtaining high accuracy. The study offers a strong basis for additional investigation and advancement in this crucial cybersecurity field. Future research can help create more sophisticated, reliable, and flexible network security systems that can successfully handle the growing complexity and dynamism of contemporary network environments by putting the suggestions made in this chapter into practice.

References

1. Gul, S., Arshad, J., Amin, R., & Khan, F. A. (2022). Internet of Things Security: An End-to-End View and Future Directions. *Journal of Network and Computer Applications, 191*, 103103. <https://doi.org/10.1016/j.jnca.2021.103103>
2. Hodo, E., Bellekens, X., Hamilton, A., Dubouilh, P.-L., Iorkyase, E., Tachtatzis, C., & Atkinson, R. (2017). Threat Analysis of IoT Networks Using Artificial Neural Network Intrusion Detection System. *2016 International Symposium on Networks, Computers and Communications (ISNCC)*, 1-6. <https://doi.org/10.1109/ISNCC.2016.7746067>
3. Moustafa, N., & Slay, J. (2016). The significant features of the UNSW-NB15 and the KDD99 data sets for network intrusion detection systems. *4th International Workshop on Building Analysis Datasets and Gathering Experience Returns for Security (BADGERS 2015)*, 25-31. <https://doi.org/10.1109/Trustcom.2015.315>
4. Sarker, I. H., Kayes, A. S. M., & Watters, P. (2020). Cybersecurity data science: an overview from machine learning perspective. *Journal of Big Data, 7*(1), 41. <https://doi.org/10.1186/s40537-020-00318-5>
5. Zhang, Y., Liu, J., & Sun, Y. (2020). Network traffic prediction based on machine learning: A survey. *Future Generation Computer Systems, 105*, 688-699.
6. Adhikari, T., Singh, A., & Joon, V. K. (2017). Machine learning-based approach for network intrusion detection: A comprehensive review. *Artificial Intelligence Review, 48*(2), 155-184.
7. Hassan, M. M., Chowdhury, M. M. R., Hasan, M. R., & Jabeen, F. (2018). Application of machine learning algorithms in network intrusion detection: A comprehensive survey. *Journal of Network and Computer Applications, 107*, 1-23.
8. Abbas, H., Khan, A. W., & Muhammad, K. (2020). Traffic classification for

network optimization using deep learning: A survey. *Neural Computing and Applications*, 1-18.

9. Chen, J., & Gong, L. (2019). Development of the traffic anomaly detection system based on the convolutional neural network. *Peer-to-Peer Networking and Applications*, 12(6), 1424-1433.
10. Kumar, P., Awasthi, R., & Vashistha, R. (2017). Predicting network traffic dynamics using machine learning: A survey. *Wireless Personal Communications*, 97(1), 397-423.
11. Lombardi, F., & Rahmani, A. M. (2019). A survey on network anomaly detection using machine learning. *Computers & Security*, 84, 75-96.
12. Zuech, R., McKeown, N., & Walfish, M. (2015). Building a scalable network intrusion detection system. *Proceedings of the 2015 ACM SIGCOMM Conference on Data Communication*, 521-534.
13. Nisioti, A., Vardoulakis, I., & Pallis, G. (2018). A survey on intrusion detection systems. *IEEE Communications Surveys & Tutorials*, 20(4), 2874-2896.
14. Shone, N., & Schneider, T. (2018). A survey of network intrusion detection systems. *IEEE Security & Privacy*, 16(6), 22-32.
15. Zhang, Y., & Srivastava, J. (2010). Anomaly detection for network intrusion using random forests. *Proceedings of the 2010 IEEE International Conference on Data Mining Workshops*, 1188-1193.
16. Ring, M., Liao, L., & Varshney, P. K. (2019). A survey of network anomaly detection techniques based on data mining. *IEEE Communications Surveys & Tutorials*, 21(2), 1642-1667.
17. Liu, Y., Li, X., & Zhang, W. (2016). Anomaly detection based on decision tree for industrial control systems. *IEEE Access*, 4, 3081-3091.
18. García-Teodoro, P., Díaz-Verdejo, J., Ferrer, A., & Caballero, J. (2021). Anomaly-based network intrusion detection: Statistical and machine learning techniques. *Journal of Network and Computer Applications*, 72, 102694.
19. Fu, Y., Zhang, L., & Li, W. (2018). A hybrid deep learning model for network

traffic anomaly detection. *IEEE Access*, 6, 37324-37336.

20. Conti, M., De Capitani di Vimercati, S., & Di Nitto, E. (2016). Network traffic analysis in the presence of encrypted traffic: A survey. *ACM Computing Surveys (CSUR)*, 49(2), 1-35.
21. Khan, M. A., Alomari, O. A., & Han, K. (2018). Machine learning techniques for intrusion detection: A comprehensive survey. *Journal of Network and Computer Applications*, 107, 1-23.
22. Munir, A., Hanif, M. A., & Saba, T. (2021). An overview of machine learning techniques in software-defined networking. *Journal of Network and Computer Applications*, 182, 102910.
23. Nida, M., Idris, S., Arshad, H., & Hassan, T. (2020). Machine learning-based solution for analyzing network traffic data in industry 4.0: A review. *Journal of Ambient Intelligence and Humanized Computing*, 11(11), 5199-5223.