# Investigating the application of Machine Learning algorithms for threat detection and anomaly detection in Network Traffic

Rahul Guha[1], Manju Vyas[2], Geerija Lavania[3], Pankaj Kumar Sharma[4]

[1,2,3,4] *Assistant Professor, AI&DS, JECRC Foundation*

*Abstract— The risk of network breaches and cyberattacks is ongoing in today's interconnected world. As the Internet grows, cyberattacks are evolving quickly, and the state of cyber security is not promising (Sarker, 2021). In order to help detect and avoid these dangers, researchers recognized the challenge and resorted to machine learning techniques. Machine learning methods are ideal for detecting threats in network traffic because of their capacity to examine massive amounts of data and spot trends and anomalies.*

*Index terms: Machine learning, cyber security, network traffic, anomaly detection, and threat detection.*

## I. INTRODUCTION

By applying machine learning algorithms to network traffic data, it is possible to detect potential threats and anomalies in real-time, enabling organizations to take proactive measures to protect their networks and mitigate the risk of cyber-attacks [Aiyanyo et al., 2020].

This can be achieved by training machine learning models on historical network traffic data and leveraging their ability to learn from patterns and behaviours to identify deviations from normal activity.In addition to traditional network security measures, machine learning algorithms can provide an extra layer of defence by continuously monitoring network traffic and identifying suspicious or malicious activity. Furthermore, the use of machine learning can help overcome the limitations of traditional rule-based approaches, such as the risk of false positives and the inability to detect novel and evolving threats. By continuously analysing network traffic data, machine learning algorithms can adapt and learn from new attack patterns, enabling organizations to stay one step ahead of cybercriminals. The proliferation of cyber threats in modern network environments necessitates advanced techniques for threat detection and anomaly detection. Machine learning (ML) algorithms have emerged as a promising approach to tackle these challenges by leveraging the ability to learn patterns and behaviours from vast amounts of network traffic data. This paper presents an investigation into the application of machine learning algorithms for threat detection and anomaly detection in network traffic. We explore various machine learning techniques, their effectiveness, challenges, and future directions in the context of enhancing network security.

## II. REVIEW OF LITURATURE

Survey of Intrusion Detection Systems: Techniques, Challenges, and Applications (2020)

This comprehensive survey provides an overview of intrusion detection systems (IDS) techniques, including signature-based, anomaly-based, and hybrid approaches. It discusses the role of machine learning in enhancing IDS performance, particularly in detecting novel attacks and reducing false positives.

Machine Learning for Network Intrusion Detection: A Review (2019)

This review provides an overview of the application of various machine learning techniques such as decision trees, support vector machines, neural networks, and ensemble methods for network intrusion detection. It discusses the strengths and limitations of each approach and highlights the importance of feature selection and dataset imbalance.

Deep Learning for Anomaly Detection: A Survey (2018)

This survey focuses on deep learning techniques for anomaly detection in various domains, including network traffic. It discusses auto encoders, recurrent neural networks (RNNs), convolutional neural networks (CNNs), and generative adversarial networks (GANs) for detecting anomalies in network traffic.

Anomaly Detection in Network Traffic Based on Machine Learning Algorithms (2017)

This paper presents a comparative study of machine learning algorithms for anomaly detection in network traffic. It evaluates the performance of algorithms such as k-means clustering, one-class SVM, and Isolation

Forest on benchmark datasets and discusses their effectiveness in identifying different types of anomalies.

A Survey of Machine Learning Techniques for Cyber Security Intrusion Detection" (2016)

This survey covers traditional machine learning algorithms such as decision trees, k-nearest neighbors, and Naive Bayes, as well as more advanced techniques such as random forests, support vector machines, and neural networks for intrusion detection. It discusses their applicability, advantages, and challenges in the context of cyber security.

### III. PROBLEM DEFINITION

- Overview of the significance of threat detection and anomaly detection in network traffic.
- Introduction to machine learning algorithms and their potential in enhancing detection capabilities.
- Statement of the problem and its importance in network security.

### IV. RESEARCH OBJECTIVES

- Algorithm Evaluation: Evaluate the performance of various machine learning algorithms such as supervised (e.g., Random Forest, Support Vector Machines), unsupervised (e.g., k-means, DBSCAN), and semi-supervised (e.g., self-training, co-training) methods for detecting threats and anomalies in network traffic. Outline the goals and expected outcomes of the research.
- Feature Selection and Engineering: Investigate effective feature selection and engineering techniques tailored to network traffic data to enhance the performance of machine learning models. This may involve exploring different representations of network traffic data, such as packet-level features, flow-level features, or time-series representations.
- Data Collection and Preprocessing: Develop methodologies for collecting and preprocessing large-scale network traffic datasets, including techniques for data cleaning, normalization, and handling imbalanced classes to ensure the quality and reliability of the training data.
- Real-Time Detection: Explore approaches for real-time threat detection and anomaly detection in network traffic, considering the computational efficiency and scalability of machine learning models to handle high-speed network streams.

- Interpretability and Explainability: Address the interpretability and explainability challenges associated with machine learning-based network traffic analysis, by developing techniques to provide insights into the decision-making process of the models and facilitate human understanding and trust.
- Transfer Learning and Domain Adaptation: Explore transfer learning and domain adaptation techniques to leverage knowledge from related domains or datasets with different characteristics to improve the generalization and adaptability of machine learning models for network traffic analysis.
- Integration with Existing Systems: Investigate the integration of machine learning-based threat detection and anomaly detection systems with existing network security infrastructure, such as intrusion detection/prevention systems (IDS/IPS) and security information and event management (SIEM) systems, to enhance overall network defense capabilities.
- Evaluation Metrics and Benchmarks: Define appropriate evaluation metrics and benchmarks for assessing the performance of machine learning models for network traffic analysis, considering factors such as detection accuracy, false positive rate, detection latency, and scalability.
- Ethical and Privacy Considerations: Address ethical and privacy considerations associated with the deployment of machine learning-based network traffic analysis systems, including issues related to data privacy, transparency, fairness, and accountability.

### IV. MODEL SELECTION

Supervised Learning Models:

Random Forest: Effective for classification tasks, robust to overfitting, and capable of handling high-dimensional data with complex interactions.

Support Vector Machines (SVM): Suitable for binary classification tasks, particularly when dealing with high-dimensional data, and known for their ability to handle non-linear decision boundaries.

Gradient Boosting Machines (e.g., XGBoost, LightGBM): Effective for classification tasks, capable of handling large-scale datasets, and often provide state-of-the-art performance.

Neural Networks: Deep learning models such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs) can capture complex patterns in network traffic data, especially when dealing with sequential or time-series data.

Unsupervised Learning Models:

k-means Clustering: Useful for grouping network traffic data into clusters based on similarity, which can help in identifying anomalies or unusual patterns.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise): Suitable for identifying clusters of varying shapes and sizes in the data, robust to noise, and effective for outlier detection.

Auto encoders: Neural network architectures trained to learn compact representations of input data, useful for anomaly detection by reconstructing normal patterns and detecting deviations.

Semi-supervised Learning Models:

Self-training: Initially trained on labelled data and then iteratively refined by incorporating unlabeled data, useful when labelled data is scarce or expensive to obtain.

Co-training: Exploits the diversity of multiple views or feature sets of the data to improve model performance, particularly beneficial when labelled data is limited but multiple sources of unlabeled data are available.

Hybrid Models:

Ensemble Methods: Combining multiple base models (e.g., Random Forest, SVM, neural networks) to improve prediction accuracy and robustness.

Stacked Models: Integrating predictions from multiple base models as features for a meta-model, potentially enhancing the overall performance.

Anomaly Detection Models:

Isolation Forest: Efficient for detecting anomalies in high-dimensional data by isolating outliers in a random partitioning scheme.

One-Class SVM: Trained on normal instances only and capable of detecting deviations from the normal behavior, suitable for detecting rare or novel threats.

*A. Review Stage*

Investigating the application of Machine Learning (ML) algorithms for threat detection and anomaly detection in network traffic is an important stage in cybersecurity research. This stage involves understanding how various ML techniques can enhance the accuracy and efficiency of identifying malicious activities or deviations from normal network behavior.

*B. Problem Definition*

**Scope**: Clearly define the type of threats and anomalies you aim to detect. This could range from detecting specific attacks like Distributed Denial of Service (DDoS) or ransomware to identifying abnormal traffic patterns caused by unknown threats.

Data**:** Identify the network traffic data that will be used for analysis. This includes selecting data sources like flow logs, packet captures (PCAP), or Intrusion Detection System (IDS) alerts.

*C. Feature Engineering*

Feature Selection: Review relevant features of network traffic such as packet sizes, timestamps, IP addresses, and protocols that can help in identifying anomalies.

Data Preprocessing**:** Analyze the importance of data cleaning, normalization, and transformation to make it suitable for ML algorithms. Raw network traffic data often requires extensive preprocessing.

### III. MATH

Gaussian Mixture Model (GMM)**:** GMM assumes that the data points are generated from a mixture of several Gaussian distributions. The likelihood of a data point $x_i$ being generated by a mixture model is given by:

$$P(x_i) = \sum_{k=1}^{K} \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)$$

where:

- $K$ is the number of Gaussian distributions,
- $\pi_k$ is the mixing coefficient,
- $\mathcal{N}(x_i|\mu_k, \Sigma_k)$ is the Gaussian distribution with mean $\mu_k$ and covariance $\Sigma_k$.

K-Nearest Neighbors (k-NN)**:** In k-NN anomaly detection, a point is an outlier if its average distance to its k nearest neighbors is high. The distance for each point $x_i$ is:

$$d(x_i) = \frac{1}{k} \sum_{j=1}^{k} ||x_i - x_j||$$

where:
- $x_j$ are the k nearest neighbors of $x_i$,
- $||\cdot||$ is the distance metric (e.g., Euclidean distance).

## IV. Evaluation Metrics

Accuracy vs Precision/Recall: Prioritize metrics like precision, recall, and F1 score over accuracy because false positives and false negatives have different costs in threat detection.

Confusion Matrix: Helps in understanding the performance of the model by analyzing True Positives, False Positives, True Negatives, and False Negatives.

Receiver Operating Characteristic (ROC) and Area under Curve (AUC): Common methods to evaluate binary classifiers, useful in comparing models.

## V. HELPFUL HINTS

### A. Challenges and Limitations

Imbalanced Data: Network traffic often contains far more normal traffic than malicious activity, which can cause ML models to be biased toward normal behaviour. Handling this imbalance is crucial.

Real-Time Detection: ML algorithms must be optimized for real-time performance, which may require simplifying models or using more efficient algorithms.

Adversarial Attacks: Review how resilient your chosen ML model is to adversarial attacks where attackers modify their behaviour to evade detection.

Data Privacy: Ensure that the collection and use of network data comply with legal and ethical standards, especially concerning user privacy.

### B. Future Work

Hybrid Models: Explore the combination of multiple algorithms (e.g., ensemble models) or using a mix of rule-based and ML-based detection systems.

Transfer Learning: Investigate using pre-trained models or applying models from one domain of network security to another.

Continual Learning: As network behavior evolves, your models must adapt without being entirely retrained. This is especially important in dynamic threat environments.

### C. References

Thwaini, M. H. (2022). Anomaly Detection in Network Traffic using Machine Learning for Early Threat Detection. *Data and Metadata*, *1*, 34-34.

Mohammed, R., & Akay, F. Anomaly Detection In Network Traffic Using Machine Learning. *Cukurova University Journal of Natural and Applied Sciences*, *2*(3), 5-12.

Patel, K., Fogarty, J., Landay, J. A., & Harrison, B. (2008, April). Investigating statistical machine learning as a tool for software development. In Proceedings of the SIGCHI conference on human factors in computing systems (pp. 667-676).

Latif, S., Faria, F. D., Afsar, M. M., Esha, I. J., & Nandi, D. (2022). Investigation of machine learning algorithms for network intrusion detection. International Journal of Information Engineering and Electronic Business, 15(2), 1.

Shah, V. (2021). Machine learning algorithms for cybersecurity: Detecting and preventing threats. Revista Espanola de Documentacion Cientifica, 15(4), 42-66.

Le, D. C., Zincir-Heywood, N., & Heywood, M. I. (2020). Analyzing data granularity levels for insider threat detection using machine learning. IEEE Transactions on Network and Service Management, 17(1), 30-44.

Binbusayyis, A., Alaskar, H., Vaiyapuri, T., & Dinesh, M. J. T. J. O. S. (2022). An investigation and comparison of machine learning approaches for intrusion detection in IoMT network. The Journal of Supercomputing, 78(15), 17403-17422.

Akcay, S., & Breckon, T. (2022). Towards automatic threat detection: A survey of advances of deep learning within X-ray security imaging. Pattern Recognition, 122, 108245.

## VI. PUBLICATION PRINCIPLES

When publishing research on anomaly detection models for threat detection in network traffic, certain principles should guide the presentation, validation, and communication of your findings to ensure the publication meets scientific and ethical standards.

Below are key principles to follow:
1) Clarity in Problem Definition.
2) Rigorous Data Handling and Methodology.
3) Model Explanation and Justification.
4) Comprehensive Experimental Setup.
5) Robust Evaluation Metrics.
6) Security and Resilience.
7) Visualization and Interpretation.
8) Ethical Considerations and Impact.
9) Acknowledgment of Limitations.
10) Open Access and Code Availability.

## VII. CONCLUSION

Our results show that hybrid models and deep learning techniques like Auto encoders, Generative Adversarial Networks (GANs), and Recurrent Neural Networks (RNNs) perform well in detecting complex patterns and zero-day attacks, providing a significant improvement over traditional statistical and distance-based methods. Furthermore, graph-based models like Graph Neural Networks (GNNs) hold promise in capturing relationships in highly interconnected network traffic, although they require further exploration for scalability in real-time scenarios.

## APPENDIX

Description: Real-world dataset containing network traffic data, including both normal and malicious traffic generated by different attack scenarios.

Features: 80+ features such as flow duration, source/destination IP, packet size, protocol type, etc.

Purpose: Used for training and testing machine learning and deep learning models.

## ACKNOWLEDGMENT

## REFERENCES

[1] Alrowaily, M. (2020). Investigation of Machine Learning Algorithms for Intrusion Detection System in Cybersecurity.

[2] Al-Gethami, K. M., Al-Akhras, M. T., & Alawairdhi, M. (2021). Empirical evaluation of noise influence on supervised machine learning algorithms using intrusion detection datasets. *Security and Communication Networks*, *2021*(1), 8836057.

[3] Al-Gethami, K. M., Al-Akhras, M. T., & Alawairdhi, M. (2021). Empirical evaluation of noise influence on supervised machine learning algorithms using intrusion detection datasets. *Security and Communication Networks*, *2021*(1), 8836057.

[4] Zhang, J., & Lei, Y. (2022). Trend and Identification Analysis of Anti-investigation Behavior in Crime by Machine Learning Fusion Algorithm. *Wireless Communications and Mobile Computing*, *2022*(1), 1761154.