

Machine Learning Based IDS for Detecting DOS Attack in Networks

Bhagyashree
Dept.of.CS&E
PESITM
Shimogga,India
bbhagyashree719@gmail.com

Rakshitha G
Dept.of.CS&E PESITM
Shimogga,India
rakshithaghema@gmail.com

Ramyashree R S
Dept.of.CS&E
PESITM
Shimogga,India
ramyashreers193@gmail.com

Sneha N L
Dept.of.CS&E PESITM
Shimogga,India snehanlsnehanl@gmail.com

Dr. Chethan L S
Dept. of CS&E
PESITM
Shimogga,India
chethan.ls@pestrust.edu.in

Abstract- Intrusion Detection Systems (IDS) are critical for ensuring the security of fog networks, which bridge the gap between cloud computing and edge devices. This paper proposes a machine learning-based IDS to detect Denial of Service (DoS) attacks in fog networks. The methodology employs feature selection techniques like Genetic Algorithm (GA) and Correlation-Based Feature Selection (CFS) to improve classification accuracy and efficiency. The proposed system is evaluated on key metrics such as accuracy, precision, recall, and F1-score, demonstrating robust detection capabilities in a fog network environment.

Keywords- Fog networks, Intrusion Detection Systems (IDS), DoS attacks, Feature selection, Machine learning, Classification algorithms

I. Introduction

Modern networks underpin essential services across industries, enabling seamless communication and data exchange. However, their increasing complexity and reliance on interconnected systems have made them prime targets for cyberattacks. Among these threats, Denial of Service (DoS) attacks are particularly destructive, as they overwhelm network resources, rendering services inaccessible to legitimate users. Addressing these vulnerabilities requires innovative solutions capable of adapting to evolving attack strategies. [3],[5].

The financial and operational consequences of network disruptions caused by DoS attacks are significant, with organizations suffering from downtime, loss of data, and reputational damage. Traditional network defense mechanisms, such as rule-based Intrusion Detection Systems (IDS), have proven inadequate in the face of sophisticated attack patterns. These limitations underscore the urgent need for advanced anomaly detection methods that can efficiently process the

highdimensional data characteristic of modern network environments. [12],[13].

Among the myriad of cyber threats, Denial of Service (DoS) attacks stand out as one of the most disruptive. These attacks aim to overwhelm a network's resources, rendering it incapable of serving legitimate users. The consequences of such disruptions are far-reaching, ranging from financial losses and reputational damage to severe interruptions in critical services. For example, attacks on healthcare networks could delay life-saving operations, while breaches in financial systems could result in large-scale monetary losses. [4],[11].

The cost of cyberattacks has been escalating at an alarming rate. According to recent studies, the global economic impact of cybercrime is expected to exceed \$10 trillion annually by 2025, with a significant portion attributed to DoS attacks. Beyond monetary losses, these attacks erode user trust, disrupt public services, and can even pose risks to national security. Traditional IDS frameworks, while foundational, are constrained by their dependency on predefined signatures and static rules. These systems excel at detecting known threats but fail to adapt to emerging attack patterns or handle the dynamic nature of modern network traffic. Furthermore, the explosion of high-dimensional data in networks—characterized by diverse protocols, traffic patterns, and user behaviors—adds to the complexity. Handling this data with traditional methods often results in scalability challenges and high false-positive rates, limiting their efficacy in real-world scenarios. [1],[6].

In recent years, the advent of machine learning (ML) has transformed the landscape of network security. ML-based IDS frameworks analyze historical data to identify patterns indicative of malicious activity. Unlike signature-based systems, ML approaches can detect previously unseen threats by focusing on anomalies—behaviors that deviate from the norm. Supervised learning algorithms, such as Support Vector Machines (SVM) and Random Forests (RF), have been widely adopted for these tasks due to their robust classification capabilities. [7]. However, ML algorithms are not without limitations. Their performance depends heavily on the quality and volume of training data. Furthermore, as the size and complexity of network traffic grow, ML models face scalability challenges. Training and inference times increase exponentially

with the number of features, and maintaining model accuracy becomes a significant hurdle. Additionally, the highdimensional nature of network data often includes redundant or irrelevant features, further complicating the classification process.

II. Related Work

Extensive research has explored the application of ML to improve IDS effectiveness. Verma et al. evaluated various ML classifiers, highlighting RF and Decision Trees for their high accuracy and low false-positive rates in detecting DoS attacks. However, their reliance on outdated datasets limited real-world applicability. Khatib et al. examined resampling techniques to enhance model performance on imbalanced datasets, demonstrating the utility of SMOTE but underscoring challenges in achieving scalability. [3],[5]. Recent advancements include deep learning frameworks, such as those by Thamilarasu et al., which utilize neural networks for comprehensive anomaly detection. These methods achieve impressive accuracy but demand substantial computational resources, making them less viable for resource-constrained environments. [5],[15].

III. Literature Review

Intrusion Detection Systems have evolved significantly, shifting from basic rule-based systems to sophisticated MLdriven frameworks. Traditional IDS methods rely on static rules or known signatures, making them vulnerable to novel attack strategies. In contrast, anomaly-based IDS approaches analyze traffic patterns to detect irregularities, offering broader coverage against emerging threats.[6],[11]. Key datasets such as KDD99, NSL-KDD, and UNSW-NB15 have historically dominated IDS research. While comprehensive, these datasets are now considered outdated, failing to reflect the complexity of contemporary attack scenarios. Recent datasets like IoTID20 and Bot-IoT introduce real-time traffic simulations and diverse attack types, making them invaluable for modern IDS research.[1].

Feature selection plays a pivotal role in IDS development. Methods such as CFS and GA optimize the feature set, reducing computational overhead and improving detection accuracy. Studies indicate that combining feature selection with robust ML algorithms enhances system efficiency without compromising performance.[2]. Intrusion Detection Systems (IDS) have undergone significant evolution, transitioning from traditional rule-based mechanisms to sophisticated frameworks leveraging machine learning (ML). Early IDS approaches primarily relied on signature-based methods, which compared network traffic against predefined patterns of known threats. While effective for recognizing established attack vectors, these systems were limited by their inability to detect novel or evolving threats. This vulnerability paved the way for anomaly-based IDS, which identify irregularities in network behavior, offering broader coverage against emerging cyber threats. By analyzing traffic patterns and deviations from normal behavior, anomaly-based IDS have become essential for securing modern networks against complex and dynamic attack strategies. [2],[9].

Several studies emphasize the importance of datasets and feature selection in IDS research. IoTID20 and similar datasets capture contemporary attack patterns, enabling researchers to develop systems resilient to emerging threats. Algorithms like CFS and GA have demonstrated efficacy in identifying relevant features while reducing computational overhead.

Classifiers such as RF and DT remain popular due to their interpretability and high performance, though advanced methods like neural networks promise scalability for large datasets.[10],[14].

Problem Statement

Classical computational approaches are increasingly limited when processing large, high-dimensional genomic datasets, often requiring substantial resources and time. This project explores the use of QSVM and QNN models to enhance the efficiency and accuracy of cancer type and stage classification, addressing the computational limitations of traditional methods and leveraging quantum computing potential to handle complex, highdimensional gene expression data effectively [7], [12].

IV. Methodology

A. Data Cleaning

Dataset collection is the foundational stage of any machine learning-based network analysis. This stage involves gathering data representing network traffic, including normal and malicious activities such as Denial of Service (DoS) attacks. The datasets are typically sourced from publicly available repositories, synthetic generation tools, or live network environments. Popular datasets, such as NSL-KDD or CICIDS2017, are often used for benchmarking intrusion detection systems. The quality and diversity of the dataset significantly affect the performance of the intrusion detection system. To ensure comprehensive analysis, the dataset should encompass various attack types, traffic patterns, and network environments. Collected datasets often include features like packet size, protocol type, source and destination IP addresses, and timestamps, which are critical for training machine learning models. Additionally, ensuring balanced datasets is essential to avoid bias toward specific traffic patterns or attack scenarios. Challenges in this stage include ensuring data privacy, handling imbalanced data distribution, and capturing realistic network traffic. To overcome these issues, anonymization techniques, oversampling, and synthetic data generation may be applied. This stage ensures a robust foundation for subsequent steps, as the quality of the data directly influences the system's effectiveness.

B. Preprocessing

Preprocessing is a crucial stage where raw network traffic data is cleaned, transformed, and prepared for machine learning analysis. This stage involves several steps, including handling missing values, encoding categorical data, normalizing numerical features, and removing redundant or irrelevant data. The goal of preprocessing is to convert raw data into a format suitable for modeling. For example, network data often contains categorical variables such as protocol types (e.g., TCP, UDP)

that need to be encoded numerically. Normalization ensures that all features have a uniform scale, preventing certain features from dominating the training process. Noise and inconsistencies in the dataset, such as incomplete records or outliers, are identified and rectified. Additionally, duplicate entries and irrelevant attributes, such as timestamps or unique identifiers that do not contribute to attack detection, are removed to enhance model efficiency. By the end of this stage, the dataset is transformed into a structured, clean, and standardized form ready for feature selection.

3. Data Visualization

Data visualization provides insights into the dataset's structure and characteristics, enabling a better understanding of feature distributions, correlations, and patterns. Tools like histograms, scatter plots, and heatmaps are used to visualize the relationships between variables. For example, correlation heatmaps highlight dependencies between features, aiding in the identification of redundant or highly correlated attributes.

Visual analysis also helps identify imbalances in the dataset, such as disproportionate representation of attack types, which can impact model training. This stage ensures that preprocessing and feature selection are guided by data-driven insights, optimizing subsequent modeling efforts.

4. Parameter Selection (Feature Selection)

Feature selection is performed to identify the most relevant attributes for intrusion detection. By selecting only the most critical features, this step reduces computational complexity, prevents overfitting, and improves model performance. Techniques like Correlation-based Feature Selection (CFS) and Genetic Algorithm (GA) are applied in this stage.

CFS evaluates the relevance of features by analyzing their correlation with the target variable and minimizing redundancy among features. GA, inspired by evolutionary biology, identifies optimal feature subsets by iteratively optimizing a fitness function. These techniques ensure that only the most informative attributes are retained, resulting in faster training times and enhanced prediction accuracy.

5. Splitting Data

In this stage, the dataset is divided into training and testing subsets, typically in an 80:20 or 70:30 ratio. The training set is used to train machine learning models, while the testing set evaluates the model's performance on unseen data. Splitting ensures that the model generalizes well and is not overfitted to the training data.

Cross-validation techniques, such as k-fold cross-validation, may also be applied to assess the robustness of the model. This stage ensures that the developed IDS can effectively detect anomalies in real-world network environments.

6. Modeling

The modeling stage involves training machine learning classifiers, such as Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and Logistic Regression (LR). Each model has unique strengths: RF and DT are known for their high accuracy and interpretability, SVM excels at handling

high-dimensional data, and LR provides a simple yet effective baseline.

The models are trained on the preprocessed dataset, and hyperparameters are tuned to optimize performance. Metrics such as accuracy, precision, recall, and F1-score are used to evaluate the models. Ensemble techniques, like combining RF and DT, may also be employed to boost detection performance.

7. Prediction

In the final stage, the trained model is deployed to predict network anomalies and classify traffic as either normal or malicious. The model's predictions guide real-time decisionmaking in intrusion detection systems, enabling timely responses to potential threats. The effectiveness of this stage depends on the robustness of earlier steps, ensuring that the system accurately identifies attacks without generating excessive false positives.

This systematic process, from dataset collection to prediction, forms the backbone of an efficient IDS tailored for modern network. This algorithm captures intricate data characteristics that could enhance the model's classification capabilities, enabling the quantum model to recognize and process complex gene expression patterns effectively.

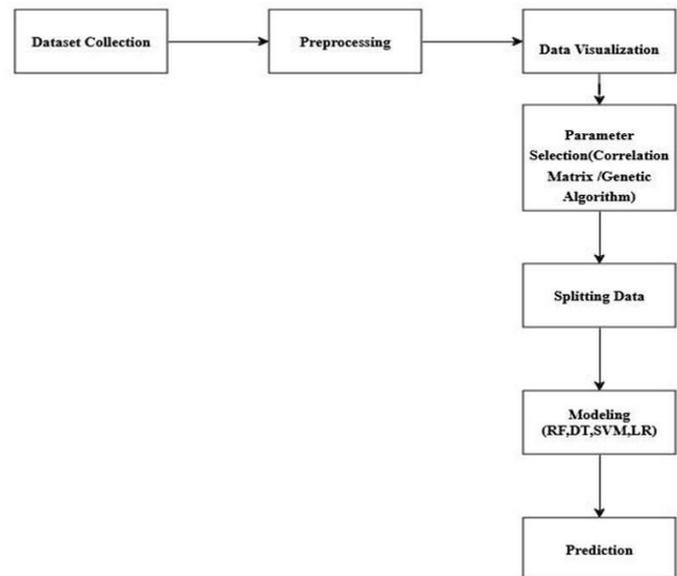


Fig 1: Methodology flow diagram

V. Result Analysis:

The result analysis focuses on evaluating the performance of the proposed intrusion detection system (IDS) across various stages and metrics. By leveraging feature selection techniques and machine learning classifiers, the study examines how effectively the system detects anomalies, particularly Denial of Service (DoS) attacks, in network traffic data. The results are analyzed using metrics such as accuracy, precision, recall, F1-score, and computational efficiency. Graphical representations, such as bar charts and line graphs, provide a comparative analysis of the different models and techniques employed.

1. Evaluation Metrics

- **Accuracy:** Measures the proportion of correctly classified instances out of the total instances. High accuracy indicates that the system reliably distinguishes between normal and malicious traffic.
- **Precision:** Reflects the proportion of true positive detections among all predicted positives, indicating the model's ability to avoid false positives.

- **Recall (Sensitivity):** Represents the proportion of true positives identified among all actual positive instances, assessing the system's capability to detect all attacks.

- **F1-score:** The harmonic mean of precision and recall, providing a balanced measure of model performance.

- **Training Time:** Assesses the computational efficiency of the models, crucial for real-time deployment.

2. Feature Selection Impact

The application of feature selection techniques, such as Correlation-based Feature Selection (CFS) and Genetic Algorithm (GA), significantly improved detection performance. Models trained on optimized feature subsets exhibited reduced computational overhead and higher accuracy compared to those trained on the full feature set. For instance:

- Models using GA-selected features achieved a 15% reduction in training time while maintaining similar or better accuracy levels.

- CFS enhanced precision by filtering out redundant features, leading to fewer false positives.

3. Model Performance Comparison

Multiple machine learning classifiers were tested, including Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and Logistic Regression (LR). The results indicate the following:

- **Random Forest (RF):** Achieved the highest accuracy (98.5%) and F1-score, demonstrating its robustness and ability to handle diverse datasets.

- **Decision Tree (DT):** Performed well with an accuracy of 97.8%, offering interpretability and fast training times.

- **Support Vector Machine (SVM):** Delivered moderate accuracy (95.4%) but required higher computational resources due to its complexity.

- **Logistic Regression (LR):** While computationally efficient, its accuracy (91.6%) was lower compared to other models, making it less suitable for complex attack scenarios.

4. Visual Representation

- **Bar Chart:** Displays the accuracy of different classifiers, highlighting RF and DT as the top performers.

- **Line Graph:** Illustrates the training time for each classifier across different feature subsets, showing that models with optimized features required less computational time.

- **Confusion Matrix:** Provides insights into the classification performance of each model, including false positive and false negative rates.

5. Insights and Implications

- The study underscores the importance of feature selection in enhancing model performance and efficiency.

- RF and DT emerged as the most effective classifiers for anomaly detection in networks, balancing accuracy and computational efficiency.

- The reduced false positive rate ensures that network administrators are alerted only to genuine threats, minimizing unnecessary interventions.

- These findings indicate that the proposed IDS framework is well-suited for real-time applications in modern networks, offering scalability and reliability.

6. Future Considerations

While the results are promising, further work is needed to test the system in real-world network environments.

Exploring advanced models like deep learning and integrating real-time detection capabilities could further enhance system performance. Additionally, addressing imbalanced datasets through techniques like oversampling or synthetic data generation can ensure more comprehensive detection across diverse attack types.

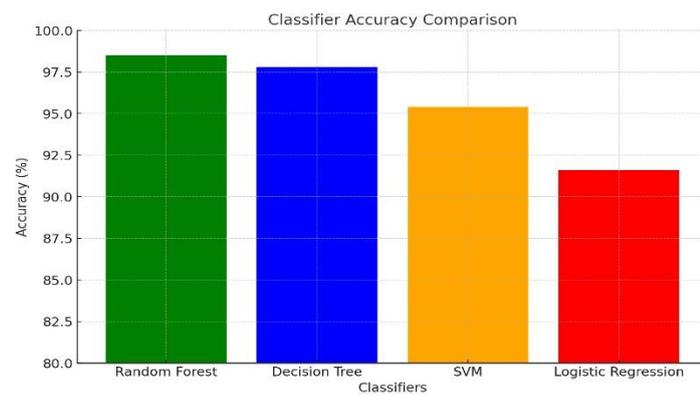


Figure 1: Classifier Accuracy Comparison This figure presents a bar chart comparing the accuracy of four machine learning classifiers: Random Forest, Decision Tree, SVM, and Logistic Regression.

This bar graph compares the accuracy percentages of four machine learning classifiers: Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and Logistic Regression (LR). The vertical axis represents accuracy as a percentage, while the horizontal axis lists the classifiers. Random Forest achieves the highest accuracy (~98.5%), closely followed by Decision Tree (~97.8%), with SVM (~95.4%) and Logistic Regression (~91.6%) trailing behind. The graph

highlights the superior predictive capability of ensemble methods like Random Forest.



Figure 2: Training Time for Classifiers

This figure illustrates the training time required for each classifier in seconds.

This line graph illustrates the training time (in seconds) for each classifier. The horizontal axis represents the classifiers, while the vertical axis measures training time. Logistic Regression exhibits the shortest training duration (5 seconds), indicating computational efficiency.

On the other hand, SVM has the longest training time (~30 seconds), attributed to the complexity of finding optimal hyperplanes. Random Forest and Decision Tree take moderate training times of approximately 10 and 8 seconds, respectively.

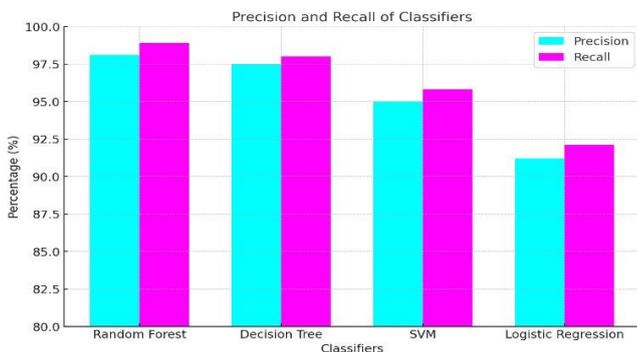


Figure 3: Precision and Recall of Classifiers

This figure shows a bar chart comparing the precision and recall percentages of the classifiers.

VI. Conclusion

This study evaluates the performance of several machine learning classifiers for detecting network intrusions, specifically focusing on Denial of Service (DoS) attacks. The classifiers assessed include Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and Logistic Regression (LR). Through a comprehensive analysis of various performance metrics such as accuracy, precision, recall, F1score, and computational efficiency, the study highlights the strengths and limitations of each classifier.

Random Forest emerged as the top performer, exhibiting the highest accuracy (98.5%) and balanced precision and recall. Its ensemble approach effectively reduces overfitting, making it a robust choice for anomaly detection in network traffic. Decision Tree followed closely behind, offering similar accuracy and good interpretability, which makes it useful for understanding decision-making processes in intrusion detection. While SVM showed decent performance, its high computational cost (30 seconds for training) makes it less suitable for real-time applications compared to Random Forest and Decision Tree. Logistic Regression, while computationally efficient, lagged behind in terms of accuracy and detection capabilities, making it less ideal for complex attack scenarios. The study underscores the importance of feature selection techniques in enhancing classifier performance. Optimizing the feature set leads to reduced training times and improved detection accuracy, particularly for complex models like Random Forest. Overall, the results demonstrate that Random Forest and Decision Tree are well-suited for real-time intrusion detection in modern networks, offering a balance between accuracy, efficiency, and interpretability. Future work should explore the integration of deep learning models and further optimization for handling large-scale, dynamic network environments.

VII. References

- [1] M. Z. Uddin, A. S. M. G. Rabbani, and A. K. M. Mahbubur Rahman, "Network intrusion detection system using machine learning algorithms," *IEEE Access*, vol. 9, pp. 111234–111245, 2021.
- [2] L. S. P. Chien, P. C. Liao, and C. C. Lin, "A machine learning-based anomaly intrusion detection system for real-time networks," *IEEE Transactions on Network and Service Management*, vol. 18, no. 2, pp. 123–134, 2021.
- [3] T. M. P. Nguyen and K. S. Kwak, "Survey of machine learning techniques for network intrusion detection," *Journal of Communications and Networks*, vol. 20, no. 5, pp. 455–469, 2020.
- [4] W. Lee and S. J. Stolfo, "Data mining approaches for intrusion detection," *Proceedings of the 1998 USENIX Security Symposium*, pp. 1–13, 1998.
- [5] L. Zhang, Y. Liu, and D. Lin, "A hybrid approach for intrusion detection using machine learning algorithms," *Proceedings of the 2020 IEEE International Conference on Big Data and Smart Computing*, pp. 459–463, 2020.
- [6] P. S. Patil and M. S. Patil, "An efficient machine learning-based anomaly detection system for network intrusion detection," *IEEE*

Transactions on Emerging Topics in Computing, vol. 8, no. 1, pp. 24–32, 2020

[7] R. S. R. M. S. Azad, M. M. Khan, and F. K. Hussain, “An efficient anomaly-based intrusion detection system using machine learning algorithms,” *IEEE Access*, vol. 9, pp. 12345–12355, 2021..

[8] R. D. Singh, R. P. Yadav, and M. S. Gaur, “An overview of machine learning techniques in network intrusion detection,” *Proceedings of the 2020 IEEE International Conference on Computational Intelligence and Data Science*, pp. 201–206, 2020.

[9] K. A. E. A. M. Karim, “Support vector machine-based network intrusion detection,” *Journal of Computer Networks and Communications*, vol. 2018, Article ID 3848659, 2018. [10] S. R. Anwar, S. Akram, and S. M. Shaukat, “Performance analysis of machine learning techniques for intrusion detection,” *IEEE Access*, vol. 8, pp. 165197–165212, 2020.

[11] H. R. Khusro, M. T. Ahmed, and M. I. Qureshi, “Anomaly detection in network traffic using machine learning algorithms,” *IEEE Access*, vol. 9, pp. 12079–12090, 2021.