

Network Intrusion Detection

1st Amer Mohamed

Autonomous Systems A

Hochschule Hamm-Lippstadt

Lippstadt, Deutschland

mohamed-ahmed-mohamed-ali.amer@stud.hshl.de

Abstract—

I. INTRODUCTION

Networks safety is becoming a crucial topic in all types of enterprises and organizations. The importance of Network security and data confidentiality are growing rapidly due to the rapid evolution in technology and AI. One of the approaches taken to ensure network safety are network intrusion detection techniques [1]. Many techniques have been introduced through the evolution of technology. The goal of intrusion detection systems is to be able to detect unusual activities, unauthorized personnel, or misuse from insiders and external penetrators in a network, preferably in real-time. [2]. One of the emerging new technologies for network security is the use of machine learning and deep learning models to identify network attacks and intruders in a network. The abundance of big data has led to the ability to train machine learning and deep learning models efficiently. More research is being done on making those models extremely accurate and efficient for large deployment. [3]

II. BACKGROUND

The main two types of Network detection system are Deployment method based and Detection method based. From the detection method perspective, it is further divided into two more categories "Signature-based intrusion detection (SIDS)" and "Anomaly detection-based intrusion detection (AIDS)".

The SIDS is based on the idea of defining a specific unique signature for network attacks. Those signatures are stored in a database. The system from there matches those signatures with the activity in the network and detects if there is a probable attack on the service. This type of approach lacks the ability to detect new types of attack as it lacks its signature and requires a huge carefully selected database which increases the computing resources needed for this algorithm [3].

The AIDS approach, also called the "behavior-based IDS," is based on the idea of defining a clear profile of normal users. Any deviation from this normal profile will be considered as an anomaly [4]. The biggest advantages of

using the AIDS approach are its ability to detect novel and new types of attacks. Though the only drawback of using this approach is the hard nature of classifying the difference between a normal and an abnormal profiles specially with the rising popularity of different IOT devices [5]. In this paper, we aim to explore Random Forest machine learning approach along with data preprocessing techniques to improve the accuracy of network intrusion detection systems. The goal is to identify the most effective model parameters for detecting malicious activities within network traffic.

III. DATASET

The dataset used in this study is the KDD Cup 1999 Intrusion Detection Dataset, which was created by simulating attacks on a U.S. Air Force LAN to capture raw TCP/IP traffic. It contains 41 features (3 qualitative and 38 quantitative features) per connection and the target variable named class is labeled as normal and anomalous behaviour [6]. The features can be grouped by type into 4 categories.

A. Features

- **Basic Features of Individual TCP Connections:** These features are extracted from the basic TCP connection. Many attacks can be determined just by analysing how the connection behaves [6], [7].
- **Content Features within a Connection:** These are features that inspect the payload or command content of a connection to detect any suspicious behaviour. They go beyond the basic properties of a connection and look deeper into what is actually being transmitted during the session. [6], [7]
- **Time based Traffic features:** These features consider connections to the same host in the past two seconds. [6], [7]
- **Host-based Traffic Features:** It is similar to Time based traffic features but these use a larger time window to detect patterns in connections. [6], [7]

B. Types of attacks

- **Denial of Service Attack (DOS):** It is a type of attack where the attacker overloads computing and memory resource. That makes the service too busy to fully handle legitimate requests and denies legitimate users access to the machine [7], [8]

Identify applicable funding agency here. If none, delete this.

- **User to Root Attack (U2R):** This type of attack occurs when the attacker starts out with access to a normal legitimate account on the system and then becoming able to exploit vulnerability to gain root access to the system [7], [8]
- **Remote to Local Attack (R2L):** It occurs when an attacker can send packets to a machine over a network where the attacker does not have access as a user to that machine [7], [8]
- **Probing Attack:** It is an attempt to gather information about a network of computers for the purpose of breaching through their security and gaining root access [7], [8]

C. Potential Issues in the Dataset

Previous research shows that this specific dataset has some issues. They experimented with various machine learning models all of which showed a very high accuracy of approximately 98% on the training data set while having inconsistent accuracy fluctuations through epochs for the validation data set. In other terms, it will not translate well into real life deployment. The first important deficiency in the KDD data set is the huge number of redundant records. Analysing KDD train and test sets, it was found that about 78% and 75% of the records are duplicated in the train and test set respectively. This will cause the model to be biased towards the more frequent records and prevent it from learning the novel records. [8], [9]. The previous research made on the original KDD dataset introduced a newer version dataset called NSL-KDD. The researchers who made this addressed all the current issues observed in KDD data set and made a new dataset that performed better with machine learning models in binary classification and multi-class classification models [9].

D. Data set Preprocessing

By analysing the features provided in the KDD dataset, two features had an always zero value which were "num_outbound_cmds, is_host_login" features. Those two features were dropped in the data preprocessing because they did not provide any importance in the model training. Moreover, three categorical columns were included in the dataset. After experimenting with encoding algorithms. "HotEncoding" the features proved the best metrics in model evaluation.

IV. MACHINE LEARNING MODEL

A. Overview

For the network intrusion detection dataset, A random forest classifier algorithm was used. A random forest classifier is based on classification tree algorithm. A classification tree splits data recursively based on feature thresholds to create a tree structure that predicts class labels. Its aim is to divide the data set into pure regions where almost all samples belong to a single class [10]. The algorithm works by partitioning the feature space into rectangles and then fitting a simple constant in each [11].

B. The mathematical model

The classification tree which is the parent method of Random forest classification tree works by having a root node that represents the first input and the entire data to be used, then it keeps on branching. each internal node represents decisions made depending on the input at an instance. Each leaf of the tree represents class labels or the final prediction. The splitting of nodes are usually based on a mathematical relation like "Recursive binary splitting" shown in equation (1) [9]–[11].

At each node of a classification tree, the dataset is split into two regions using a feature x_j and threshold s , as follows:

$$\begin{aligned} R_1(j, s) &= \{\mathbf{x} \mid x_j < s\}, \\ R_2(j, s) &= \{\mathbf{x} \mid x_j \geq s\} \end{aligned} \quad (1)$$

where:

- $\mathbf{x} \in \mathbb{R}^p$ is the feature vector
- x_j is its j -th component, and s is the splitting threshold.
- x_j is the j -th feature
- s is the threshold that defines the split.

In other words the algorithm tries to find the best pair of (j,s) that result in the purest split [10], [11].

The purity of the split is calculated by "Gini Index Equation" shown in equation (2) [11].

$$G(t) = 1 - \sum_{k=1}^K p_k^2 \quad (2)$$

where:

- K is the number of classes,
- p_k is the proportion of samples in node t that belong to class k .

The lower the value of $G(t)$ indicate purer nodes thus, better splits.

More algorithms for getting the purest split are sometimes used like "Cross-entropy" and "Misclassification error" [11].

Random forest is built upon the classification trees. It is a group of classification decision trees that merges to get more accurate values. In simpler terms, it is the same as using multiple decision trees machine learning models and getting the best output out of all of them [12]. In Fig. 1 is a small visualization of the random forest algorithm with two decision trees [12].

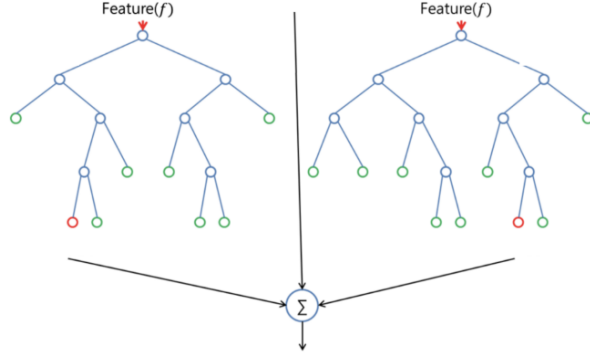


Fig. 1. Random Forest visualization with 2 decision trees

One of the downsides of the random forest algorithm is that it has more tendency to overfit the data when the tree is deep. That is why it is important to utilize hyper parameters to ensure not having a too deep tree structure [12].

V. IMPLEMENTATION OF RANDOM FOREST CLASSIFIER

The supplied dataset was used to train and test the model due to the lack of labelling on the validation dataset. 80% of the data was used for the training and the rest for validation. the metrics captured showed a very high suspicious accuracy. Although both validation and training accuracy difference was too low, which shows that the model was working great. a very high accuracy of approximately 0.99 seems too high to generalize on real life. This shows that the dataset had issues as discussed in section 3. The data are too simple and too synthetic for the machine learning model to generalize and there is data leakage in the features. the training and validation accuracy are shown in Fig. 2.

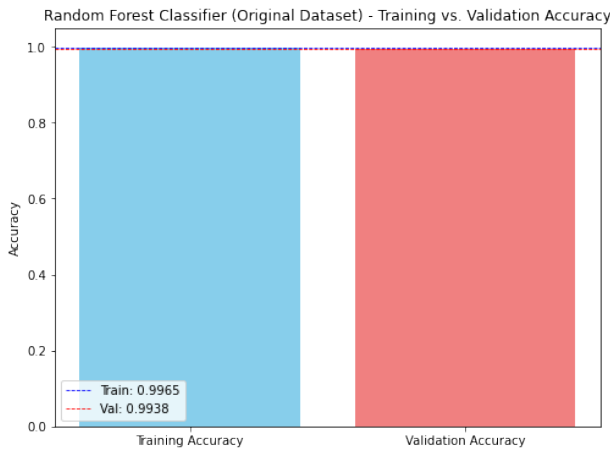


Fig. 2. Training and validation accuracy on KDD1999 dataset using random forest classifier machine learning algorithm.

A. Modified Dataset

Redoing the same process with the modified dataset showed different metrics. Due to the fact that the modified dataset had a valid labelled validation dataset. When splitting the training dataset and using 80% for the training and the rest for validation. The model did well in terms of accuracy, F1 score, recall and precision with the splitted validation data. but when using the entire training dataset for training the model and taking the validation data provided to test the model. the accuracy decreased form 99% to 76%. Though an extensive search grid was used to tune the hyper parameters, just a slight difference in the metrics where found. The accuracy achieved 76% but it can be considered that the model can be actually implemented in a real-life environment due to the fact that the model achieved 97% recall. which means that 97% of attacks can be easily recognized but on the other hand the model showed bad recall for normal connections. In other terms, The model would flag 57% of normal connections as attacks which might cause some issues for real-world deployment.

VI. DEEP LEARNING APPROACH

A. Original Dataset

A deep learning algorithm was used that included 4 layers including the output layer. batch normalization for further data scaling before each layer and an aggressive drop out between layers the model performed nearly normally on the original dataset shown in Fig. 3.

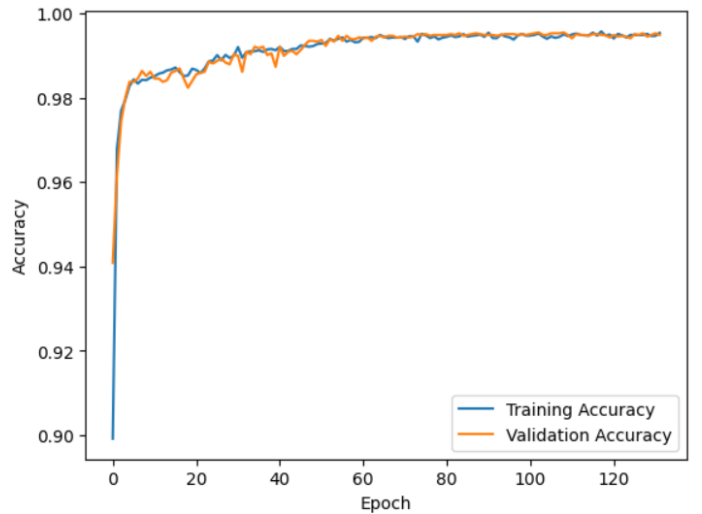


Fig. 3. Training and validation accuracy on KDD1999 original dataset after optimizing deep learning algorithm.

This particularly shows that the original dataset can still be used but with more complex deep learning models that utilizes strong penalties on the loss function, dropping out a random percentage of the neurons in training and optimizing the learning rate can help the model being more stable. Since accuracy is really high, this probably indicates data leakage or that the data is too simple for a Deep learning algorithm. The data was gathered in 1999 by simulating a network intrusion

attack. Protocols then were too simple thus, the intrusion detection can be achieved by way simpler methods.

B. Modified Dataset

After trying the same model and parameters used on the original dataset. Splitting the training data with 80 to 20 ratio. The model showed good overall performance but with many fluctuations for the validation accuracy shown in Fig. 4.

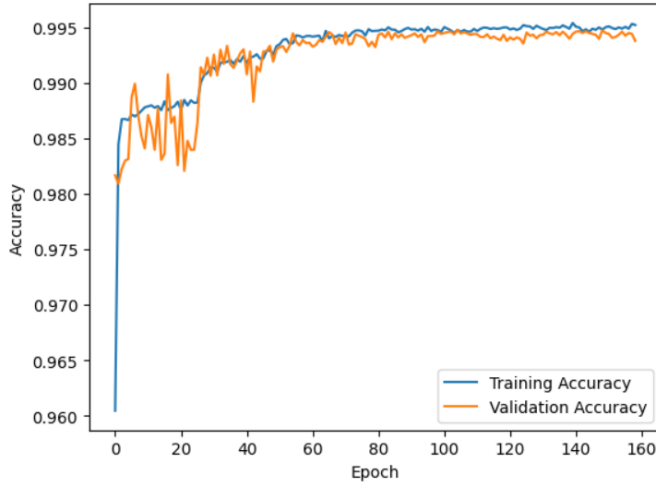


Fig. 4. Training and validation accuracy on NSL-KDD1999 modified dataset with deep learning using the same parameters of the model applied on the original dataset.

This showed that the dropout and regularization methods used for the original dataset made the model unstable because the new dataset was already modified to have stable data with no bias and the old dataset issues. Hence, less normalization and strict behaviour adjustment to the model are needed. More optimization of the parameters are needed to reach good accuracy and decrease the fluctuation with the modified dataset.

When experimented to use the full training data for training and using the validation data for testing the model still performed poorly like in section 5.A. This further proves that the data provided even after modification tends to be not perfect for real life deployment. Moreover, the new data set supplied did not mention explicitly the difference between train and test data and the way of getting this data which might point that the testing data was collected differently in a different environment that the old data from the model could not generalize it.

VII. CONCLUSION

Network Intrusion have been always a big worry due to the fast paced technologies evolving over the years and the need to feel secure while using technology is crucial for humans. That is why developing a model for network intrusion detection that would perform perfectly is really needed. In the approaches

used in this paper, It was deduced that the need for more up to date problem-free data is still needed for a machine learning or deep learning model to perform better in real life scenarios of network attacks. More research and algorithms are still needed to provide the system with the needed safety and reliability.

REFERENCES

- [1] S. Kumar, Survey of current network intrusion detection techniques, Washington Univ. in St. Louis, pp. 1–18, 2007.
- [2] B. Mukherjee, L. T. Heberlein and K. N. Levitt, "Network intrusion detection," in *IEEE Network*, vol. 8, no. 3, pp. 26–41, May–June 1994, doi: 10.1109/65.283931. keywords: Intrusion detection; Computer networks; Protection; Computer security; Data security; Computer science; Computer crime; Information security; Real time systems; Prototypes,
- [3] Z. Ahmad, A. Shahid Khan, C. Wai Shiang, J. Abdullah, and F. Ahmad, "Network intrusion detection system: A systematic study of machine learning and deep learning approaches," *Trans. Emerging Telecommun. Technol.*, vol. 32, no. 1, e4150, 2021. [Online]. Available: <https://doi.org/10.1002/ett.4150>
- [4] W. Ma, "Analysis of anomaly detection method for Internet of Things based on deep learning," *Trans. Emerg. Telecommun. Technol.*, vol. 31, no. 6, e3893, 2020. [Online]. Available: <https://doi.org/10.1002/ett.3893>
- [5] Y. Mehmood, F. Ahmad, I. Yaqoob, A. Adnane, M. Imran, and S. Guizani, "Internet-of-Things-based smart cities: Recent advances and challenges," *IEEE Commun. Mag.*, vol. 55, no. 9, pp. 16–24, Sep. 2017. [Online]. Available: <https://doi.org/10.1109/MCOM.2017.1600514>
- [6] KDD Cup 1999, "KDD Cup 1999 Data," 1999. [Online]. Available: [Accessed: May 13, 2025].
- [7] KDD Cup, "KDD Cup 1999 Data," [Online]. Available: <http://kdd.ics.uci.edu/databases/kddcup99/task.html>, [Accessed: May 13, 2025].
- [8] M. Tavallae, E. Bagheri, W. Lu and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," 2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications, Ottawa, ON, Canada, 2009, pp. 1–6, doi: 10.1109/CISDA.2009.5356528. keywords: Testing; Intrusion detection; Data security; Statistical analysis; Computer security; Computer aided manufacturing; Learning systems; Computational intelligence; Computer networks; Application software,
- [9] A. R. Tapsoba and T. Frédéric OUEDRAOGO, "Evaluation of supervised learning algorithms in binary and multi-class network anomalies detection," 2021 IEEE AFRICON, Arusha, Tanzania, United Republic of, 2021, pp. 1–6, doi: 10.1109/AFRICON51333.2021.9570886. keywords: Training; Supervised learning; Support vector machine classification; Predictive models; Prediction algorithms; Feature extraction; Classification algorithms; Intrusion Detection System (IDS); Supervised Learning Algorithms (SLA); Recursive Feature Elimination (RFE); AUC - ROC Curve; NSL-KDD,
- [10] F. Pedregosa et al., "Tree-based models," *Scikit-learn: Machine Learning in Python*, [Online]. Available: <https://scikit-learn.org/stable/modules/tree.html> [Accessed: May 20, 2025].
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York, NY, USA: Springer, 2009.
- [12] M. S. Brown, "Random Forest Algorithm: A Complete Guide," *Built In*, 2023. [Online]. Available: <https://builtin.com/data-science/random-forest-algorithm>