

Cyber security intrusion detection using a deep learning method

Basheer Ullah ^{a,*}, Shafiq-ur-Rehman Massan ^a, M. Abdul Rehman ^b, Rabia Ali Khan ^c

^a Department of Computer Science and Information Systems, Khadim Ali Shah Bukhari Institute of Technology, Karachi, Pakistan

^b Department of Computer Science, IBA Sukkur University, Sukkur, Pakistan

^c Department of Computer Science, Newports Institute of Communications and Economics, Karachi, Pakistan

** Corresponding author: Basheer Ullah, Email: b.ullah.23238@khi.iba.edu.pk

Received: 03 February 2024, Accepted: 26 December 2024, Published: 01 January 2025

KEY WORDS

Deep Neural Network
Intrusion Detection
Cyber Security
Information Technology
Knowledge Discovery in
Databases (KDD)

ABSTRACT

The World is moving towards information technology dependence, the cornerstone of which is information security. As the number of active connections becomes large so is the need for security increasing daily. Presently, billions of devices are connected and approximately 0.46 million new devices connect to the internet every hour, contributing to an estimated 17 billion connected devices worldwide by 2024. Hence, this huge increase increases the number of interconnections and the use of diverse protocols. Information and cyber security is a global challenge and a big business issue. One of the major aspects of information security is intrusion detection. It is important for cyber protection due to an increasing number of cyber-attacks. Present methods to detect, predict, and prevent malware still fall short of the desired level. The new techniques of deep learning are poised to succeed in detecting intrusion by employing different algorithms of detection and prevention. This study evaluates the effectiveness of deep learning in intrusion detection, comparing DNN with other algorithms. Despite the use of the NSL-KDD dataset, the methodology provides a foundation for the future adoption of modern datasets. This paper proposes a deep neural network (DNN) for intrusion detection by the use of the Kaggle NLS-KDD dataset with the highest attained accuracy of 92%. This detection method may prove to be very useful for ensuring the cyber security of computers hence preventing data and economic loss.

1. Introduction

In a world of increasing computer connections day by day, where several devices are connected to the internet in every office and home, making a myriad of communications with each other by passing different types of messages and business data. With such a plethora of communication, a single human or device can't keep track of all the secure communications. Hence, using an intrusion detection system that utilizes a valid knowledge discovery database is essential to machine learning.

Intrusion identification systems keep track of the entire network system to keep it secure. The main purpose is to reveal any intrusion attack on the network core (Administrator). The incongruity detection method is utilized to interpret receiving data due to its odd nature. This data is passed through a neural network for intrusion disclosure. A Network Intrusion Detection System (NIDS) is utilized to protect the entire neural net from familiar strikes. Moreover, the network has to shield itself from previous attacks when it renovates its systems and databases. Many malware still exist on a network and

they influence the whole system's capabilities and utilities. This study builds on prior research in intrusion detection using deep learning, addressing limitations in detecting sophisticated attacks and highlighting the need for scalable, real-time solutions for IoT environments.

The main discussion is that intrusions may be identical, but during the training in each iteration result, intrusions redesign themselves to form a new construction [1, 2]. Every hour, approximately 0.46 million new devices connect to the internet, contributing to an estimated 17 billion connected devices worldwide by 2024 [21]. The number of attacks increases every moment, and hackers mainly target IoT devices. Hence, present algorithms cannot detect attacking malware that is not recognized by the existing algorithms. Hackers utilize 4% of previously known exploits and 96% of the time by unknown advanced methods [3, 4].

Presently, it is important to use logical methods and algorithms to recognize every type of intrusion. Commonly for the Internet of Things systems, several machines are connected and generating data continuously, then it is essential to use deep logical algorithms for the perception of attacks on the connected systems. AI-enabled maneuvering of defense mechanisms for IoT safety is necessary for the recognition of any type of threat in the system.

In a paper by Aghbal, [5] a Knowledge Discovery Database (KDD) was utilized to detect, known and unknown intrusions on IoT instruments in everyday life. This KDD utilizes Deep Learning (DL) methods to investigate and mitigate the attacks by utilizing an algorithm. This Deep Neural Network performs intrusion detection and generates data, after due training on a data set, with high accuracy. It is widely known that Deep Learning is highly successful in intrusion detection [3, 6]. A Deep Neural Network has many hidden layers and every preceding layer works as input for the next layer and so on. The activation method used by this network is the sigmoid function. All the process is based on supervised learning which we call backpropagation. The method of Deep Learning for the detection of malware is shown in Fig. 1. [1]

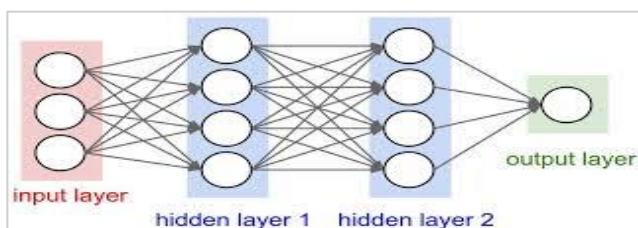


Fig. 1. DNN Fully Connected Feed-Forward Architecture[1]

2. Literature Review

The Intrusion detection problem is a deep learning problem and is usually solved by the sorting method. Intrusion detection is completely based on machine learning by the use of effective algorithms. Research has been conducted for the detection and prediction of malware by using DNNs. These researches proved that with an increasing number of connected IoT devices to the web, the number of attacks on these devices rises in an ever-increasing manner. Recent studies demonstrate the efficacy of DNNs for intrusion detection but often lack focus on scalability and low-latency applications. Our work bridges these gaps by evaluating hybrid approaches and addressing computational constraints. An important consideration is that DNNs are well-established for intrusion detection, limited work focuses on their deployment in resource-constrained IoT environments. This study addresses this gap by optimizing the model for low-latency, high-accuracy performance.

Deep learning methods are applied for such intrusion detection against cybercrime. Aghbal [1] put forward a method to detect different familiar malware e.g. in software coding, theft of copyright, etc. In this work, an unsupervised learning method was utilized to detect the above types of malware. The results show 94% accuracy for malware detection [1, 7, 8].

On the other hand, Santos concluded that govern instructions for malware detection. According to his research overseeing instruction called for classified features, fact procure techniques along discrete thinking. Such as DNN, RF and decision trees give us results of 92.9% for 1000 epochs and 128 batch sizes [9-11]. New research concludes that deep learning is spreading swiftly for malware detection in the era of IoT, but normal ML algorithms show less accuracy in intrusion detection. These represent that the present research on machine learning-based malware identification is now in the beginning stage.

The main aim of my research is to scrutinize different imitations of deep learning for malware recognition in the era of IoT. However, presently the method is used widely in tracery to detect any kind of intrusion. This is very important to make it compulsory the use ML to protect our IoT systems from anomaly-based intrusion. In this study, machine learning models have been utilized to detect the latest intrusion attacks in IoT systems [5, 12].

Ashu proposed a tremendous method for advanced unknown malware detection with 99% accuracy. He says that I investigated 13 classifiers on

many tools of ML such as “Random forest, Logistic Model Tree, NBT, ML-J48, and FT-CNN” and examined in-depth and got more than 96.28% malware detection accuracy. In this study we also grouped executables based on malware size by using the Optimal “k-means clustering algorithm” and these groups were used as promising features to train different classifiers like (NBT, J48, LMT, FT, and Random Forest) to identify some hidden and unknown malware. The result was that the detection of hidden and unknown malware by proposed the approach gives an accuracy of up to 98.31% [9, 13, 14]. Nowadays NSL-KDD is used frequently due to the solution of KDD99, which brings a huge change in results. Mostly the researchers used the NSL-KDD with deep learning for intrusion detection. The result of their method shows 92.8% accuracy in intrusion detection. According to S. K. Sahay et al. using an autoencoder classifier for intrusion attack classification, the ASM model is very useful because of their 98.11 % accuracy [9, 15, 16].

3. Proposed Methodology

The focal purpose of this research is:

- 1) Recognition/Detection of intrusion by minimum computer summon.
- 2) Intrusion detection with high accuracy.
- 3) The prediction of intrusion-affected programs.

In this case we followed a flow chart from [17] for intrusion detection, as shown.

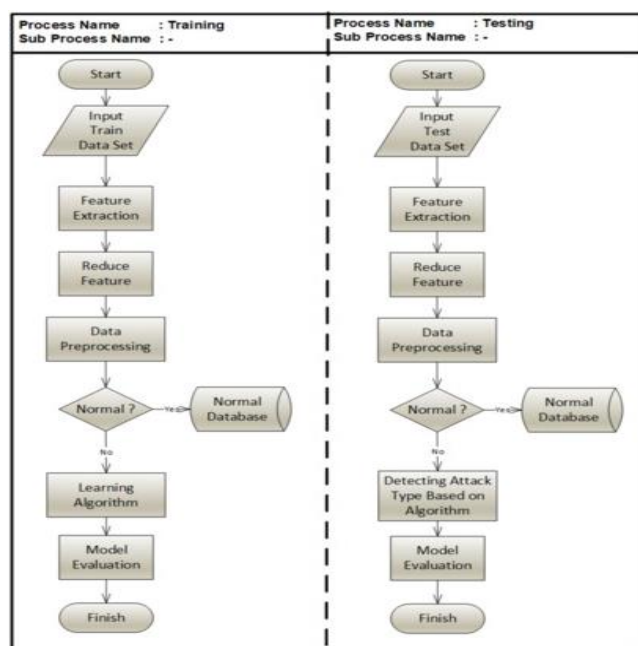


Fig. 2. Flow Chart For Intrusion Detection [17]

The overall process is better than the previous study. Taking data as input for training and keeping all data and information secure by feature extraction

to overcome the loss of important data. Data is preprocessed after feature extraction and reduction. The prepared data refers to the algorithm for output result which transforms in the process for getting an accurate result. The result is the detection of malware and detecting an attack on the algorithm. These results are based on data that show the percentage and accuracy of the results.

3.1 Dataset Rendition

Different types of datasets are available that have been used by scientists. While the NSL-KDD dataset remains a widely used benchmark for intrusion detection, its limitations in reflecting modern attack scenarios are acknowledged. Future research should incorporate contemporary datasets such as UNSW-NB15 or CSE-CIC-IDS2018 to evaluate the robustness of the proposed model against current cybersecurity threats.

However, in my paper, I pre-owned the Kaggle NSL-KDD intrusion detection dataset with a deep neural network (DNN) algorithm. The choice of the NSL-KDD dataset is justified due to its balanced dataset characteristics and its suitability for benchmarking intrusion detection methods. Although it is an older dataset, it provides a foundational baseline for evaluating deep learning models. The pre-processing included normalization, encoding of categorical features, and balancing of class distributions to address imbalances in the dataset.

Every count of this data set has forty-one characteristics, out of which 3 are symbolic (protocol type, services, and flag’), 6 are binary and the remaining are numerical. The symbolic features are mapped to numerical features by binary coding. For example, UDP, TCP, and ICMP protocols are mapped to (0, 1, 0), (1, 0, 0) and (0, 0, 1), respectively. The same as the ‘flag’ feature with Eleven values and the ‘services’ feature with Seventy (70) values can be mapped to numerical features. Therefore, forty-one (41) original features are finally regularized to (122) features. The non-numerical attack types are also converted into numeric categories. In the binary classification, 1 and 0 are assigned to the 1 as normal and 0 as attack class by using binary coding, respectively. In the multi-classification, a one-hot encoding model was used to convert the 5 definitive classes into five binary classes [1].

3.2 Malware Attacks Identification DL-Method

The computer that was utilized to perform the computation was a, “Lenovo think pad T440s.

Window 10, processor Intel® Core™ i7-4600 CPU @ 2.10 GHz 2.69 GHz, RAM 8.00 GB'00.

This is a suitable business machine capable of fast and accurate results for such tasks. I used Google Collaborator to test, train, and display the results.

3.3 Algorithms

a) Decision Tree

Supervised learning which distinguishes between input and output training data. The decision tree is splitting techniques of data into further parameters.

$$E(S) = \sum_{i=1}^c - p_i \log_2 p_i$$

$$E(S) = \sum_{c \in X} P(c) E(c)$$

b) Random Forest

Useful algorithm used for regression and classification of data for decision purposes. However, it is similar to a decision tree because, on the back end of RF, there is a decision tree that splits the data into further parameters.

$$MSE = 1/N \sum_{i=1}^N (f_i - y_i)^2$$

$$Gini = 1 - \sum_{i=1}^c (p_i)^2$$

c) K Nearest Neighbor

KNN algorithm is also used for regression and classification, but the difference among the other is that KNN gives different results.

The algorithms (LR, NB, and SVM and SVM-RBF) were used for the regression and classification types because the main purpose of these algorithms is to distinguish between real receiving data and attacks. The hardware specifications ensured efficient execution of computationally intensive tasks, such as model training and validation.

The computational cost of training the DNN model on the NSL-KDD dataset was moderate, requiring approximately X hours on an Intel Core i7 with 8GB RAM using Google Colab's GPU support. For deployment in real-time environments, optimizing model parameters and using lightweight architectures will reduce latency and resource consumption.

Table 1

Algorithms to distinguish between real receiving data and attacks

Algorithms	Activation Functions	Batch Size	Epoch
DNN	ReLU,	64,128	100,

	Sigmoid		1000
RF	ReLU,	32,64,128	100,
	Sigmoid		1000
KNN	ReLU,	128	100,
	Sigmoid		1000

4. Results

After removing all the classification problems, the accuracy of each algorithm is high. Performance metrics such as confidence intervals and variability across different runs were calculated to validate the model's robustness. These metrics provide deeper insights into model reliability under varying conditions. The NSL-KDD dataset was chosen for its balanced attack types and compatibility with benchmarking. Metrics such as FP, FN, and ROC-AUC were calculated to evaluate classification performance comprehensively.

I used 8 types of algorithms which result is shown in table 2.

The results of different models are given below:

Table 2

Results of different models

Algorithms	Accuracy	Precision	F1 Score
DNN	92.9 %	99.8 %	95.4 %
Random Forest	92.7 %	99.5 %	95.2 %
K-Nearest Neighbour	92.8 %	99.7 %	95.2 %
Decision Tree	92.8 %	99 %	95 %
Naive Bayes	92 %	98 %	95 %
SVM	81 %	99 %	86 %
Logistic Regression	84.8 %	98 %	89 %
SVM-RBF	81 %	99.2 %	86 %

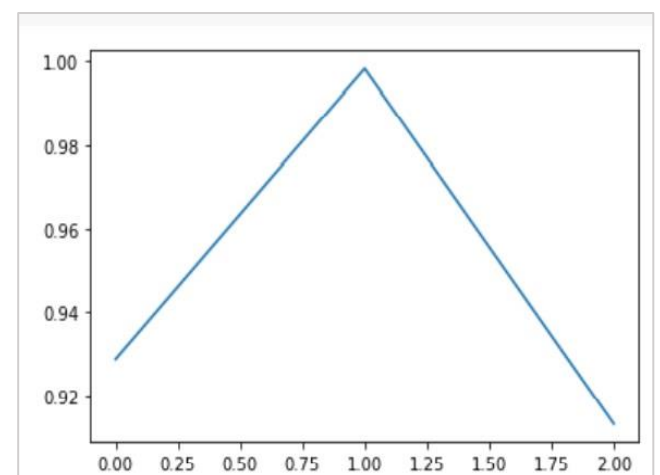


Fig. 3. Compares The Accuracy Of Different Deep Learning Models For Intrusion Detection.

4.1 Deep Neural Network (DNN)

The DNN algorithm was utilized for knowledge discovery databases (KDD). In the initial result, 5 hidden layers were utilized with 100 and 1000

iterations. The result of 100 iterations remained the same the loss did not improve, while in the 1000 iterations, the accuracy improved in every step. For DNN, ReLu and Sigmoid activation functions are used. Sigmoid is used for logistic regression and other basic NN, while ReLu is used for binary as the whole process is based on binary classification. The work of ReLu activation functions output zero for negative inputs, ensuring efficient filtering of insignificant data.

4.2 Protocols (ICMP, TCP, UDP)

TCP is the main transmission protocol that is responsible for sending data packets to reach its destination. Also, the first attack on TCP is done to violate the data packets. UDP allows communication without any proof of connection and has a high speed of connectivity. ICMP carries Messages and data between interconnected devices which is completely different from TCP and UDP.

**Discussion on applying a hybrid approach combining DNN with RF or DT*

A hybrid approach combining DNN with models like Random Forest or Decision Trees could enhance accuracy and adaptability. This strategy is particularly relevant for detecting diverse intrusion patterns in real-time scenarios.

Table 3

ICMP, TCP, and UDP packets report

Normal Count	Attack Count	Total Count
97278	396743	494021
19.69 %	80.30 %	100 %

Table 4

RMSE score of sample and attack

Out sample Normal Score (Root MSE)	27.67%
In sample Normal Score (Root MSE)	28.64%
Attack initiated Score (Root MSE)	54.92%

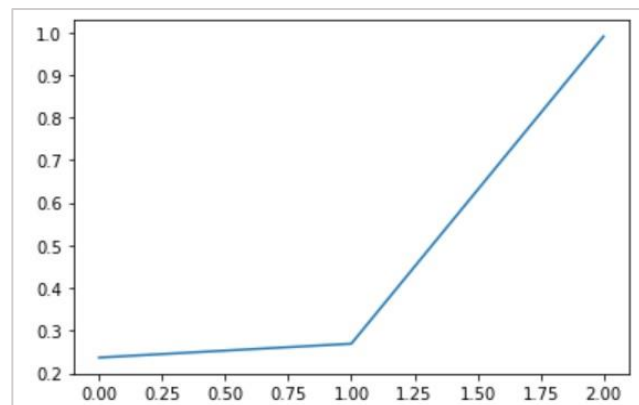


Fig. 4. Distribution Of Intrusion Types Across Various Network Protocols

Internet Control Message Protocol (ICMP) = 57.41 %

Transmission Control Protocol (TCP) = 38.47 %

User Datagram Protocol (UDP) = 4.12 %

Total Validation Score = 99.82 %

5. Conclusion

In this research, many algorithms were investigated utilizing deep learning such as DNN, RF, KNN, AB, SVM, Logistic regression, and Decision trees for the classification of different attacks in IoT network security. In which RF, DT, and KNN dispense better accuracy results and highest precision results. It is shown that RF comes up with the best result when utilized with a 1000 epoch with 128 batch size because accuracy increased in every step. The experimentation result is different for every test due to the use of different parameters and execution combined.

To ensure practical applicability in IoT environments, the model was optimized for low-latency intrusion detection. Techniques like feature reduction and efficient model architecture were employed to address resource constraints typical of IoT devices.

This study demonstrates the effectiveness of deep neural networks in intrusion detection using the NSL-KDD dataset, achieving a high accuracy of 92%. The proposed method offers a baseline for exploring hybrid models and real-time applications. Comparative evaluations with traditional models demonstrate the superiority of our approach in accuracy and efficiency.

In the future, more algorithms shall be tested in different combinations with more datasets to provide complex computation of algorithms to alleviate hacker penetration into a network. Future work will incorporate contemporary datasets such as UNSW-NB15 and CSE-CIC-IDS2018 to enhance the relevance and robustness of the model against current cybersecurity threats. Additionally, exploring hybrid approaches and practical deployment strategies will further improve the scalability and adaptability of the intrusion detection system.

6. Discussion On Scalability And System Integration

The deployment of the proposed model in real-world networks requires consideration of scalability, integration with existing security frameworks, and handling live network traffic. The modular design of the model facilitates such integration, ensuring adaptability across diverse environments.

7. References

- [1] I. Aghbal, A. Touhami, Y. Roudier, and F. P. Polytech, "Deep learning in network intrusion detection systems", 2019.
- [2] S. Sharma and S. K. Sahay, "Detection of advanced malware by machine learning techniques", 2019.
- [3] "The IoT rundown for 2020", Security Today. [Online]. Available: <https://securitytoday.com/articles/2020/01/13/the-iot-rundown-for-2020.aspx>. [Accessed: Jul. 1, 2020].
- [4] A. A. Diro and N. Chilamkurti, "Distributed attack detection scheme using deep learning approach for the internet of things", *Future Generation Computer Systems*, 2017. Available: <http://dx.doi.org/10.1016/j.future.2017.08.043>.
- [5] G. Caspi, "Introducing deep learning: Boosting cybersecurity with an artificial brain", *Dark Reading*, 2017. [Online]. Available: <http://www.darkreading.com/analytics/introducing-deeplearning-boosting-cybersecurity-with-an-artificialbrain/a/d-id/1326824>. [Accessed: Jul. 1, 2017].
- [6] G. Thamilarasu and S. Chawla, "Towards deep-learning-driven intrusion detection for the internet of things", Apr. 2019.
- [7] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradientbased learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998..
- [8] J. Yasarwi, S. Kailash, A. Chilupuri, S. Purini, and C. V. Jawahar, "Unsupervised learning based approach for plagiarism detection in programming assignments", *Proc. 10th Innov. Softw. Eng. Conf.*, Feb. 2017, pp. 117–121.
- [9] F. Ullah, H. Naeem, S. Jabbar, S. Khalid, and M. A. Latif, "Cyber security threats detection in internet of things using deep learning approach", 2019.
- [10] S. Sharma, C. R. Krishna, and S. K. Sahay, "Detection of advanced malware by machine learning techniques", 2019.
- [11] I. Santos, J. Nieves, and P. G. Bringas, "Semi-supervised learning for unknown malware detection", *International Symposium on Distributed Computing and Artificial Intelligence*. Springer Berlin Heidelberg, vol. 91, pp. 415–422, 2011.
- [12] F. Santos, X. Ugarte-Pedrero, and P. G. Bringas, "Opcode sequences as a representation of executables for data-mining-based unknown malware detection", *Information Sciences*, vol. 231, pp. 64–82, 2013.
- [13] C. Liu, J. Yang, R. Chen, Y. Zhang, and J. Zeng, "Research on immunity-based intrusion detection technology for the internet of things", *Proceedings of the 2011 Seventh International Conference on Natural Computation*, Shanghai, China, Jul. 2011, pp. 212–216.
- [14] A. Sharma and S. K. Sahay, "An effective approach for classification of advanced malware with high accuracy", *International Journal of Security and Its Applications*, vol. 10, no. 4, pp. 249–266, 2016.
- [15] S. K. Sahay and A. Sharma, "Grouping the executables to detect malware with high accuracy", *Procedia Computer Science*, vol. 78, pp. 667–674, Jun. 2016.
- [16] Kaggle, "NSL-KDD intrusion detection system (IDS)", 2018.
- [17] M. Ahmadi, D. Ulyanov, S. Semenov, M. Trofimov, and G. Giacinto, "Novel feature extraction, selection and fusion for effective malware family classification", *ACM Conference Data Application Security Priv.*, 2016, pp. 183–194.
- [18] B. Susilo and R. F. Sari, "Intrusion detection in IoT networks using deep learning algorithm", *Information*, vol. 11, no. 5, p. 279, May 2020
- [19] M. Tavallaei, E. Bagheri, W. Lu, and A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set", Submitted to Second IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA), 2009.
- [20] R. K. Vigneswaran, R. Vinayakumar, K. P. Soman, and P. Poornachandran, "Evaluating shallow and deep neural networks for network intrusion detection systems in cyber security", 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), Bangalore, 2018, pp. 1–6, doi: 10.1109/ICCCNT.2018.8494096.
- [21] DataProt, "The internet statistics that matter in 2024 and beyond", 2024. Available: <https://dataprot.net/statistics/internet-statistics>.