# An Ensemble-Based Classifier Network Intrusion Detection

**Vasumathi AK [a], Banupriya V [b], Viswanath Kani T [c]**

[a] PG Scholar Department of Computer Science and Engineering, Vivekanandha College of Engineering for Women(Autonomous). Elayampalayam, Tiruchengode, Namakkal-637205, Tamilnadu, India.
[b] Assistant Professor  Department of Computer Science and Engineering, Vivekanandha College of Engineering for Women(Autonomous). Elayampalayam, Tiruchengode, Namakkal-637205, Tamilnadu, India.
[c] Assistant Professor Department of Computer Science and Engineering, Vivekanandha College of Engineering for Women(Autonomous). Elayampalayam, Tiruchengode, Namakkal-637205, Tamilnadu, India.

_____

**Abstract:** An Intrusion Detection System monitors the flow of data on a network through the computer for the search of malicious activities which are majorly known as threats and viruses. There are two types of intrusion detection, one, Signature-based detection where the intrusion detection collects the information, analyses it, and then compares them to the attack signatures stored in the database. The second one is an Anomaly-based intrusion detection system that learns normal and anomalous behaviour by analysis in various benchmark datasets. Common challenges for Intrusion Detection Systems are large amounts of data to process, low detection rates, and high rates of false alarms. Considering anomaly pattern as detecting a point in time where the behaviour of the system is unusual and significantly different from past behaviour. In such context anomaly detection mean detecting the behaviours that deviate from normal behaviours. An ensemble based classifier method is considered using Naïve Bayes and Multivariate Linear Regression algorithms. To get a better accuracy rate of the intrusion for real time data packets are received in the system. On the experimental results achieved, we are proposing the Naïve Bayes and Multivariate techniques as an efficient method for network intrusion detection.

**Keywords:** Network Intrusion Detection, NSL KDD Dataset, Classifiers, Ensemble, Naïve Bayes, Multivariate Linear Regression, Feature Selection.

1. Introduction

As with the increasing applications of the computer and network nowadays, the security problems in information system are becoming complicated. Though we have a lot of methods in machine learning, there are issues in handling huge and imbalanced dataset machine learning based intrusion detection systems for real-time data packets, it faces a lot of challenges to process the entire data. Hence, it is necessary to identify types of intrusions through network traffic behavior. The Intrusion detection system is divided into two categories: Host Based Intrusion Detection System and Network Based Intrusion Detection System.

Host Based Intrusion Detection System can detect the internal changes, where, an attack can occur inside the computer due to certain applications and would spread inside the system. But, Network Intrusion Detection System detects malicious activities in the packets of the network as they enter the network or by unusual behavior on the network by their attacks. The advantage of a network intrusion detection system is monitoring a large network.

The intrusion detection system has many methods and classifiers based on the datasets and attacks. The machine learning techniques have been adapted for the detection rate. It reduces false positives and increases the true positives trying to give an accuracy rate for the real time packets for large amount of data.

We have considered the data set of the NSLKDD dataset as there are many challenges in evaluating the performance of Network intrusion detection is the unavailability of network based dataset. In this paper, we are using the machine learning models such as Naïve Bayes and Multivariate linear regression models with the ensemble method technique using the NSLKDD dataset.

## 1.1 Dataset

In the past few decades, data mining and machine learning techniques have been extensively researched in developing intrusion detection systems using different intrusion detection datasets. In 1998 DARPA Intrusion Detection Evaluation Program was prepared which contains different attacks between internet protocols. It had been managed by MIT Lincoln Labs. The main objective of the work was to survey and evaluate research in intrusion detection.

NSL-KDD Data: Are utilized in many of the machine learning methods. NSL-KDD dataset does not include redundant records within the train set; hence the classifiers will not be biased towards frequent records. The amount of records within the NSLKDD dataset train and test sets is reasonable, which makes this dataset affordable to run the experiments on an entire set. Data files in NSL-KDD are KDDTrain+.ARFF, KDDTrain+.TXT, KDDTrain+_20Percent.ARFF, and KDDTrain+_20Percent.TXT, even KDDTrain+.ARFF, KDDTrain+.TXT, KDD-Train+_21Percent.ARFF, and KDDTrain+_21Percent.TXT. There are 41 features of NSL-KDD

Table 1: Features of NSL-KDD Dataset

| Number | Feature | Type of feature | Number | Feature | Type of feature |
|---|---|---|---|---|---|
| 1 | Duration | Numeric | 22 | Is_guest_login | nominal |
| 2 | Protocol_type | Nominal | 23 | Count | numeric |
| 3 | Service | Nominal | 24 | Srv_count | numeric |
| 4 | Flag | Nominal | 25 | Serror_rate | numeric |
| 5 | Src_bytes | Numeric | 26 | Srv_serror_rate | numeric |
| 6 | Dst_bytes | Numeric | 27 | Rerror_rate | numeric |
| 7 | Land | Nominal | 28 | Srv_reerror_rate | numeric |

| 8 | Wrong_fragment | Numeric | 29 | Same_srv_rate | numeric |
|---|---|---|---|---|---|
| 9 | Urgent | Numeric | 30 | Diff_srv_rate | numeric |
| 10 | Hot | Numeric | 31 | Srv_diff_host_rate | numeric |
| 11 | Num_failed_logins | Numeric | 32 | Dst_host _count | numeric |
| 12 | Logged_in | Nominal | 33 | Dst_host_srv_count | numeric |
| 13 | Num_compromised | Numeric | 34 | Dst_host_same_srv_rate | numeric |
| 14 | Root_shell | Numeric | 35 | Dst_host_diff_srv_rate | numeric |
| 15 | Su_attempted | Numeric | 36 | Dst_host_same_src_port_rate | numeric |
| 16 | Num_root | Numeric | 37 | Dst_host_srv_diff_host_rate | numeric |
| 17 | Num_file_creations | Numeric | 38 | Dst_host_serror_rate | numeric |
| 18 | Num_shells | Numeric | 39 | Dst_host_srv_serror_rate | numeric |
| 19 | Num_access_files | Numeric | 40 | Dst_host_rerror_rate | numeric |
| 20 | Num_outbound_cmds | Numeric | 41 | Dst_host_srv_rerror_rate | numeric |
| 21 | Is_host_login | Nominal | | | |

## 1.2 Network Traffic Attacks

Table 2 Network Traffic Attacks

| Attack group | Attacks |
|---|---|
| Probe | ipsweep, mscan, nmap, portsweep, saint, satan |
| DoS | apache2, back, land, mailbomb, Neptune, processtable, pod, udpstorm, smurf, teardrop |
| U2R | buffer_overflow, httptuneel, loadmodule, perl, rootkit, xterm, ps, sqlattack |
| R2L | ftp_write, imap, guess_passwd, named, multihop, phf, sendmail, snmpgetattack, snmpguess, spy, warezclient, worm, warezmaster, zsnoop, xlock |

The network traffic attacks are basically of these four types of attack. We are considering these attacks and considering a binary value for the attack and normal data packets.
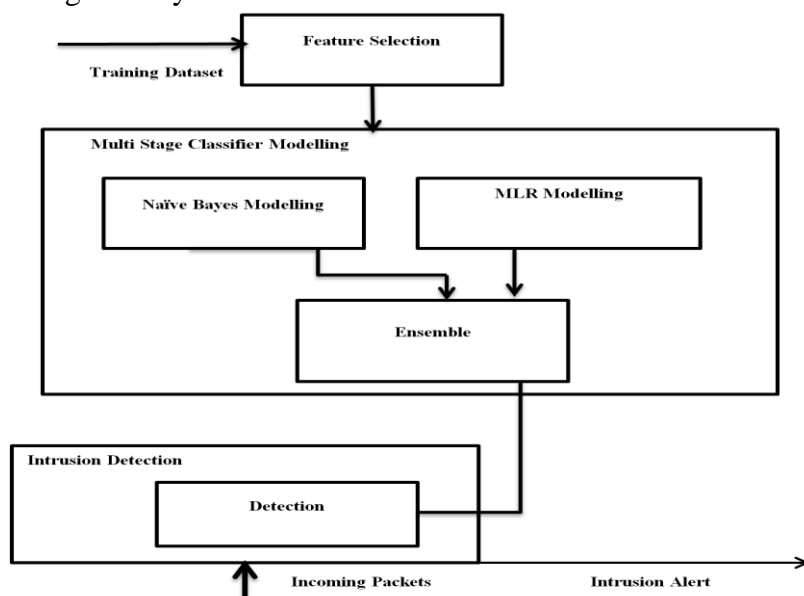
## 2. Literature Review

The study done by Srinivas Mukkamala, Guadalupe Janoski, Andrew Sung (2002) was on Support Vector machine (SVM) and Neural Network. They gave high priorities to fast training and scalability but then, restricted it to the binary classifications for few of the features in the dataset used. There was a comparative results of technique with other existing Intrusion detection system are evaluated using binary class NSL KDD dataset and multi class NSL KDD dataset in the study done by Raman Singh, Harish Kumar, R.K.Singla (2015), where the network traffic dataset was huge and imbalanced. It could give the highest accuracy rate of 97.6% precision. In the opinion given by the research done by K. A. Taher, B. Mohammed Yasin Jisan, and M.M. Rahman (2019), they presented different machine learning models using feature selection methods, but the analysis of the result was network traffic detection rate was 94.02% with the known attacks. Anomaly based intrusion detection system using Ad Hoc network was proposed by A. A. Korba, M. Nafaa and Y. Ghamri-Doudane (2016), it explained the clustered network topology. The idea was to enable the cluster head node to detect the malicious node. U. S. Musa, M. Chhabra, A. Ali and M. Kaur (2020) tried to give the review of emergence of machine learning new techniques for intrusion detection. By giving a comparison table, they are highly giving the dataset comparison and algorithms comparison. In the opinion of L. Koc and A. D. Carswell, (2010) the use Hidden Naïve Bayes algorithms for the intrusion detection system, for the attacks where they have challenges for a huge feature of

datasets, is better than the Naive Bayes algorithm. Using the KDD'99 dataset, they need to apply data mining methods and concluded that future work can be done to balance the imbalanced data in the dataset. Using Naïve Bayes along with Ada Boost to enhance network anomaly intrusion detection was proposed by W. Li and Q. Li (2010) where they could improve the detection rate of positive rate keeping the false positive rate low. By introducing class information to Feature extraction methods by Chen, Long-Shen & Syu, J.-S (2015) by keeping the classification performance and tried to reduce the computations time, but it could not be generalized for multiclass classification. By improving the performance, using ensemble methods proposed by Ngoc Tu Pham, Ernest Foo, Suriadi Suriadi (2018) using feature selection method. The techniques reduce the number of irrelevant features and classification accuracy. The drawback is that, only one dataset was used to evaluate the built classifier. To reduce the false alarm rate and true positive rate by using the Real Valued Negative Selection algorithm proposed by F.Selahshoor, H. Jazayeriy and H.Omranpour (2019). There were overlapping samples where the expressions were not accurate. In the paper proposed by J. Olamantanmi Mebawondu, Olufunso D. Alowolodu, Jacob O. Mebawondu, Adebayo O. Adetunmbi (2020) using, supervised learning paradigm method was adopted for real time intrusion detection, but it failed to test the performance of the model using different attributes. Compared to Support Vector Machine and Discriminative Restricted Boltzmann Machine the detection rate is much higher in the case of Dos, Probe, U2R and normal traffic detection this was the proposed paper given by M. Raihan-Al-Masud and H. A. Mustafa (2019). It needs to improve the detection rate of R2L. A study done by M. Azizjon, A. Jumabek and W. Kim (2020) 1D CNN based Network intrusion detection, in which they have used normalization on imbalanced data. Where, convolution neural network is used with a 3-layer model, which outperformed by achieving the highest accuracy. But, the imbalanced data problem led to poor performance of the neural network model.

## 3. System Architecture

The proposed system consists of the Feature Selection method and ensemble based machine learning methods. As shown in figure 1, Feature selection is responsible to extract most feature attributes to identify the instances. The ensemble based machine learning algorithms used are Naïve Bayes and the Multivariate Linear Regression Model. The real time data has been captured and trained using the Feature selection.

Figure 1 System Model for Network Intrusion Detection



## 3.1 Feature Selection

Feature Selection is an important process where we can automatically or manually select those features which will contribute to our prediction variable, even the output in which we are interested, as there might be irrelevant features that can decrease the accuracy of the model. The major benefit of Feature selection before modelling the data is, it reduces overfitting, improves accuracy, and reduces training time.

## 3.2 Naïve Bayes Modelling

Naïve Bayes classifier is one of the simplest and effective classification algorithms of Bayesian network classifiers, which helps in building fast machine learning models, which can make quicker predictions. It is a probabilistic classifier that is; it predicts based on the probability of an instance or object.

A Bayesian classifier maps the A features which consist of {a1, a2,.., an} into C classes that consist of {c1,c2,.., cn} on a dataset D which consists of {E1, E2,.., Et} instances and can be defined as the equation(1)

$$c(E) = arg \max_{c \in C} P(c)P(a_1, a_2, \dots, a_n|c).$$

(1)

Then, with the consideration of Naïve assumptions of the independence of the attributes given in the class as in equation (2), Naïve Bayes classifier is illustrated in figure 2

$$P(E|c) = P(a_1, a_2, \dots, a_n|c) = \prod_{i=1}^{n} P(a_i|c).$$

(2)

Figure 2 Naïve Bayes Structure



## 3.3 Multivariate Linear Regression Modelling

When we have multiple features and we need to train a model that can predict the value given for those features, multivariate linear regression can be used here. This is a similar method as a simple linear regression model with multiple independent variables contributing to determine and complex computation due to the added variables.

$$Y_i = \alpha + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + \dots + \beta_n x_i^{(n)}$$

$Y_i$ is the estimate of $i^{th}$ component of dependent feature $y$, where we have $n$ independent feature and $x_i^j$ denotes the $i^{th}$ component of $j^{th}$ independent feature. Also, the cost function is as follows

$$E(\alpha, \beta_1, \beta_2, \ldots, \beta_n) = \frac{1}{2m} \sum_{i=1}^{m} (y_i - Y_i)$$

### 3.4 Ensemble Method

Ensemble learning helps improve machine learning results by combining several methods, in our model we ensemble the Naïve Bayes and Multivariate Linear Regression. This method approaches the production of better predictive performance compared to a single model. The main advantage to use the ensemble method is, it gives an improvement in predictive accuracy.

### 4. Experiment and Result

The data is captured from the real time data from the internet as shown in figure 3. The data is trained by the algorithms, the Feature selection method is used to train for the features in particular for the attacks, then train the model using the Naïve Bayes algorithm and Multivariate algorithm. We use the Support Vector Machine where the data is partitioned into two classes: normal and attack, the objective of our experiment is to separate normal and attack patterns. In our experiment, all attacks are classified as +1, and normal data classified as -1 in our experiment we have considered 11 features of the NSL KDD data set and trained the real time packets captured and classify them and detection will be alerted as the data are captured.
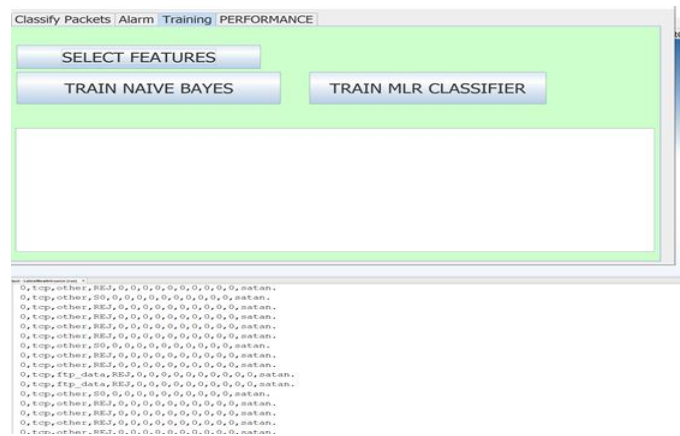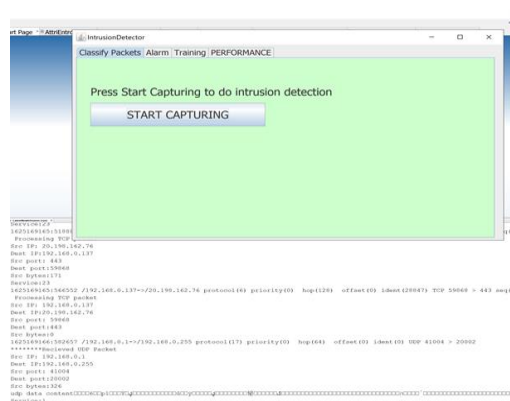
Figure 3 Feature Selection

Figure 4 Capturing the Real Time Data

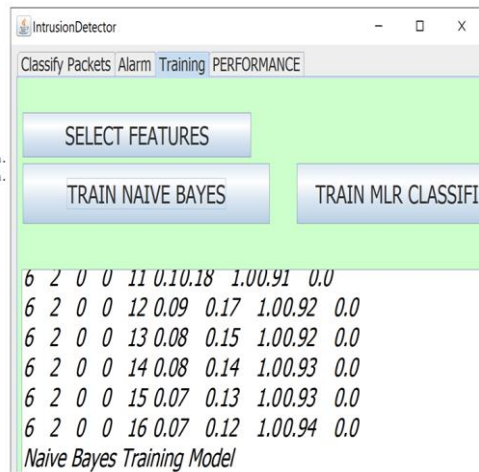Figure 5 Training the Naïve Bayes Model



Figure 6 Training the Multivariate Linear Regression Model
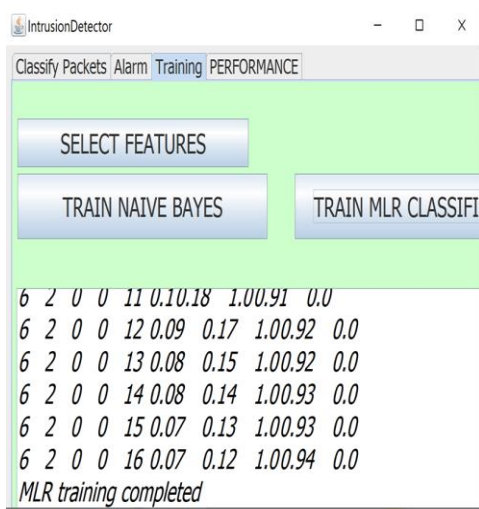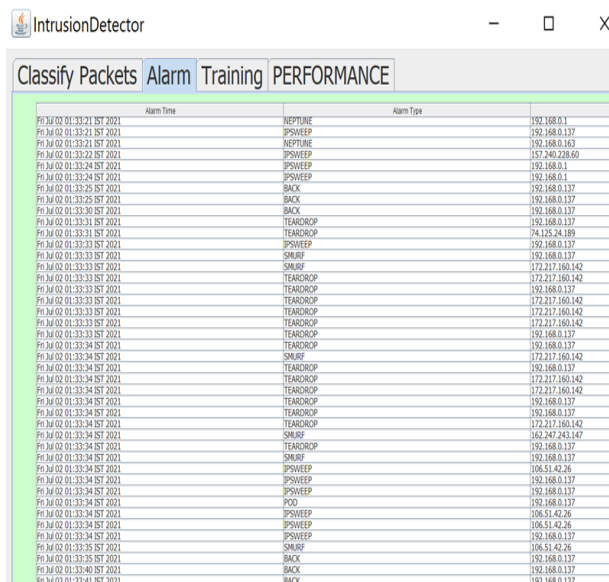


Figure 7 Alarms for the Intrusion

The result of the experiment is shown graphically. We can achieve 98% of accuracy rate. We can achieve a better accuracy than compared to other algorithms by using the ensemble method. Figure 8,9,10 shows the accuracy rate, true positive value, and False Positive value. Our experiment results in giving a high True Positive value and low False Positive value.
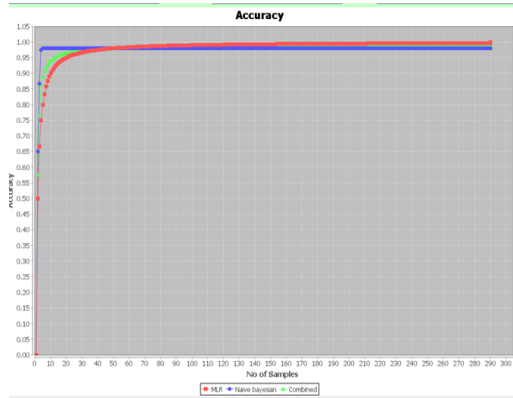
Figure 8 Accuracy Rate
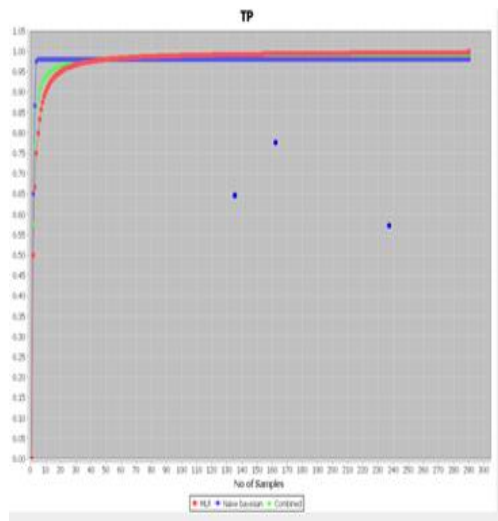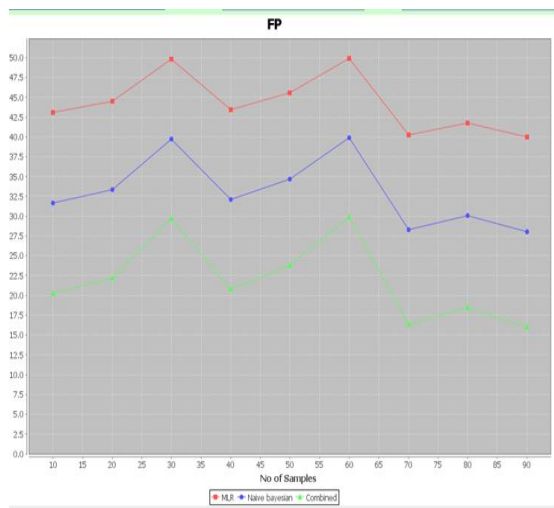


Figure 9 True Positive Values



Figure 10 False Positive Values

## 5. Conclusion

We have performed an adaptive ensemble method of machine learning to measure the intrusion detection, using the NSL KDD dataset feature for the real time packets for intrusion evaluation. The classifications were performed on a binary attack/normal basis.

The ensemble method delivers highly accurate performance, with the Naïve Bayes and Multivariate linear regression algorithms showing positive rates and accurate results. The performance is for the few features of the NSL KDD data set that have been trained, the data were trained to deliver accurate results. Comparative to [1] the result is high. The True positive value rate is high and the False positive rate is low. The accuracy rate is around 98% for the set of features trained in the model. Future work can be done for more number of features which can be checked for accuracy with a faster rate of detection for the real time data.

Future work can be taken forth for training more features from the dataset considered for training.

**References**

1. Mukkamala, S.et al(2002). "Intrusion detection using neural networks and support vector machines", IJCNN'02 (cCat.No.02CH37290) 2 (2002): 1702-1707 vol.2.
2. Raman Singh, Harish Kumar, R.K.Singla, An intrusion detection system using traffic profiling and online sequential extreme learning machine, Expert Systems with Application, Volume 42, Issue 22, 2015, Pages 8609-8624, ISSN 0957-4174, https://doi.org/j.eswa.2015.07.015.
3. K. A. Taher, B. Mohammed Yasin Jisan and M.M. Rahman, " Net-work Intrusion Detection using Supervised Machine Learning Technique with Feature Selection", 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST), Dahka, Bangladesh, 2019,pp. 643-646, doi: 10.1109/ICREST.2019.8644161.
4. A.A.Korba, M. Nafaa and Y. Ghamri-Doudane, "Anomaly-based intrusion detection system for ad hoc networks," 2016 7th Interna-tional Conference on the Network of the Future (NOF), 2016, pp. 1-3, doi: 10.1109/NOF.2016.7810132.
5. U. S. Musa, M. Chhabra, A. Ali and M. Kaur, "Intrusion Detection System using Machine Learning Techniques: A Review," 2020 Inter-national Conference on Smart Electronics and Communication (ICOSEC), 2020, pp. 149-155, doi: 10.1109/ICOSEC49089.2020.9215333.
6. L. Koc and A. D. Carswell, "Network Intrusion Detection Using a HNB Binary Classifier," 2015 17th UKSim-AMSS International Con-ference on Modelling and Simulation (UKSim), 2015, pp. 81-85, doi: 10.1109/UKSim.2015.37.
7. W. Li and Q. Li, "Using Naive Bayes with AdaBoost to Enhance Network Anomaly Intrusion Detection," 2010 Third International Conference on Intelligent Networks and Intelligent Systems, 2010, pp. 486-489, doi: 10.1109/ICINIS.2010.133.
8. Chen, Long-Shen & Syu, J.-S. (2015). Feature Extraction based Ap-proaches for Improving the Performance of IntrusionDetection Sys-tems. Lecture Notes in Engineering and Computer Science. 1. 289-291.
9. Ngoc Tu Pham, Ernest Foo, Suriadi Suriadi. Improving perfor-mance of intrusion detection system using ensemble methods and feature selection. ACSW '18: Proceedings of the Australasian Computer Science Week MulticonferenceJanuary 2018 Article No.: 2 Pages 1-6https://doi.org/10.1145/3167918.3167951.
10. F.Selahshoor, H. Jazayeriy and H.Omranpour, " Intrusion Detection system using Real-Valued Negative Selection Algorithm with Opti-mised Detectors," 2019 5th Iranian Conference on Signal Processing and Intelligent System (ICSPIS), Shahrood, Iran, 2019, pp. 1-5,doi: 10.1109/ICSPIS48872.2019.9066040.
11. J. Olamantanmi Mebawondu, Olufunso D. Alowolodu, Jacob O. Mebawondu, Adebayo O. Adetunmbi, "Network intrusion detec-tion system using supervised learning paradigm", Scientific African, Volume 9, 2020, e00497, ISSN 2468-2276, https://doi.org/10.1016/j.sciaf.2020.e00497.

12. M. Raihan-Al-Masud and H. A. Mustafa, " Network Intrusion Detection System Using Voting Ensemble Machine Learning", 2019 IEEE International Conference on Telecommunications and Photonics (ICTP), Dhaka, Bangladesh, 2019, pp. 1-4,doi: 10.1109/ICTP48844.2019.9041736.

13. M. Azizjon, A. Jumabek and W. Kim, " 1D CNN based network intrusion detection with normalization on imbalanced data" , 2020 International Conference on Artificial Intelligence in Information and Communication (ICAIIC), Fukuoka, Japan, 2020, pp. 218-224, doi: 10.1109/ICAIIC48513.2020.9064976.