

Determination of Attacks occurred on Dataset based on Intrusion Access System by using Machine Learning Techniques

Suwarna Bokade¹, Roshani Junghare², Shravani Tawale³, Neha Dhumane⁴, Nitin Mohurle⁵, Prof. Rashmi Ghate⁶

^{1,2,3,4,5,6}Department of computer engineering, Bapurao Deshmukh of College of Engineering, Sevagram

Abstract - The effectiveness of any Intrusion Detection System (IDS) system is a complex problem because of its incompatibility with multi-featured network data distribution or measurement data. To remove this situation, several a variety of intrusion access methods have been suggested and shown with varying degrees of accuracy This is why the file Selecting an effective and efficient IDS) Systems are a very important aspect of data security. In this work we have created two models for differentiation. One is based on Support Vector Machine (SVM) and the other is based on Random Forest (RF). To finding dangerous work or breaking the law is often reported, collected locally using secure data and Event management system and can also block packets.

Index Terms - IDS, Machine Learning techniques, Network Based Attacks, Various types of Attacks, Various types of classifier.

1.INTRODUCTION

Accessibility systems monitor dangerous activity networks, and they are discarded on false alarms. Therefore, companies need to adjust them properly when they first install their IDS products. This means setting up a login program to see how Compares normal traffic on the network with malicious work, and also monitors network packets entering the system detect malicious activity in which they are involved and at the same time send out alerts. Access to Partition Program:

IDS are divided into 5 types:

Network Access Program (NIDS):

Network Access Programs are installed on the network to monitor traffic on all devices on the network. It detects vehicles moving in the subnet and resembles a

road passing through the subnet through known attacks. When an attack is detected or a strange character appears, a warning is sent to the controller. An example of NIDS in a subnet where firefighters are found to find out if someone is trying to crack a firewall.

Internal Access to Immunization Program (HIDS):

Internal access control programs (HIDS) work on private hosts or network devices. HIDS scans incoming and outgoing packages. It only leaves the device and notifies a supervisor when a suspicious or dangerous job is found. It takes a snapshot of existing program files and compares them with a previous image. When system files are modified or deleted, a warning is sent to the administrator for confirmation. An example of the use of HIDS can be found on complex machines which is not expected to change its structure. The process of intrusion detection model is as following figure:

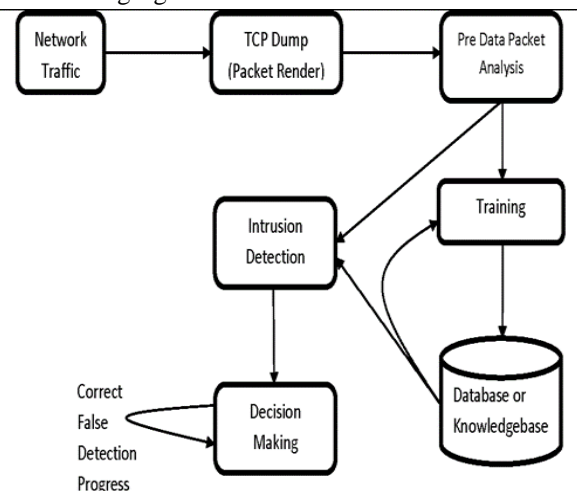


Fig 1: Architecture of Intrusion Detection

Protocol-based Intrusion Detection System (PIDS): A protocol-primarily based intrusion detection system includes a program or agent that resides at the stop of the server, retaining and interpreting the protocol among the consumer / device and the server. It seeks to defend the internet server via continuously monitoring https protocol and adopting the http protocol because https is not always encrypted and its net presentation layer must be in this connection, it is able to be used between applications before installing https.

Applications Access Applications (APIDS):

Application access (APIDS) program is a program or agent that resides within a group of servers. A particular application detects interference by monitoring and interpreting communications within the system. For example, it clearly looks at the SQL protocol displayed by Middleware because it works with a web server database.

Hybrid Access Detection System:

Hybrid intrusion detection systems are created by using combining or greater intrusion detection techniques. In a hybrid access detection system, the host agent or machine facts is blended with the network statistics to develop a complete assessment of the network system. The hybrid intrusion detection system works well compared to other intrusion detection systems. Introduction An example of a hybrid IDS.

2. TECHNIQUES FOR IMPLMENTING HIDS SYSTEM

Most intrusion prevention programs use one of three access methods: signatures, uncontrolled statistics, and explicit protocol analysis. We introduced how to use data to improve IDS in this study, known as IDS. Data processing has two components: sampling and job selection. First used sample data and a combination of Genetic Algorithm and RF to improve sample size. In job selection, Genetic Algorithm and RF combinations are also used to identify the best performance subset. Use RF to upgrade IDS to create isolation Data for intrusion detection.

IDS is a method of testing RF tested based on data usage; The IDS works on all subject-selected indicators much better than the RF classificatory,

which demonstrates the importance of data usability in IDS. In addition, it reveals that RF segregation is a very robust division by comparing conventional machine learning methods, so the combined effect of data usage and RF classification makes IDS almost always the best, especially for non-compliant identification with lower records, those doses, analysis, back, and staff. However, enhancements are still possible, such as long-term costs in the field of data usage and online processing support.

Because the proposed use of data can effectively reduce the impact of unequal distribution of samples on IDS and indicate encouraging performance, additional areas of uncommon acquisition such as fraud detection can also be used. In addition, the search process may be improved as it takes more time to train classifiers. We first applied that we collect CICIDS dataset and then we train and test the dataset: The training and testing technique is shown in following figure:

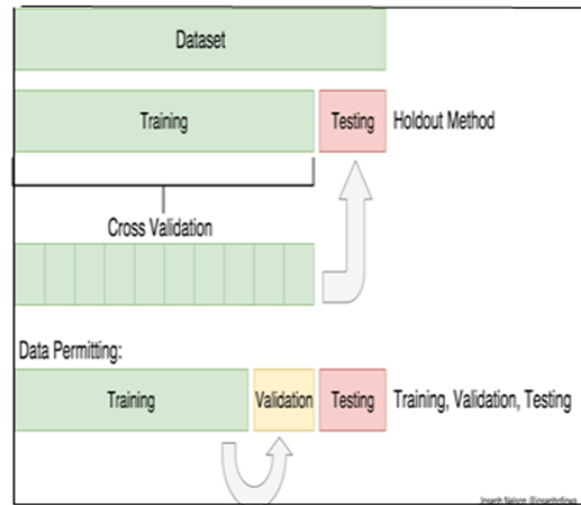


Fig 2: Process of Training and Testing Dataset We split dataset i.e. 80% training data and 20% tested dataset. The training set is the one in which we trained and measure our model equally to fit the parameters and the tested data is only used to test the performance of the model. Training data releases are available for modeling and test data is unrecognized data for predictability points monitors all incoming and outgoing traffic flowing to or from network devices. RF elegance feature is used as a subcommittee. The motive of the function selection is to find a hard and fast of features which could enhance the performance of the receiver.

3 RECOGNITION OF ATTACKS ON DATASET

According to the KDD database there are 21 types of attacks divided into four groups (DOS, R2L, U2R, and PROBE) with different numbers of cases and database appearances. Based on an in-depth analysis of KDD data the distribution event for different types of attacks has been saved. In other words, 79% of the data extracted to initiate DOS attacks with 19% of standard traffic while 2% of other types of entry (U2R, R2U and PROBE).

In networks, the general behavior of users exceeds the undesirable behavior, which makes the data-sharing of normal and unethical behavior unequal. To improve IDS acquisition performance, A hybrid information expansion approach is based totally on gadget getting to know algorithms. The records processing method consists of two elements: data sampling.

1. Data sampling: In this section, the iForest acquisition method is used for data collection, GA is used to increase the size of the global sample, and the performance of the RF section of student sample data is used as a test indicator. The purpose of the data model is to obtain an appropriate training database and reduce data inequalities.
2. Feature Selection: In this paper, the GA and RF integration technique is used to select homes. As a records version, GA is used as a seek method to determine subsets of an electoral responsibility, and RF elegance feature is used as a subcommittee. The motive of the function selection is to find a hard and fast of features which could enhance the performance of the receiver. The cause of the svm set of rules is to create a dotted line or boundary line that may divide n-dimensional space into a circle so that future.

4. MODEL CLASSIFIER

The steps of the proposed model can be summarized as follows:

- 1 Upload the database and export it to Distributed Database (RDD) and Data Frame in Apache Spark.
- 2 Pre-data.
- 3 Feature selection.
- 4 Train SVM and training database.
- 5 Explore and test KDD model and database.

4.1 Random Forest Classifier: The Random Forest is a well-known form of mechanized learning of supervised learning technology. It can be used for both tax problems and ML backlogs. In terms of mass learning, a process that involves many variables in solving a complex problem and improving model performance.

As the name suggests, "Random Forest is a subdivision that contains multiple decision trees in the lower set of a given database and takes a measure to improve the approximate accuracy of that database. "Instead of relying on a deciduous tree, the unplanned forest takes one of each tree.

We find out accuracy, precision, recall of dataset by using Random Forest algorithm and RF is best classifier than SVM algorithm. The Accurate classification of the proposed model is superior to models that use RF partitions where the parameters are selected therefore better performance of general performance than SVM we propose a novel framework called the hybrid internal detection system. The cause of the svm set of rules is to create a dotted line or boundary line that may divide n-dimensional space into a circle so that future.

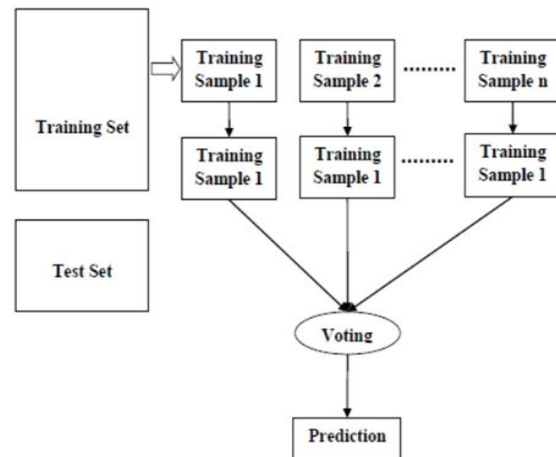


Fig3: Working of RF Classifier

4.2 Support Vector Machine Classifier: Vector Machine or SVM is one of the most widely studied algorithms used to diagnose and retrieve problems. The cause of the svm set of rules is to create a dotted line or boundary line that may divide n-dimensional space into a circle so that future new information points are without problems positioned in the precise orbit.

5. PERFORMANCE AND RESULT ANALYSIS

The performance of overall experiment is as follows:

5.1: Collection of CICIDS Dataset: We put CICIDS (Friday morning working hours) dataset from local area network. This dataset has various attacks such as DOS, R2L, SQL injection, etc.

5.2 Data Sampling: Data sampling is done by training dataset and testing dataset. Real-world data often contains sounds, lost values, and perhaps in an unusable format that cannot be used directly by machine learning models. Data processing is required for data refining operations and optimization of machine learning models, which also increases the accuracy and efficiency of machine learning models.

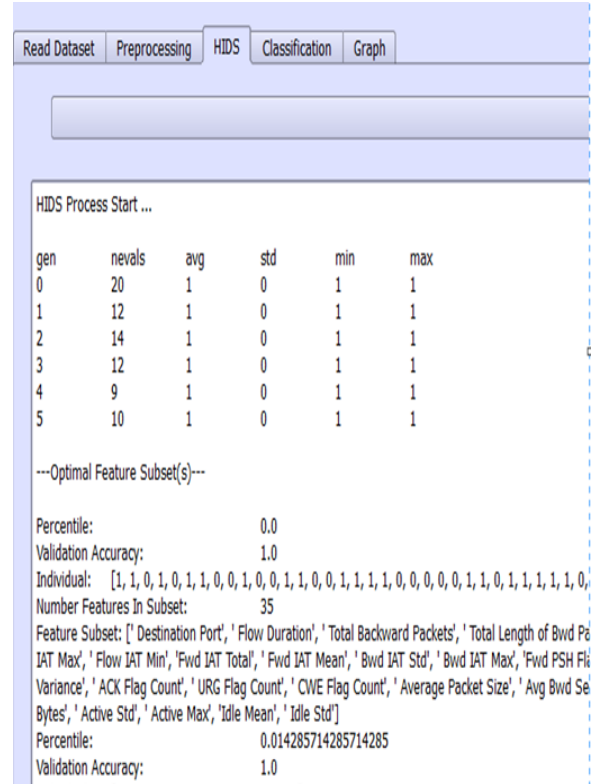
5.3 Feature Selection: Selection of feature in machine learning to find a set of features that allow one to create useful models for learning objects. There are some features selected out of 34 features. There are some features is shown in following table that we selected. Feature selection in machine learning to find a set of features Allow objects to create useful patterns for learning. Some of the 34 features were selected. The following table shows some of the features we have selected.

Feature name	Type	Description
Flow duration	continuous	duration of the flow in microsecond
total Fwd Packet	continuous	total packets in the forward direction
total Bwd packets	continuous	total packets in the backward direction
total Length of Fwd Packet	continuous	total size of packet in forward direction
total Length of Bwd Packet	continuous	total size of packet in backward direction
Fwd Packet Length Min	continuous	minimum size of packet in forward direction
Fwd Packet Length Max	continuous	maximum size of packet in forward direction
Fwd Packet Length Mean	continuous	mean size of packet in forward direction
Fwd Packet Length Std	continuous	standard deviation size of packet in forward direction
Bwd Packet Length Min	continuous	minimum size of packet in backward direction
Bwd Packet Length Max	continuous	maximum size of packet in backward direction
Bwd Packet Length Mean	continuous	mean size of packet in backward direction

Table 1: Selected Features from CICIDS Dataset

5.4 Implementation of Hybrid Intrusion Detection: We are proposing a new framework called Hybrid Intrusion Access System (H-IDS) to detect various attacks. In this system, we use both random-based identification methods and signatures to achieve more accurate identification Parameters.

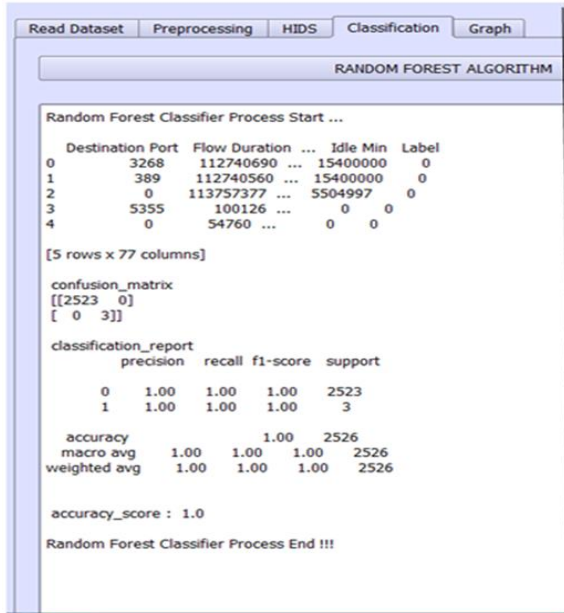
The exact classification of the proposed model is better than the models using the selected RF segmentation, so we use the genetic algorithm for the selection of features for better performance of the generalization compared to SVM. And finding the dataset is malicious or not. The process of Intrusion detection system is shown in following result:



SCR 01: Process of Hybrid Intrusion Detection System

5.5 Classification of RF and SVM classifier:

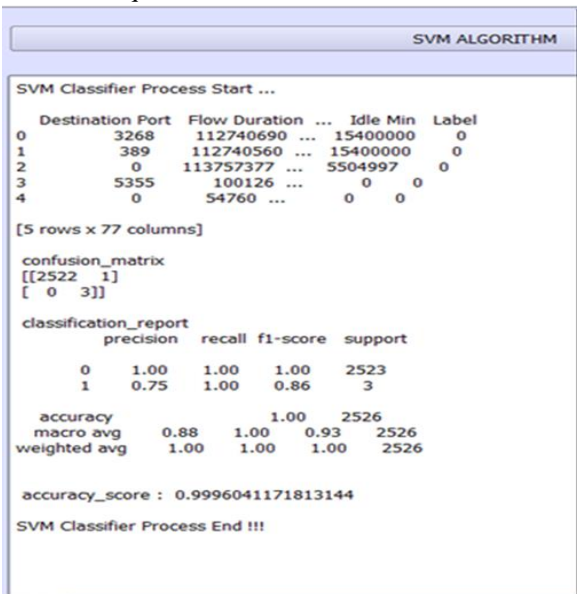
The Accurate classification of the proposed model is superior to models that use RF partitions where the parameters are selected therefore better performance of general performance than SVM we propose a novel framework called the hybrid internal detection system. DOS, probe, U2R, R2L11. Random fore (RF) is a composite separator we use both random-based identification methods with the forest planning problem Random Forest offers you the opportunity to belong to a class. SVM takes you off the line, you still need to turn it into opportunities in some way if you need opportunities. In those cases, where SVM works, it works better than Random Forest. The result of RF algorithm is as follows:



SCR02: Process of IDS using RF Algorithm

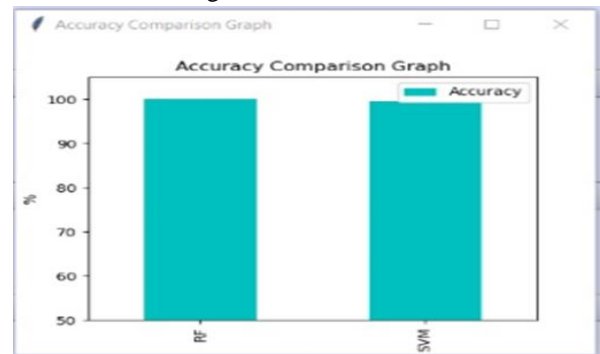
Support vector machines (SVMs) have become one of the most popular ML algorithms used to detect intrusion due to its aesthetic nature and ability to deal with curses of size. As cited by various researchers, the scale still affects the performance of SVM-based IDS. We find out that SVM classifier is very difficult for detection of intrusion on this dataset. RF algorithm used for partitioning and retrieval. But more often, they are used. We find out RF has more time than SVM for detecting intrusion on dataset. Therefore, we find SVM is better than RF algorithm for detect malicious data.

The time required for SVM is less than RF classifier.



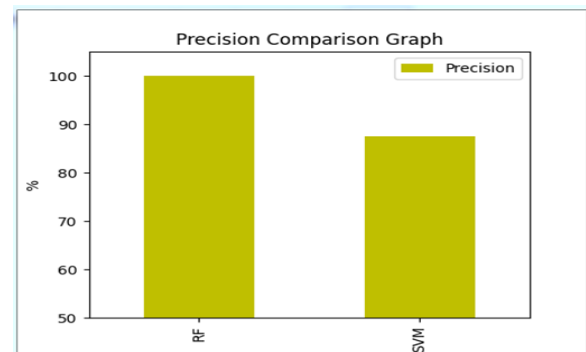
SCR03: Process of IDS using SVM algorithm

Accuracy comparison Graph the following Graph shows that the result of comparison between RF and SVM Classifier which is used for intrusion detection on dataset that shows Accuracy, precision, Recall, etc. The Accurate classification of the proposed model is superior to models that use RF partitions where the parameters are selected therefore better performance of general performance than SVM. We find accuracy of RF is 99.95 and SVM is 99.45 for detecting intrusion on the CICIDS dataset. The following graph shows that comparison between RF and SVM algorithm i.e., Accuracy, Precision, Recall etc. of intrusion detection system. we see is that the computer The complexity of support vector machines (SVMs) is much higher than that of random forests (RF). This means that SVM training is longer training than RFI when the amount of training data is high. This is shown in following:



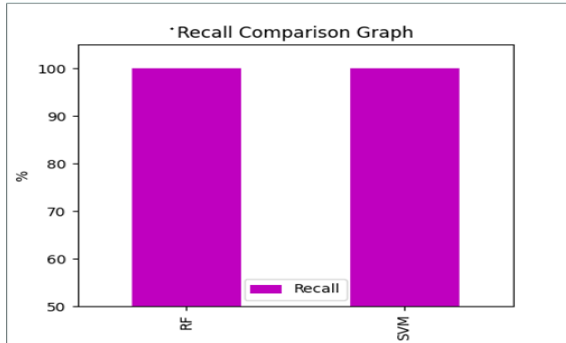
SCR04: Accuracy Comparison Graph

Precision Comparison Graph: We find out precision of RF partition is greater than SVM partition. The accuracy of the measurement system, which is related to reproduction and multiplication, is the rate at which repeated measurements under fixed conditions show similar results. The precision result is shown is as follows:



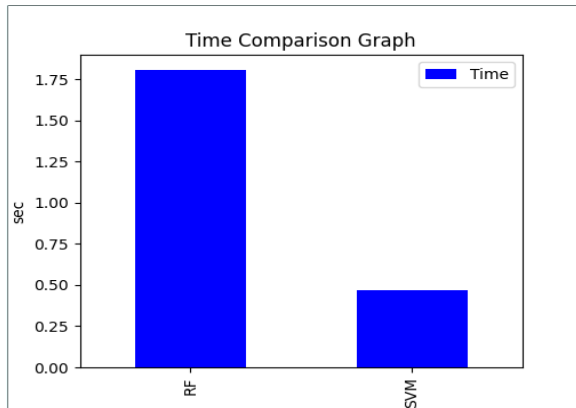
SCR06: Precision Comparison Graph.

Recall Comparison Graph: We find out RF has greater recall than SVM algorithm. while recollection (also known as sensitivity) is part of the appropriate conditions found Both accuracy and recall are based on correlation. We find out RF has greater recall. The comparison of recall of RF and SVM algorithm is shown in following graph:



SCR05: Recall comparison graph

Time Comparison Graph: We find out RF has more time than SVM for detecting intrusion on dataset. Therefore, we find SVM is better than RF algorithm for detect malicious data. The time required for SVM is less than RF classifier.



SCR07: Time Comparison Graph

F1 Comparison Graph: In the data model, chromosome sequence = $x, 1, 2, \dots, K$ is the variety of classes of network behavior, zero.1, zero.2, zero.3, 0.4, 0.5, zero.6, 0.7, 0.8, zero. Nine a is the gene at the chromosome that determines the ratio of rectangular outliers to iforest. Inside the classification problem, the health characteristic is generally set to the accuracy of the type. Fitness characteristic is taken into consideration as f1 rating. The f1 rating is a harmonic feature that takes under consideration both accuracy and recall. The F1 score count is shown as follows.

$$F1 \text{ score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

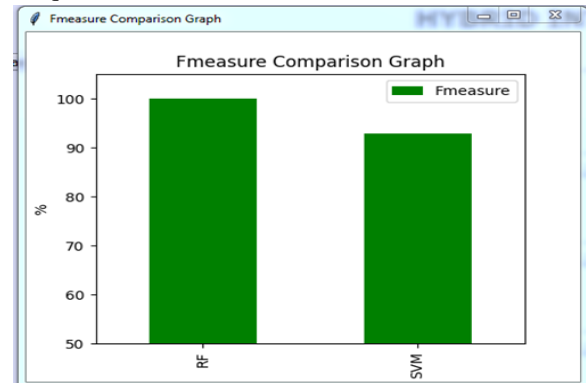
$$\text{Recall} = \frac{TP}{TP + FN}$$

$$+ (3)$$

$$\text{Memory} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$+ (4)$$

Among them, real fine (TN) is the number of authentic bizarre records categorized as inequalities, real poor (TN) is the variety of real ordinary statistics classified as ordinary, and fake positive (FPP) is classified as proper normal. Number of information. The wide variety of actual anomaly facts categorized as inequalities and fake bad (FN) are common.



SCR08: F1 Score Comparison Graph

6. CONCLUSION

We have presented a method of data optimization to develop IDS in this research, known as IDS. The optimization of information consists of components: sampling and choice of features. iforest is used for statistics sampling and integration of GA and RF for the optimization of the pattern ratio. In the choice of capabilities, GA and rf integration is once more hired to become aware of the fine feature subset. Use RF to develop IDS to carry out classification. The intrusion detection data CICIDS was appreciated for IDS. We also conclude that there are many types of attacks on dataset. and RF classifier is best for detection of intrusion on the CICIDS dataset than the SVM.

REFERENCES

- [1] A feature selection model for network intrusion detection system based on PSO, GWO, FFA and GA algorithms O Almomani - Symmetry, 2020.
- [2] Gillala Rekha1, Shaveta Malik2, Amit Kumar Tyagi3, a), Meghna Manoj Nair3,"Intrusion

- Detection in Cyber Security: Role of Machine Learning and Data Mining in Cyber Security”, Volume 5, Issue 3, Page No 72-81, 2020.
- [3] Beigh BM, Peer MA (2012) Intrusion detection and prevention system: classification and quick review. *ARNP J Sci Technol* 2(7):661–675
- [4] Vilela, Douglas W. F. L.; Lotufo, Anna Diva P.; Santos, Carlos R. (2018). Fuzzy ARTMAP Neural Network IDS Evaluation applied for real IEEE 802.11w data base. 2018 International Joint Conference on Neural Networks (IJCNN). IEEE.
- [5] S Krishnaveni, P Vigneshwar, S Kishore, B Jothi, Anomaly-Based Intrusion Detection System Using Support Vector Machine, springer-2020.
- [6] Saurabh Kumar, VIT University "Intrusion Detection System using Random Forest", April 2019. Iqbal H. Sarker 1,2,, Yoosef B. Abushark 3, IntruDTree: A Machine Learning Based Cyber", Received: 31 March 2020; Accepted: 15 April 2020; Published: 6 May 2020.
- [7] Yirui Wu ,1 Dabao Wei,1 and Jun Feng, "Network Attacks Detection Methods Based on Deep Learning Techniques: A Survey", Published 28 Aug 2020.
- [8] Zhihong Tian, Chaochao Luo, Jing Qiu, Xiaojiang Du, Mohsen Guizani, "A distributed deep learning system for web attack detection on edge devices”, *IEEE Transactions on Industrial Informatics* 16 (3), 1963-1971, 2019.
- [9] H. Liu, B. Lang, M. Liu, and H. Yan, “Cnn and rnn based payload classification methods for attack detection,” *Knowledge-Based Systems*, vol. 163, pp. 332–341, 2019.
- [10] Nancy Agarwal 1 and Syed Zeeshan Hussain1, "A Closer Look at Intrusion Detection System for Web Applications", Published 14 Aug 2018.
- [11] Stacy Stanford, Roberto Iriondo, Pratik Shukla "Best Public Datasets for Machine Learning and Data Science", January 6, 2021.
- [12] SN Mighan, M Kahani, A novel scalable intrusion detection system based on deep learning, spriner, 2020.
- [13] W. –C. Lin, Shih-Wen K. Chih-Fong T. (2015). Intrusion detection system based on combining cluster centers and nearest neighbors. *Knowledge-Based Systems* 78 (pp. 13-21). Elsevier.
- [14] Nabila Farnaaz and M.A Jabbar. (2016). Random Forest Modeling for Network Intrusion Detection System. *International Multi-conference on information processing (IMCIP) 12* (pp. 213-217). Elsevier.
- [15] Kayvan A. Saadiah Y. Amirali R. and Hazyanti S. (2016). Anomaly Detection Based on Profile Signature in Network using Machine Learning Techniques. *IEEE TENSYPMP*. (pp. 71-76). IEEE.