



Intrusion Detection System using PCA with NN Approach

M.Tech. Scholar Megha Tomar, Asst. Prof. Avinash Pal , Trapti Ozha(HOD), Director Durgesh Mishra
 Department of Computer Science and Engineering Sri Aurobindo Institute of Technology, Indore,India
 meghatomar1008@gmail.com , avinash.pal@sait.ac.in , Trapti.ozha@sait.ac.in , durgesh.mishra@sait.ac.in

Abstract: On the internet, harmful behaviours that harm a single machine or the entire network may occur. They may be random. Computer connection is growing at an unprecedented rate, making it tougher to stay up. The internet may pose security risks, just as it might in person. IDS software monitors a network for suspicious or malicious behaviour. IDS is a technology that helps detect system assaults and identify their perpetrators. In the past, machine learning (ML) approaches were employed to improve IDS accuracy and intruder detection outcomes. In this post, we discuss how to build an IDS using PCA and NN classification method. This is a strategy for building a good IDS. PCA may be used to minimise the dimensionality of the data and organise it, whereas random forest can categorise it. The suggested procedure will be applied on the KDD (Knowledge Discovery Dataset). It's evident that the offered approach is more accurate than SVM, Naive Bayes, and Decision Tree. One can be confident of this because: Our technique yielded the following results: Performance time is 2.25 minutes, accuracy is 98.35%, and error rate is 0.12.

Keywords: IDS, Knowledge Discovery Dataset, PCA, Random Forest.

I. INTRODUCTION

Because of the rapid acceleration of technological progress, the internet's presence in people's daily lives is rapidly expanding and becoming more pervasive. The internet has become so ubiquitous that almost everyone's life is reliant on it in some way, shape, or form. These days, having an internet connection and knowing how to use it are both necessities for everyone. As a result, protecting the system from potentially hazardous behaviour is becoming more and more necessary as the number of individuals utilising the internet for their own personal activities continues to rise.

There are many different kinds of attacks that may be observed on the system or the network. The purpose of these assaults is to either steal information from a system or modify the data that is already there on any system [1]. Attackers will employ a wide number of methods to get access to the system's data in order to misuse it. Some of these methods include denial of service (DoS), probe, sniff, r2l, and more methods that are quite similar. As a direct

result of this, an intrusion detection system was built and installed to guard the system from attacks of this kind. System intrusion detection systems, often known as IDS, are responsible for monitoring any attacks made against the system and working to protect it from any dangers. Intrusion Detection System, Version 1.1:

Intrusion is a term that refers to the act of entering a system without authorisation and inflicting harm to the information that is held inside of the system [1]. Any system that has this infiltration risk having its hardware compromised as a result of the intrusion. The phrase "keeping the system from being compromised" has come to place a great deal of importance on the word "intrusion," which has developed into a very important noun. Depending on the circumstances, the intrusion detection system (IDS) may either be used to control any intrusions that take place inside a system or keep track of any intrusions that take place within the system. In spite of the fact that several forms of intrusion detection systems have been used in the past, the accuracy of each method has been called into doubt over the course of the previous few years. In order to establish the system's level of accuracy, it is necessary to investigate two separate metrics, namely the detection rate and the false alarm rate [2]. It is essential that the system be developed in such a manner that the rate of false alarms is reduced to an absolute minimum while simultaneously increasing the detection rate. As a direct consequence of this, the IDS makes use of the PCA in combination with the random forest.

The IDS may take one of two forms in nature; both of these forms are relevant to its usefulness, and they are as follows:

In this system, the network traffic is analysed, and any intrusions that occur as a consequence of the traffic are detected and examined. These actions take place in response to the monitoring of the network traffic.

Host-based intrusion detection systems, often known as HIDS, are systems that are used in the process of detecting

network intrusions. These systems monitor system files that are accessible over the network.

An additional point to consider is that there is a subset of IDS types. Those variants that depend on signature identification and anomaly detection are the ones that are seen the most often.

Signature-based: In this instance, the system identified a number of specific patterns that malicious software employs in order to conceal its true nature. The patterns that have been uncovered are referred to as signatures. When it comes to detecting fresh assaults, this does not perform as well as it does when recognising already committed crimes; the signature detection procedure is where it falls short.

A kind of detection called as anomaly-based detection is one that has been developed expressly for the purpose of identifying unknown forms of attack. This system, which incorporates ML into its operation, is used to create the model.

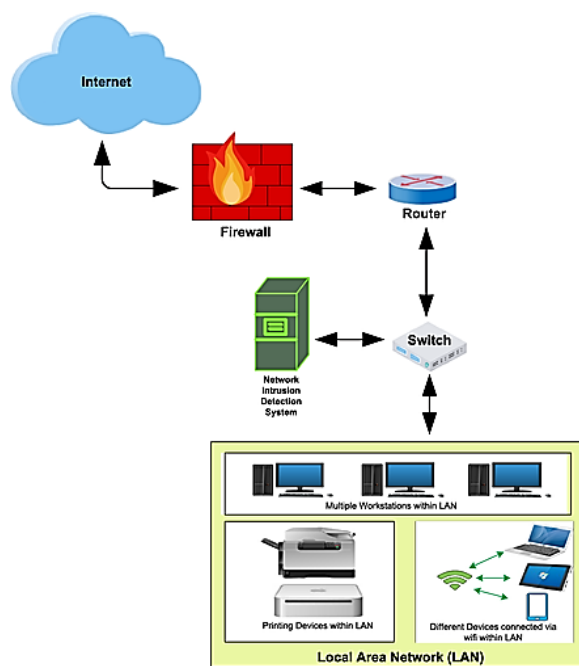


Figure 1. Intrusion Detection System[2]

1.2 Random Forest:

When it comes to solving classification issues, RF is among the most effective approaches that can be found in the field of machine learning. Random forests are classified as supervised classification algorithms [2], which is the category that the random forest belongs to. This procedure is carried out in two distinct stages: the first stage is concerned with the formation of the forest using the dataset that has been provided, and the second stage is concerned with the prediction using the classifier that was generated in the very first stage.

The following is a pseudocode representation of how to generate a random forest:

1. Pick a few characteristics, k , from a total of m to represent k/m .
2. By applying the split point to k different characteristics, get node d
3. In order to get the daughter nodes, use the optimal split.
4. Continue with steps 3, 2, and 1 until we reach the first node.
5. Establish a forest by repeating the procedures one through four in order to establish a forest.

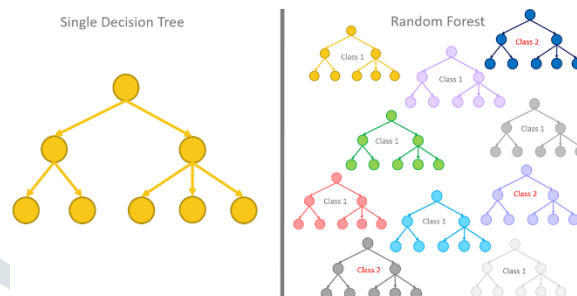


Figure 2. Random Forest Model.

1.3 PCA:

One of the methods that is utilised, in particular for the goal of lowering the size of the dimensions of a given dataset, is known as the principal component analysis technique. In addition to being one of the most effective and accurate methods currently available for reducing the dimensionality of data, it is also capable of producing the outcomes that are needed [3]. With the use of this method, the properties of a particular dataset may be whittled down to an acceptable amount of attributes, which are referred to as the primary components.

This method considers the whole input as a dataset, which since it has a big number of attributes also contains a large dimension. Since the dataset contains a large number of characteristics, this method has a large dimension. This strategy helps to reduce the amount of the dataset by positioning all of the data points along the same axis. Following the repositioning of the data points along one of the axes, the principal components analysis is performed.

II. LITERATURE REVIEW

The authors claim that a solution to the IDS was discovered by using the SVM and Nave Bayes algorithms, with the SVM algorithm showing to be better to both of the other approaches that were tried. Their experiment was carried out by utilising information from the KDD dataset, and they provide their results in terms of the detection rate as well as the false alarm rate. [4]

The writers of this paper carried out three independent experiments, each of which is described in extensive detail in the authors' respective paragraphs. Both in the analysis and in the design, they made use of the feature selection technique. In addition, the naive Bayes algorithm, the adaptive boost algorithm, and the partial decision tree were presented. They investigated each and every option for detecting infiltration that was accessible. [5]

The writers of this article have arrived at the conclusion, with the assistance of this article, that the approach of artificial neural networks with feature selection would

produce higher results when compared to the Support vector machine method. The NSL-KDD dataset was used throughout the course of the project. It turned out that the approach that was offered was effective. [6]

The purpose of this study is to provide an overview of intrusion detection systems that employ a machine-learning algorithm, as described by the authors of this paper. On the basis of the findings that they acquired, the authors offered a performance comparison of several machine learning methods. They examined the survey's detection rates as well as its rates of false alarms so that they could evaluate it. [7]

The authors have created a method for detecting intrusions that takes use of logistic regression and belief propagation methods. This approach was developed by the writers. As was said earlier, the proposed technique has shown that it provides a quicker average detection time in comparison to the approaches that have been used in the past. [8]

The authors created a technique of in-depth learning that was used to extract features from the dataset. This approach was successful in accomplishing this task. They made an effort to extract features from a dataset in order to make the dataset more useful for utilisation, and as a result of their efforts, they arrived at the realisation that they might provide improved input to the intrusion detection system. [9]

They used a process known as machine learning to carry out the research for this segment, which focused on intrusion detection systems. In their research, the authors examined all of the machine learning algorithms that have been used up to this point. Based on their findings, they determined that the machine learning algorithms that were supplied by Md. Nasimuzzaman Chowdhury and the ANN algorithms that were provided by Alex Shenfield, Aladdin Ayesah, and David Day were the most successful. [10]

The authors of this study evaluated a variety of different machine learning strategies that may be used into an intrusion detection system. They examined a variety of approaches, some of which were the SVM, the Extreme learning machine, and the random forest, amongst others. The authors come to the conclusion that the Extreme machine learning methodology is superior to every other method by a significant margin in their conclusion. [11]

The authors of this study made an effort to improve the quality of the dataset before making it accessible for examination by the intrusion detection system. In order to improve upon the dataset, which they have previously discussed, they have implemented a fuzzy rule-based feature selection approach. They used the KDD dataset, and the findings of the IDS revealed a dynamic rise in the total number of outcomes. [12]

III. PROBLEM DOMAIN

The computer networks that are connected to the internet are vulnerable to a broad variety of malicious activities. The intrusion into the system with the intention of breaking the information has been deemed to be the most significant

problem that has been seen in this sector. The construction of an intrusion detection system allows for the identification of this incursion; nevertheless, in order for this system to be useful, it must be accurate and efficient in its detection of intruders. For the purpose of intrusion detection, machine learning algorithms were used. These approaches included SVM, Naive Bayes, and other variations of the methodology. In spite of this, the data suggest that there is perhaps some space for development in a variety of areas, including but not limited to accuracy, detection rates, and the occurrence of false alarms. It is possible that approaches that have been utilised in the past may need to be replaced with new methods. Some of these new methods are SVM and Nave Bayes. In addition, the study suggests that the dataset may be improved in some way if certain procedures were applied to it. It is required to in order to improve the overall quality of the data that is fed into the suggested system.

IV. PROPOSED SOLUTION

The purpose of the intrusion detection system is to improve the overall performance of the system, which is negatively influenced by the presence of intruders. This apparatus is able to recognise suspicious activity on the grounds and alert the appropriate personnel. The proposed method makes an effort to address the problems that have surfaced as a consequence of the work that has been done before. One of the proposed components of the system is known as principal component analysis, while the other is known as random forest. It is recommended that both of these components be included in the system. The dimension of a dataset can be reduced in size through the application of the technique of principal component analysis; the quality of the dataset will be improved as a result of this approach due to the fact that the dataset will include the appropriate characteristics as a result of this method. The SVM approach will not be employed for intruder detection because the random forest algorithm, which has a greater detection rate and a lower false alarm rate than SVM, will be used instead.

4.1 Algorithm for the proposed solution:

The coordination degree of the original attribute for the split node standard has been replaced with compatibility as the replacement attribute.

1. Attribute compatibility

Let the modulus for the primary decision set be "Pr," and the modulus for the secondary set be "Se." The definition of attribute compatibility is as follows:

$$CO(X \rightarrow D) = \frac{|P_r| - |S_e|}{|X|} \quad (1)$$

X is the subset for non-empty C in this context. When the impact of the secondary set can be demonstrated to have over the primary set, we term this situation "strict compatibility." The primary text and the additional set both include errors or inconsistencies. The expression puts the finishing touches to the secondary set.

$$CO(X \rightarrow D) = \frac{|Pr|}{|X|} \quad (2)$$

here X is the subset for non-empty C. In this, the wide compatibility of the second set is seen.

Algorithm for The Base Classifier Improvement:

In the first step, the active attribute of the data set is given its starting value by marking every condition attribute.

The second step is to compute the modulus for each condition attribute that is included in either the main or secondary set.

Phase 3: The computation of compatibility between all conditional attributes is carried out in this step by applying equation (1). If other characteristics with a similar compatibility are observed, equation (2) should be used.

Step 4: To divide the sample, pick the node with the most comprehensive compatibility for splitting to use as the split node, and then remove the tag that is now active.

Step 5: continue picking the active attribute for splitting up to the point when we receive the active attribute when we reach the leaf node.

The sixth and last step is the generation of the basic classifier.

4.2 Flowchart for The Proposed Algorithm:

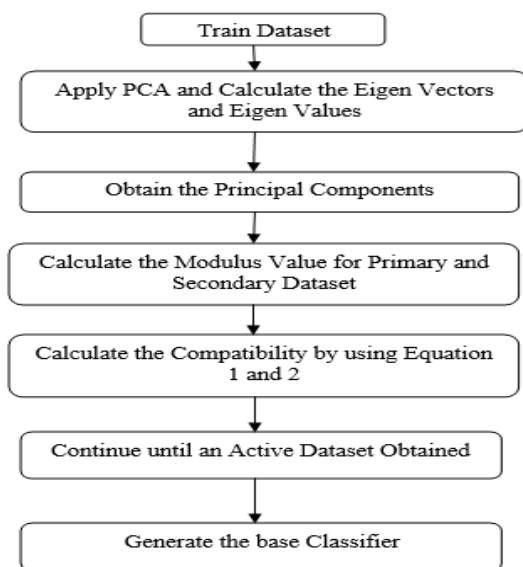


Figure 3. Flowchart for The Proposed Approach

V. RESULTS

It was feasible to acquire good results with the experiment that was carried out to test the recommended approach when the KDD dataset was used as the foundation for the experiment. The use of principal component analysis (PCA) in combination with convolutional neural networks (NN) fared very well in comparison to other prevalent approaches such as SVM, Naive Bayes, decision trees, and NN. In the next part, a table will be shown that compares the performance time (in minutes), accuracy rate (in percentage), and error rate (in percentage) of a number of different approaches:

Table 2. Result Comparison with other Classifiers

Method	Performance time (min)	Accuracy rate (%)	Error rate (%)
SVM	4.57	84.34	2.67
Naïve Bayes	9.12	80.85	3.49
Decision Tree	12.36	89.91	0.78
PCA with NN	2.25	98.35	0.12

When compared to other algorithms such as SVM, Naive Bayes, and Decision Tree, it is clear that the technique that has been provided here is superior in terms of performance. PCA combined with Random Forest performed much better based on three different factors. Its representation may be seen in table 1.

VI. CONCLUSION

In recent years, there has been an increase in the number of systems that are connected to the internet, which has been followed by an increase in the number of security concerns. The solution that was recommended deals with the detection of intruders successfully via the use of the internet, and it does so in a cost-effective manner. The approach that was recommended performed much better than other algorithms that have been utilised in the past, such as SVM, Naive Bayes, and Decision Tree. The approach that was presented has the potential to dramatically raise, in a number of different ways, not only the detection rates but also the rates of false errors. The knowledge discovery dataset is the one that has been used for the purpose of demonstrating this example. Using the strategy that we suggested, we were able to acquire the following results: The performance time (in minutes) is 3.24, the accuracy rate (in percentage) is 96.78, and the error rate (in percentage) is 0.21. The performance time (in minutes) is 3.24, the accuracy rate (expressed as a percentage) is 96.78, and the error rate (expressed as a percentage) is 0.21.

Reference:

1. Jafar Abo Nada; Mohammad Rasmi Al-Mosa, 2018 International Arab Conference on Information Technology (ACIT), A Proposed Wireless Intrusion Detection Prevention and Attack System
2. Kinam Park; Youngrok Song; Yun-Gyung Cheong, 2018 IEEE Fourth International Conference on Big Data Computing Service and Applications (BigDataService), Classification of Attack Types for Intrusion Detection Systems Using a Machine Learning Algorithm
3. Le, T.-T.-H., Kang, H., & Kim, H. (2019). The Impact of PCA-Scale Improving GRU Performance for Intrusion Detection. 2019 International Conference on Platform Technology and Service (PlatCon). Doi:10.1109/platcon.2019.8668960
4. Anish Halimaa A, Dr K.Sundarakantham: Proceedings of the Third International Conference on Trends in Electronics and Informatics (ICOEI 2019) 978-1-5386-9439-8/19/\$31.00 ©2019 IEEE "MACHINE LEARNING BASED INTRUSION DETECTION SYSTEM."
5. Rohit Kumar Singh Gautam, Er. Amit Doegar; 2018 8th International Conference on Cloud Computing, Data Science & Engineering (Confluence) " An Ensemble Approach for Intrusion Detection System Using Machine Learning Algorithms."

6. Kazi Abu Taher, Billal Mohammed Yasin Jisan, Md. Mahbubur Rahma, 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST) “Network Intrusion Detection using Supervised Machine Learning Technique with Feature Selection.”
7. L. Haripriya, M.A. Jabbar, 2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA) “Role of Machine Learning in Intrusion Detection System: Review”
8. Nimmy Krishnan, A. Salim, 2018 International CET Conference on Control, Communication, and Computing (IC4) “Machine Learning-Based Intrusion Detection for Virtualized Infrastructures”
9. Mohammed Ishaque, Ladislav Hudec, 2019 2nd International Conference on Computer Applications & Information Security (ICCAIS) “Feature extraction using Deep Learning for Intrusion Detection System.”
10. Aditya Phadke, Mohit Kulkarni, Pranav Bhawalkar, Rashmi Bhattad, 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC) “A Review of Machine Learning Methodologies for Network Intrusion Detection.”
11. Iftikhar Ahmad , Mohammad Basher, Muhammad Javed Iqbal, Aneel Rahim, IEEE Access (Volume: 6) **Page(s):** 33789 – 33795 “Performance Comparison of Support Vector Machine, Random Forest, and Extreme Learning Machine for Intrusion Detection.”
12. B. Riyaz, S. Ganapathy, 2018 International Conference on Recent Trends in Advanced Computing (ICRTAC) “An Intelligent Fuzzy Rule-based Feature Selection for Effective Intrusion Detection.”

