

Intrusion Detection Using a Hybrid Sequential Model



Abhishek Sinha, Aditya Pandey, and P. S. Aishwarya

1 Introduction

1.1 Context

A network intrusion could be any unauthorized activity on a computer network. Virus attacks, unauthorized access, theft of information, and denial-of-service attacks were the greatest contributors to computer crime. Detecting an intrusion depends on the defenders having a clear understanding of how attacks work. Detecting an intrusion is the first step to create any sort of counteractive security measure. Hence, it is very important to accurately determine whether a connection is an intrusion or not.

1.2 Categories

There are four broad categories of intrusions in a network of systems:

1.2.1 Dos

In computing, a denial-of-service attack is a cyber-attack in which the perpetrator seeks to make a machine or network resource unavailable to its intended users by

A. Sinha · A. Pandey · P. S. Aishwarya (✉)
PES University, 100 Feet Ring Road, Banashankari III Stage, Bangalore 560085, India
e-mail: aishwarya.ps31@gmail.com

A. Sinha
e-mail: sinha.abhish3k@gmail.com

A. Pandey
e-mail: pandeyan98@gmail.com

© Springer Nature Singapore Pte Ltd. 2021
S. Patnaik et al. (eds.), *Advances in Machine Learning and Computational Intelligence*,
Algorithms for Intelligent Systems, https://doi.org/10.1007/978-981-15-5243-4_1

temporarily or indefinitely disrupting services of a host connected to the Internet, e.g., back, land, and Neptune (Tables 1, 2, 3, 4).

1.2.2 U2r

These attacks are exploitations in which the hacker starts off on the system with a normal user account and attempts to abuse vulnerabilities in the system in order to gain superuser privileges, e.g., perl, xterm.

1.2.3 R2l

Remote to local is an unauthorized access from a remote machine, by maybe guessing the password of the local system and accessing files within the local system, e.g., ftp write, guess passwd, imap.

1.2.4 Probe

Probing is an attack in which the hacker scans a machine or a networking device in order to determine weaknesses or vulnerabilities that may later be exploited so as to compromise the system. This technique is commonly used in data mining, e.g., saint, portsweep, mscan, and nmap.

1.3 Importance of Intrusion Detection

Intrusion detection is important for both military as well as commercial sectors for the sake of their information security, which is the most important topic of research for future networks. It is critical to maintaining a high level of security to ensure safe and trusted communication of information between various organizations.

Intrusion detection system is a new safeguard technology for system security after traditional technologies, such as firewalls and message encryption. An intrusion detection system (IDS) is a device or software application that monitors network system activities for malicious activities or policy violations and produces reports to a management station.

2 Previous Work

The two approaches to an intrusion detection system are misuse detection and anomaly detection. A key advantage of misuse detection techniques is its high degree

of accuracy in detecting known attacks and their variations. Their obvious drawback is the inability to detect attacks whose instances have not yet been observed. Anomaly detection schemes, on the other hand, suffer from a high rate of false alarms. This occurs primarily because previously unseen (yet legitimate) system behaviors are also recognized as anomalies and are hence flagged as potential intrusions. Researchers have used various algorithms to solve this specific problem.

A data mining algorithm such as random forest can be applied for misuse, anomaly, and hybrid network-based intrusion detection systems. [7] uses a hybrid detection technique where they employ misuse detection followed by anomaly detection. This approach, however, has multiple limitations:

- Intrusions need to be much lesser than the normal data. Outlier detection will only work when the majority of the data is normal.
- A novel intrusion producing a large number of connections that are not filtered out by the misuse detection could decrease the performance of the anomaly detection or even the hybrid system as a whole.
- Some of the intrusions with a high degree of similarity cannot be detected by this anomaly detection approach.

The paper ‘Artificial Neural Networks for Misuse Detection’ [8] talks about the advantages of using an artificial neural network approach over a rule-based expert system approach. It also tells us about the various approaches with which these neural networks are applied to get high accuracy. A review paper on misuse detections [9] sheds light on the most common techniques to implement misuse detection. These are:

- Expert systems, which code knowledge about attacks as ‘if-then’ implication rules.
- Model-based reasoning systems, which combine models of misuse with evidential reasoning to support conclusions about the occurrence of a misuse.
- State transition analysis, which represents attacks as a sequence of state transitions of the monitored system.
- Keystroke Monitoring, which uses user keystrokes to determine the occurrence of an attack.

Landge and Wadh introduces a pattern matching model along with an artificial neural network model for their implementation of a misuse detection system. Further [9] performed intrusion detection using data mining techniques along with fuzzy logic and basic genetic algorithms. [1–5] are varied approaches, but their accuracy is not as good as the approach in [7]. The approach specified in [6] has good accuracy and also supports unsupervised learning algorithms, but it requires a very complicated state machine specification model.

3 Problem Statement

Here, our goal is to detect network intrusions and to create a predictive model that is able to differentiate between ‘good’ or ‘bad’ connections (intrusions) and classify those intrusions into known categories. To that end, we used the KDD cup dataset from 1999 which was created by simulating various intrusions over a network over the course of several weeks in a military environment.

The dataset itself consists of 42 attributes that are used in the various models according to relevance and importance. It contains approximately 1,500,000 data points. For the purposes of this project, we use a 10% representative dataset, since our machines did not have the processing power to handle the larger dataset. The dataset, as is, is unbalanced across the various result categories, and hence, we balance it by upsampling and downsampling. The dataset is also cleaned. Certain categorical variables like the ‘protocol type’ column are one-hot encoded. We found out that the ‘flag’ and ‘services’ attributes are not of much value as they are nominal variables, and hence, they are dropped.

4 Approach

Due to the various drawbacks of the individual anomaly detection models as well as the individual misuse detection models, we use a combined approach, i.e., a hybrid of the two. We combine the models serially such that the anomaly detection is followed by the misuse detection. This approach provides us with multiple advantages:

- The misuse detection acts as a verification model where it verifies whether an anomaly is actually an intrusion or not. This helps us reduce the false positives coming from the anomaly detection (primary drawback).
- The misuse detection also helps us classify the intrusions detected into various categories based on the intrusions ‘signature’ (centroid of sample data from the training dataset).

Figure 1 shows the approach for the intrusion detection system.

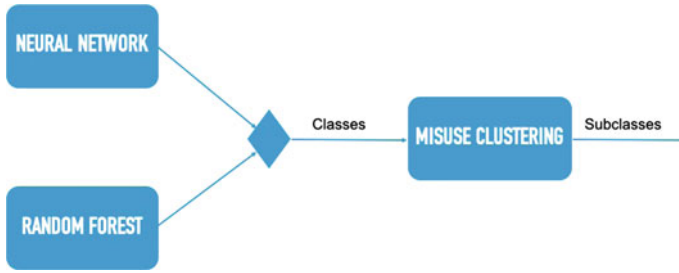


Fig. 1 Block diagram of approach

Table 1 Number of data points

Label:	Dos	Normal	Probe	Rtl	U2r
Before sampling	54,572	87,832	2131	999	52
After sampling	27,285	39,524	2131	999	86

4.1 Modeling

For anomaly detection, we use two models:

- Random forest model.
- Neural network model.

After running both these classification models parallelly, we take the combined output of the anomalies detected by either of the two. Doing so helps us reduce the rate of false negatives. The connections which are identified as anomalies by either of the two are then passed on to the misuse detection model. A false positive from the anomaly detection model classified as ‘normal’ in the misuse detection model reduces the number of false positives.

For misuse detection, we use a K-means-based clustering model. Therefore, the use of the misuse detection model on the anomalies helps trim down the number of false positives. The clustering model can also be used to further classify the 5 classes of attacks into 24 subclasses which will help in fighting or preventing the occurrence of the attack.

5 Components of the Intrusion Detection System

Our intrusion detection system consists of 3 components which have been combined to create one complete robust system. The components of the anomaly detection system are as follows.

5.1 *Neural Network*

A neural network is a system of hardware and/or software patterned after the operations of neurons in the human brain.

For the neural network, we created a sequential neural network with two hidden layers. The neural network consists of 41 input nodes and 5 output nodes. The 41 input nodes are the numerical attributes that are obtained after preprocessing, and the 5 output nodes are used to classify the data point into one of the 5 classes of connections that we have (u2r, rtl, normal, dos, probe).

The model is validated by a K-fold cross-validation with k as 2, and an average accuracy of 99.57% was obtained. In an attempt to reduce processing time without much drop in accuracy, we reduced the number of splits in the k-fold cross-validation as well as used a lower number of hidden layers in the network with minimal loss of accuracy.

5.2 *Random Forest*

A random forest is essentially a multitude of decision trees. The output obtained from a random forest model is a combination of the outputs obtained from all the decision trees.

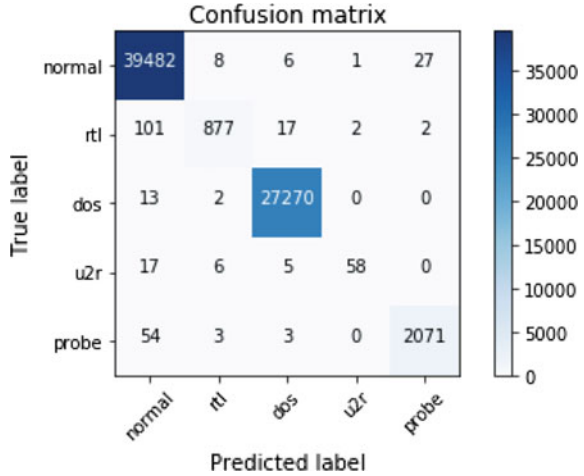
In our case, we use the mode of all the decision tree outputs as the final output from the model. Based on trial and error, the number of decision trees in our random forest was taken as 100. We use this random forest classifier to classify the connection among the 5 main classes.

Based on the importance value of variables, all the non-essential variables are dropped before the model is retrained. This helps in removing all the redundant and useless attributes in the model. Using the random forest model, an average accuracy of 99.78% is obtained. In order to reduce the processing time, the number of decision trees was reduced from an initial value of 1000–100 without resulting in any significant loss in accuracy of the model. Figure 2 shows a part of the large random forest model created (Fig. 2).

5.3 *Misuse Clustering*

We used a K-means-based clustering algorithm as the misuse detection model. This clustering algorithm is used to further classify the five classes of connections into 24 classes. During training, clusters are made for each known class and their centroids are stored. When the model receives a data point during testing, the minimum distance to any of the centroids is calculated to assign the point to the relevant class. If a

Fig. 2 Confusion matrix for 5 classes



connection arrives during testing that is falsely tagged as an attack, then the cluster it is assigned to should theoretically be ‘normal,’ thereby reducing false positives.

5.4 Combination of Models

The output from both the models in the anomaly detection stage is verified with each other, and the connections labeled as some kind of attack or connections not having the same labels are passed on to the misuse clustering model which will further classify the attack into its subclass. Another version of our final model is to use the misuse clustering model to reduce the number of false positives, thereby increasing precision and accuracy.

6 Results

6.1 Discussion of Results

Since we are using a synthetic dataset, we obtain accuracies that are quite high. However, the precision and recall for u2r in specific, in both the models are very low compared to the rest of our numbers. This can be attributed to the fact that the number of connections causing u2r attacks itself is very low, making up only around 70 of the 70,000 connections.

Table 2 Neural network results

Label:	Normal	Rtl	Dos	U2r	Probe
Precision	98.391	97.560	99.955	79.710	97.468
Recall	99.817	80.080	99.146	63.953	90.333
Accuracy	98.976	99.687	99.650	99.935	99.634

Table 3 Random forest results

Label:	Normal	Rtl	Dos	U2r	Probe
Precision	99.820	99.382	99.985	87.500	99.669
Recall	99.949	96.596	99.978	81.395	98.967
Accuracy	99.870	99.942	99.985	99.962	99.958

Table 4 Misuse clustering results

Type of classification:	5 class	24 class
Accuracy	99.651	91.31

As a result, the models have not been trained enough to be able to detect these attacks. Hence, a good proportion of u2r attacks have been misclassified. This problem can be rectified if a more balanced dataset is provided.

From the tabular columns, we can see that most of the attacks have not only been detected but also correctly classified in the anomaly detection model. In an attempt to further classify the attacks into its subcategories, we lose a bit of accuracy but it is still at an acceptable 91.3%. A further classification is necessary from the perspective of dealing with these intrusions. For example, A and B might both belong to ‘dos’ but they might have different security requirements to stop these attacks. We can also use the misuse clustering to cut down on false positives and that will only make our accuracy and precision better if in case the drop in accuracy is not preferable.

From the confusion matrix, we can also see that the number of ‘normal’ and ‘dos’ connections are significantly higher than the number of all the other connections even after downsampling these categories. Therefore, in order to achieve a balance that is acceptable, we downsampled these 2 categories and upsampled the ‘u2r’ category.

What can also be seen is that the number of misclassifications is fairly less. The only issue which can be seen is that a fraction of the misclassification is classifying one of the kinds of attacks as normal connections. The most probable reason for this is due to the fact that the number of ‘normal’ connections make up more than 55% of the total dataset. However, considering the number of false negatives versus the total number of connections, this error is almost insignificant.

Another important aspect of an intrusion detection system is its applicability to be used in real-time systems. We believe that our solution is easily deployable as once our model is trained (on previous data), and the classification happens almost instantaneously. The clustering algorithm is also able to identify new anomalies by

comparing the incoming traffic to the cluster centroids that are already there using a cutoff to determine the novelty of that particular packet.

7 Conclusion

Even though technology has advanced leaps and bounds over the past few decades, intrusion detection continues to be an active research field. This paper outlines our approach to solving this problem by using a combination of anomaly and misuse detection models. The techniques used to perform the same have been explained and illustrated. We believe that this approach is quite successful in classifying connections into their respective categories.

8 Future Work

The current approach is able to spot and classify new intrusions as such. However, the number of false positives for this is high. Furthermore, the connections detected as new intrusions could be further categorized into subclasses based on similarity (clustering) and dealt with accordingly. Different approaches could also be tried in order to obtain more accurate signatures for the different categories of intrusions which would improve the accuracy of the misuse detection model.

References

1. A. Lazarevic, L. Ertoz, V. Kumar, A comparative study of anomaly detection schemes in network intrusion detection. in *Proceedings of the 2003 SIAM International Conference on Data Mining* (2002)
2. D. Barbara, N. Wu, S. Jajodia, Detecting novel network intrusions using bayes estimators. in *Proceedings of the 2001 SIAM International Conference on Data Mining* (2001)
3. S.A. Hofmeyr, S. Forrest, A. Somayaji, Intrusion detection using sequences of system calls. *J. Comput. Security.* 6(3), 151–180 (1998)
4. A. Ghosh, A. Schwartzbard, A study in using neural networks for anomaly and misuse detection. in *Proceedings of the 8th USENIX Security Symposium*, August 23–36,(1999), pp. 141–152
5. E. Eskin, W. Lee, S.J. Stolfo, Modeling system calls for intrusion detection with dynamic window sizes. in *Proceedings DARPA Information Survivability Conference and Exposition II*, (DISCEX'01, 2001)
6. R. Sekar, A. Gupta, J. Frullo. Specification-based anomaly detection: A new approach for detecting network intrusions. in *CCS '02: Proceedings of the 9th ACM conference on Computer and communications security*, November 2002, (2002), pp. 265–274
7. J. Zhang, M. Zulkernine, A. Haque, Random-forests-based network intrusion detection systems. in *IEEE Transactions on Systems, Man, and Cybernetics—part C: Applications and Reviews*, vol. 38, No. 5 (2008)

8. J. Cannady, Artificial neural networks for misuse detection. in *National Information Systems Security Conference*, (1998)
9. R.S. Landge, A.P. Wadh, Misuse detection system using various techniques: A review. *Int. J. Adv. Res. Comput. Sci.*, Udaipur **4**(6) (2013)