

Machine Learning based IDS for Software Defined Networking

Oqbah Ghassan Abbas, Khaldoun Khorzom and Mohammed Assora
Department of Communications Engineering, HIAST, Damascus, Syria

Key words: SDN, IDS, machine learning, NSL-KDD, technique

Corresponding Author:

Oqbah Ghassan Abbas
*Department of Communications Engineering, HIAST,
Damascus, Syria*

Page No.: 3354-3358

Volume: 15, Issue 18, 2020

ISSN: 1816-949x

Journal of Engineering and Applied Sciences

Copy Right: Medwell Publications

Abstract: In the last few years, big companies have been depending more and more on Software defined Networking “SDN” to fulfill their needs for programmable networks. But like other networks, SDN has some security issues. Many technologies are used to solve such problems and machine learning is considered one of the best. Machine learning has demonstrated its ability to find data patterns when other technologies failed. This makes it a perfect choice for intrusion detection system “IDS” in general and anomaly-based detection in particular. In this research, we propose a new anomaly-based IDS that benefits from the ability of SDN to provide statistical features about flows that pass through the network and passes these features to a voting system that consists of several machine learning algorithms. This technique gives the system the ability to study the user’s behavior and predict any possible intrusion. The voting system is trained and tested using NSL-KDD and KDDCup99 datasets and the results shows increasing in detection accuracy and decreasing in false positive rate.

INTRODUCTION

The need for programmable networks has drawn the attention of data-centres and big-data companies to new paradigms like Software Defined Networking (SDN) which becomes a trend in the last few years. SDN’s main goal is to separate the control and the data planes. The controller in the control plane is able to generate and modify the flow tables to suit the network services which are running like applications within the controller. But like every new paradigm, SDN faces many challenges, especially with security^[1]. Intrusion Detection Systems (IDS) are considered as one of the best solutions for network security as it’s able to predict and alarm the network administrator about possible intrusions. In the last decade, IDSs are more interested in studying the

user’s behavior to predict intrusions. Using new technologies like machine learning which is perfect for such problems as it is able to extract new information with every interaction between the users and the network. Furthermore, it is able to find patterns in this information which helps studying the user’s behavior^[2]. The problem in using machine learning with network security is the need of data extraction from the network which increases the overhead in it. But after SDN is introduced, it is possible to benefit from its properties as it provides statistical features about the traffic that passes the network^[3]. In addition, to the possibility to add the IDS as a program within the SDN controller^[4] which makes it possible to implement an IDS without any additional cost and without increasing the network latency or decreasing its bandwidth.

Background: In this study, we briefly present SDN, IDS and machine learning.

MATERIALS AND METHODS

SDN: Software defined networking is a new paradigm in networking that was developed around 2008^[5] where SDN-Controller and SDN-Switches are the main devices in the network, instead of routers and switches. The key concept of SDN is the separation of the control plane and the data plane as the controller generates the flow tables and distributes them to the switches which are only responsible for forwarding the data packets. In addition, the switches calculate statistical information about flows and forward this information to the controller. Table 1 presents this information^[6].

The main benefit of SDN is the ability to add services to the network as applications that run within the controller without adding any devices or modifying the network’s topology^[4]. This opens a new horizon for improving networks with minimum cost.

IDS: The intrusion detection is the process of monitoring and analyzing any network event to detect any possible intrusion. It can be classified as Host-based when it runs on a network host or Network-Based when it monitors network packets^[7]. For detection strategy, IDS can be classified as signature-based detection where it looks for patterns that matches known attack pattern and anomaly-based detection where it looks for patterns that does not match normal behavior^[7]. Many metrics are used to evaluate the performance of an IDS but we will focus on detection accuracy and false positive rate. These two metrics can be calculated from the confusion matrix. Accuracy and false positive rate are given as follows:

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} \tag{1}$$

$$\text{FPR} = \frac{FP}{FP+TN} \tag{2}$$

Where:

TP = The number of attack samples that were classified correctly

TN = The number of normal samples that were classified correctly

FP = The number of normal samples that were classified incorrectly

FN = The number of attack samples that were classified incorrectly

Machine learning: Machine learning is a subfield of computer science that evolves from the study of pattern recognition and computational learning theory in artificial intelligence. It explores the construction and study of

Table 1: Features description

Feature name	Description
Duration	Length (in seconds) of the connection
Protocol_type	Type of protocol, e.g. tcp, udp, etc.
Src_bytes	Number of data bytes from source to destination
Dst_bytes	Number of data bytes from destination to source
Count	Number of connections to the same host as the current connection in the past two seconds
srv_count	Number of connections to the same service as the current connection in the past two seconds

algorithms that can learn from and make predictions on data. Machine learning is classified as supervised learning when training data is labeled, unsupervised learning when training data is unlabeled and semi-supervised when training data is a mixture between labeled and unlabeled data. Many machine learning algorithms have been developed and improved in the last two decades, from these algorithms, we will use the following:

Decision tree: The algorithm chooses the feature with the highest information gain to be the root node, then the ‘gini index’ is calculated to find the best partition, then the process is repeated till reaching the specified maximum depth.

Random forest: A number of decision trees are built depending on a different subset of the dataset for each of them, then the performance of all the trees is averaged to get the final result of the algorithm.

XGBoost: A decision tree is built by using a subset of the data set to get a level 1 decision tree, then other subsets will be used consecutively till reaching a specified level. The algorithm could start with many initial decision trees and choose the best one after a certain level.

Support vector machine: Depends on finding the best geometric separator between the classes of the dataset by finding the distances between the points of the dataset. The kernel trick could be used by using kernel functions that measure the distances between the dataset points after projecting them to another dimension.

Deep neural network: Similar to simple neural networks but has multiple hidden layers.

Previous work: Vigneswaran and Poornachandran^[8], introduce an anomaly-based IDS that works in traditional networks and depends on deep neural network model. The proposed solution gives an accuracy of 93%. They use KDDCup99 dataset which suffers from imbalance classes and redundant records which affects the reliability of the results. Ajaeiya *et al.*^[9], suggest an anomaly-based IDS that works in SDN and only uses the features provided by it. They compare the results of multiple machine learning algorithms. Random

Forest algorithm gives the best results where the true positive rate is 96.3% and the false positive rate is 0.009. The results show the efficiency of depending on the probabilistic distribution using algorithms like Random Forest. However, in their research, they do not use a standard dataset which raises some concerns about the validity of their results. Abubakar and Pranggono^[10], propose an IDS that works in SDN and consists of a signature-based IDS and an anomaly-based IDS. The anomaly-based IDS depends on deep neural network model and uses NSL-KDD dataset for training and testing. The detection accuracy is 97.4%. However, the intrusions that are detected by the signature-based part are not separated from intrusions detected by the anomaly-based part. So, the accuracy of the anomaly-based part cannot be evaluated. Tang and Mhamdi^[6], suggest an anomaly-based IDS that works in SDN and only uses the features provided by it. They compare the results of multiple machine learning algorithms. Using a deep neural network with three hidden layers gives the best results with an accuracy of 75.75%. However, the parameters of the used algorithms are not provided. Therefore, it is hard to propose a better tuning for the algorithms.

The proposed solution: In this study, we present the architecture of the proposed IDS and the datasets used to train and test the machine learning part.

Architecture: In this study, we focus on the design of the IDS. As the input of the IDS is not the whole packets, we need the first part to extract the features. These features can be extracted via the SDN-Switches. As the extracted features do not need preprocessing, we can directly pass them to the second part. The second part is a machine learning-based voting system that will produce a prediction on whether the flow is normal or abnormal.

This part consists of several machine learning algorithms that are trained using KDDCup99 and NSL-KDD datasets. Support vector machine algorithm is chosen to cover the possibility of the dataset being geometrically separable. Decision tree, random forest and XGBoost are chosen in case, we can obtain better predictions depending on the probabilistic distribution of the features. Deep neural network is also chosen because it proves its dominance in finding patterns in datasets. The third part is concern in adding rules that prevent malicious flows from passing the network. Figure 1 shows the architecture of the proposed solution.

Implementation: As we have seen, the proposed solution consists of three parts:

Extracting features: In this part, we emulate an SDN network using GNS3. The data plane consists of two hosts connected to an Openvswitch and OpenDayLight controller that runs on Ubuntu 16.04 Linux distribution and it is connected to the Openvswitch through openflow1.3 protocol. Every two seconds the controller requests the features collected by the switch using a feature-req message (openflow message) and the switch replies with a feature-reply message containing the features presented in Table 1.

The voting system: This part receives the extracted features and runs them through a set of previously trained machine learning algorithms to calculate the probability of these features being associated to a possible intrusion. The system uses the machine learning algorithms that collaborate to calculate the final decision. Every algorithm predicts new samples individually and the final prediction is made using the equation:

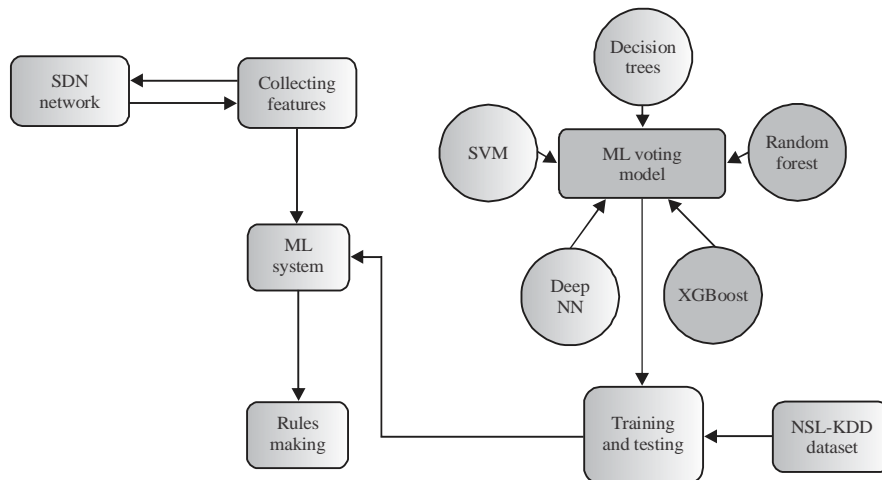


Fig. 1: The architecture of the proposed solution

$$p = \frac{\sum p_i * acc_i}{\sum acc_i} \quad (3)$$

Where:

P = The probability of a test sample being associated to an intrusion

p_i = The probability of a test sample being associated to an intrusion according to the model (i)

acc_i = The accuracy of the model (i) during the testing phase

By using Eq. 3 to calculate the final prediction, we have a voting system that gives higher priority to the algorithms that performed better during the testing phase.

Adding rules: Every time a flow is marked as an intrusion, a new rule is added to the Openvswitch that prevents this flow from passing through the network. It should be noted that adding this part to the system turns it into an Intrusion Prevention System (IPS) instead of an IDS.

RESULTS AND DISCUSSION

The experimental results: We implement a machine learning based NIDS in SDN for 2-class classification (normal and abnormal) with the aim of increasing the accuracy and decreasing the FPR. All mentioned algorithms are built, trained and tested using Keras and SKLearn libraries within PyCharm IDE. Table 2 shows the parameters used to implement each algorithm.

At first, we use KDDcup99 dataset to test these algorithms separately and the results are presented in Table 3. From the results in Table 3 we notice.

Decision tree gives better results than random forest and XGBoost which are considered to be improvements to the decision tree algorithm. We justify that by the number of the features we use. KDDcup99 provide 41 features, we only use the six features that are compatible with SDN. XGBoost and random forest use subsets of the data sets for each sub-tree and each of these subsets have a maximum of three features (best case scenario), making the sub-trees unable to learn enough from these subsets.

SVM fails to give good results compared to other algorithms which means our dataset isn't separable geometrically. DNN gives the best accuracy as this kind of algorithms proved its ability to find patterns other algorithms fail to find. Comparing with^[7,8], we have better accuracy and better FPR.

Because KDDCup99 dataset has several drawbacks, we also use NSL-KDD dataset to evaluate the machine learning algorithms and the results are presented in Table 4 from the results in Table 4, we notice.

Table 2: Parameters values for machine learning algorithms

Algorithm	Parameter: value	Values
Decision tree	Max_depth	2
Random forest	N_estimators	100
	Max_features	Sqrt
XGBoost	Objective	Binary: hinge
	Gamma	1.0
	Learning rate	0.06
	Colsample_by tree	0.3
	Max_depth	3
	N_estimators	300
	Num_round	500
SVM	Kernel	Rbf
	C: 100	100
DNN	Hidden layers	3
	Activation function	ReLU for hidden layers, sigmoid for output layer
	Loss function	binary-crossentropy
	Learning rate	1e-7
	Optimizer	Adam

Table 3: Results of the machine learning algorithms with KDDCup99

Algorithm	Accuracy (%)	FPR
Decision tree	99.40	0.0009
Random forest	98.54	0.05
XGBoost	99.05	0.01
SVM	89.68	0.012
DNN	99.50	0.0009

Table 4: Results of the machine learning algorithms with NSL-KDD

Algorithm	Accuracy (%)	FPR
Decision tree	82.72	0.05
Random forest	77.40	0.03
XGBoost	79.61	0.03
SVM	73.22	0.038
DNN	83.80	0.07

Table 5: Results of the voting system

Evaluation metric	Results
Accuracy	79.6%
FPR	0.03

The resulting pattern of Table 3 is repeated in Table 4 where DNN gives the best accuracy, followed by Decision Tree, Random Forest and XGBoost and finally SVM gives the lowest accuracy. Comparing with Tang *et al.*^[6], we have better accuracy but comparing with Abubakar *et al.*^[10] their solution gives better accuracy as their solution uses a signature-based IDS in addition to the anomaly-based IDS.

KDDCup99 gives better results than NSL-KDD, this is because the redundant samples in the KDDCup99 and the imbalance of the number of samples for each class in KDDCup99 which makes NSL-KDD more reliable.

Depending on the results of Table 3 and 4 and the fact that NSL-KDD is more reliable than KDDCup99, the voting system is implemented using the models resulted from using NSL-KDD dataset. By applying Eq. 3 to the results shown in Table 4, we get the following results:

From the results in Table 5, we notice that the voting system averages the accuracy of the used algorithms and gives the best FPR comparing to the individual results of

the used algorithms. We can justify that as the algorithms with lower accuracy come together to change the vote of the algorithms with the higher accuracy to improve the overall FPR but this happens at the expense of the overall accuracy. We can notice that; it is always possible to add or remove any algorithm to the voting system in order to trade-off between the accuracy and FPR.

CONCLUSION

In this study, we proposed a machine learning based NIDS for software defined networks. A voting system is implemented using several machine learning algorithms. This system receives the features provided by SDN, so, no more equipment or bandwidth consuming is added to the network. We test our system by using NSL-KDD dataset because it is more reliable than KDDcup99. Finally, the results show an increasing in the overall accuracy to 79.6% and decreasing in FPR to 0.03. The system maintains the ability to trade-off one at the expense of the other, by adding or removing algorithms from the voting system.

REFERENCES

01. Bindra, N. and M. Sood, 2016. Is SDN the real solution to security threats in networks? A security update on various SDN models. *Indian J. Sci. Technol.*, 9: 1-8.
02. Ushmani, A., 2014. Machine learning pattern matching. *Int. J. Comput. Sci. Trends Technol.*, 7: 4-7.
03. Da Silva, A.S., C.C. Machado, R.V. Bisol, L.Z. Granville and A. Schaeffer-Filho, 2015. Identification and selection of flow features for accurate traffic classification in SDN. *Proceedings of the 2015 IEEE 14th International Symposium on Network Computing and Applications*, September 28-30, 2015, IEEE, Cambridge, Massachusetts, pp: 134-141.
04. Hoang, D.B. and M. Pham, 2015. On software-defined networking and the design of SDN controllers. *Proceedings of the 2015 6th International Conference on the Network of the Future (NOF)*, September 30-October 2, 2015, IEEE, Montreal, Canada, pp: 1-3.
05. Hartpence, B. and R. Rosario, 2016. Software defined networking for systems and network administration programs. *USENIX J. Edu. Syst. Administration*, 2: 12-42.
06. Tang, T.A., L. Mhamdi, D. McLernon, S.A.R. Zaidi and M. Ghogho, 2016. Deep learning approach for network intrusion detection in software defined networking. *Proceedings of the 2016 International Conference on Wireless Networks and Mobile Communications (WINCOM)*, October 26-29, 2016, IEEE, Fez, Morocco, pp: 258-263.
07. Sultana, N., N. Chilamkurti, W. Peng and R. Alhadad, 2019. Survey on SDN based network intrusion detection system using machine learning approaches. *Peer-to-Peer Networking Appl.*, 12: 493-501.
08. Vigneswaran, K.R., R. Vinayakumar, K.P. Soman and P. Poornachandran, 2018. Evaluating shallow and deep neural networks for network intrusion detection systems in cyber security. *Proceedings of the 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, July 10-12, 2018, IEEE, Bangalore, India, pp: 1-6.
09. Ajaeiya, G.A., N. Adalian, I.H. Elhadj, A. Kayssi and A. Chehab, 2017. Flow-based intrusion detection system for SDN. *Proceedings of the 2017 IEEE Symposium on Computers and Communications (ISCC)*, July 3-6, 2017, IEEE, Heraklion, Greece, pp: 787-793.
10. Abubakar, A. and B. Pranggono, 2017. Machine learning based intrusion detection system for software defined networks. *Proceedings of the 2017 7th International Conference on Emerging Security Technologies (EST)*, September 6-8, 2017, IEEE, Canterbury, England, UK., pp: 138-143.