



ARTIFICIAL NEURAL NETWORK APPROACH TO CYBER THREAT DETECTION USING EVENT PROFILES

D. Siva Ranjan Das¹, B. Bhagyalaxmi², B. Bhargavi³

¹ Assoc. Professor, Department of Computer Applications, Aurora's PG College (MBA), Uppal, Hyderabad
Email: sivaranjandas@yahoo.com

² Assistant Professor, Department of Computer Applications, Aurora's PG College (MBA), Uppal, Hyderabad
Email: lxmibbaghya@gmail.com

³ Assistant Professor, Department of Computer Applications, Aurora's PG College (MBA), Uppal, Hyderabad
Email: burrabhargavi@gmail.com

ABSTRACT

Developing an automated method for detecting cyber threats is one of the main issues facing cyber security. In this paper, we describe an artificial neural network-based cyber threat detection method. The suggested solution improves cyber-threat identification by converting a large number of gathered security events into unique event profiles and utilizing a deep learning-based detection algorithm. In this study, we created an AI-SIEM system using a variety of artificial neural network techniques, such as CNN, LSTM, and FCNN, together with event profiling for data preparation. The approach helps security analysts react quickly to cyber threats by focusing on differentiating between true positive and false positive signals. The authors of this work conducted all of the experiments using two real-world datasets and two benchmark datasets, NSLKDD and CICIDS2017. We ran trials utilizing the five traditional machine-learning techniques (SVM, k-NN, RF, NB, and DT) to assess the performance comparison with current methodologies. As a result, the study's experimental findings confirm that our suggested approaches can be used as learning-based models for network intrusion detection and demonstrate that, despite being used in real-world scenarios, their performance surpasses that of traditional machine learning techniques.

Keywords: Artificial Neural Network, Cyber Threat Detection, FCNN, CNN, and LSTM

1. INTRODUCTION

Learning-based methods for identifying cyber attacks have advanced further with the advent of artificial intelligence (AI) technology, and numerous studies have found significant outcomes with them. However, it is still very difficult to defend IT systems against threats and bad activities in networks since cyber attacks are always changing. Effective defenses and security concerns were given top priority for locating trustworthy solutions due to numerous network intrusions and harmful actions [1, 2, 3, 4, 5].

In the past, there have been two main methods for identifying network breaches and



cyber threats. Within the company network, an intrusion prevention system (IPS) is implemented. Its primary way of examining network protocols and flows is signature-based. It creates relevant intrusion alarms, also known as security events, and notifies another system—like SIEM—of the alerts it generates. The collection and handling of IPS alerts has been the primary emphasis of security information and event management, or SIEM. Among the different security operations systems available for analyzing the gathered security events and logs, the SIEM is the most popular and reliable option [5]. Furthermore, security analysts try to look into suspicious alerts based on policies and thresholds, as well as find malicious activity by utilizing attack-related knowledge to analyze correlations between events and find patterns of behavior.

However, because of their high false alarm rate and the volume of security data they include, it is still challenging to identify and detect intrusions against intelligent network attacks [6], [7]. For this reason, machine learning and artificial intelligence algorithms for attack detection have received more attention in the most recent studies in the field of intrusion detection. Security analysts can investigate network attacks more quickly and automatically with the help of advancements in AI fields. These learning-based techniques necessitate using trained models to identify intrusions for unknown cyber threats, after learning the attack model from previous threat data [8], [9].

For analysts who need to quickly examine a huge number of events, a learning-based approach designed to determine whether an attack occurred in a big amount of data can be helpful. Information security solutions can be broadly classified into two types, according to [10]: machine learning-driven solutions and analyst-driven solutions. Analyst-driven solutions are based on rules that are established by analysts, who are security professionals. In the meantime, emerging cyberthreat detection can be enhanced by machine learning-driven solutions that identify uncommon or aberrant patterns [10]. However, we found that the current learning-based approaches have four major drawbacks, despite the fact that learning-based approaches are helpful in identifying cyber attacks in systems and networks. Initially, labeled data are needed for learning-based detection techniques in order to train the model and assess the produced learning models. Moreover, obtaining such labeled data at a scale that permits precise model training is not simple. Many commercial SIEM solutions lack labeled data that can be used with supervised learning models, even though labeled data is necessary [10].

Second, because they are absent from popular network security systems, the majority of the learning characteristics that are theoretically employed in each study are not generalized features in the real world [3]. As such, it is challenging to apply to real-world scenarios. Deep learning technologies have been used in recent intrusion detection research efforts, and performance has been assessed using popular datasets such as NSLKDD [11], CICIDS2017 [12], and Kyoto-Honeypot [13]. Unfortunately, because to a lack of features, many earlier research that used benchmark datasets that were correct but could not be generalized to the real world. An implemented learning model must be evaluated using datasets gathered in the real world in order to get over these restrictions.

Third, although it may result in a high false alert rate, employing an anomaly-based approach to identify network intrusion can assist in identifying unidentified cyber threats [6]. It is very



expensive and time-consuming for staff to examine false positive alarms, which are frequently triggered.

Fourth, some hackers have the ability to gradually alter their behavior patterns in order to conceal their malicious operations [10], [14]. The detection models are not adequate because attackers continuously modify their behavior, even in cases when learning-based models are feasible. Furthermore, the analysis of transient network security events has been the primary focus of practically all security systems. We anticipate that, over extended periods of time, studying the security event history connected with the formation of events can be one way to identify the malicious conduct of cyber attacks and protect against them.

2. LITERATURE SURVEY AND RELATED WORK

1. Enhanced Network Anomaly Detection Based on Deep Neural Networks

Abstract: Due to the monumental growth of Internet applications in the last decade, the need for security of information network has increased manifolds. As a primary defense of network infrastructure, an intrusion detection system is expected to adapt to dynamically changing threat landscape. Many supervised and unsupervised techniques have been devised by researchers from the discipline of machine learning and data mining to achieve reliable detection of anomalies. Deep learning is an area of machine learning which applies neuron-like structure for learning tasks. Deep learning has profoundly changed the way we approach learning tasks by delivering monumental progress in different disciplines like speech processing, computer vision, and natural language processing to name a few. It is only relevant that this new technology must be investigated for information security applications. The aim of this paper is to investigate the suitability of deep learning approaches for anomaly-based intrusion detection system. For this research, we developed anomaly detection models based on different deep neural network structures, including convolutional neural networks, autoencoders, and recurrent neural networks. These deep models were trained on NSLKDD training data set and evaluated on both test data sets provided by NSLKDD, namely NSLKDDTest+ and NSLKDDTest21. All experiments in this paper are performed by authors on a GPU-based test bed. Conventional machine learning-based intrusion detection models were implemented using well-known classification techniques, including extreme learning machine, nearest neighbor, decision-tree, random-forest, support vector machine, naive-bays, and quadratic discriminant analysis. Both deep and conventional machine learning models were evaluated using well-known classification metrics, including receiver operating characteristics, area under curve, precision-recall curve, mean average precision and accuracy of classification. Experimental results of deep IDS models showed promising results for real-world application in anomaly detection systems.

2. Network Intrusion Detection Based on Directed Acyclic Graph and Belief Rule Base

Abstract: Intrusion detection is very important for network situation awareness. While a few methods have been proposed to detect network intrusion, they cannot directly and effectively utilize semi-quantitative information consisting of expert knowledge and quantitative data. Hence, this paper proposes a new detection model based on a directed acyclic graph (DAG) and a belief rule base (BRB). In the proposed model, called DAG-BRB, the DAG is employed to construct a multi-layered BRB model that can avoid explosion of combinations



of rule number because of a large number of types of intrusion. To obtain the optimal parameters of the DAG-BRB model, an improved constraint covariance matrix adaption evolution strategy (CMA-ES) is developed that can effectively solve the constraint problem in the BRB. A case study was used to test the efficiency of the proposed DAG-BRB. The results showed that compared with other detection models, the DAG-BRB model has a higher detection rate and can be used in real networks.

3. HAST-IDS: Learning hierarchical spatial-temporal features using deep neural networks to improve intrusion detection

Abstract: The development of an anomaly-based intrusion detection system (IDS) is a primary research direction in the field of intrusion detection. An IDS learns normal and anomalous behavior by analyzing network traffic and can detect unknown and new attacks. However, the performance of an IDS is highly dependent on feature design, and designing a feature set that can accurately characterize network traffic is still an ongoing research issue. Anomaly-based IDSs also have the problem of a high false alarm rate (FAR), which seriously restricts their practical applications. In this paper, we propose a novel IDS called the hierarchical spatial-temporal features-based intrusion detection system (HAST-IDS), which first learns the low-level spatial features of network traffic using deep convolutional neural networks (CNNs) and then learns high-level temporal features using long short-term memory networks. The entire process of feature learning is completed by the deep neural networks automatically; no feature engineering techniques are required. The automatically learned traffic features effectively reduce the FAR. The standard DARPA1998 and ISCX2012 data sets are used to evaluate the performance of the proposed system. The experimental results show that the HAST-IDS outperforms other published approaches in terms of accuracy, detection rate, and FAR, which successfully demonstrates its effectiveness in both feature learning and FAR reduction

4. Data security analysis for DDoS defense of cloud based networks

Abstract: Distributed computing has become an effective approach to enhance capabilities of an institution or organization and minimize requirements for additional resource. In this regard, the distributed computing helps in broadening institutes IT capabilities. One needs to note that distributed computing is now integral part of most expanding IT business sector. It is considered novel and efficient means for expanding business. As more organizations and individuals start to use the cloud to store their data and applications, significant concerns have developed to protect sensitive data from external and internal attacks over internet. Due to security concern many clients hesitate in relocating their sensitive data on the clouds, despite significant interest in cloud-based computing. Security is a significant issue, since data much of an organizations data provides a tempting target for hackers and those concerns will continue to diminish the development of distributed computing if not addressed. Therefore, this study presents a new test and insight into a honeypot. It is a device that can be classified into two types: handling and research honeypots. Handling honeypots are used to mitigate real life dangers. A research honeypot is utilized as an exploration instrument to study and distinguish the dangers on the internet. Therefore, the primary aim of this research project is to do an intensive network security analysis through a virtualized honeypot for cloud servers to tempt an attacker and provide a new means of monitoring their behavior

3. EXISTING SYSTEM



As there is no staff available in unmanned restaurants, it is difficult for the restaurant management to estimate how the concept and the food is experienced by the customers. Existing

rating systems, such as Google and TripAdvisor, only partially solve this problem, as they only cover a part of the customer's opinions. These rating systems are only used by a subset of the customers who rate the restaurant on independent rating platforms on their own initiative. This applies mainly to customers who experience their visit as very positive or negative.

4. PROPOSED SYSTEM

In order to solve the above problem, all customers must be motivated to give a rating. This paper

Introduces an approach for a restaurant rating system that asks every customer for a rating after their visit to increase the number of ratings as much as possible. This system can be used in unmanned restaurants; the scoring system is based on facial expression detection using pretrained convolution neural network (CNN) models. It allows the customer to rate the food by taking or capturing a picture of his face that reflects the corresponding feelings. Compared to

Text-based rating system, there is much less information and no individual experience reports collected. However, this simple fast and playful rating system should give a wider range of opinions about the experiences of the customers with the restaurant concept.

5. IMPLEMENTATION

MODULES:

upload Train Dataset
Run Preprocessing TF-IDF Algorithm
Generate Event Vector
Neural Network Profiling
Run SVM Algorithm
Run KNN Algorithm
Run Naive Bayes Algorithm
Run Decision Tree Algorithm
Accuracy Comparison Graph
Precision Comparison Graph
Recall Comparison Graph
F Measure Comparison Graph

MODULES DESCRIPTION:

Proposed algorithms consist of the following module

1. **Data Parsing:** This module takes input dataset and parses that dataset to create a raw data event model
2. **TF-IDF:** using this module we will convert raw data into event vector which will contain normal and attack signatures



3. **Event Profiling Stage:** Processed data will be splitted into train and test model based on profiling events.
4. **Deep Learning Neural Network Model:** This module runs CNN and LSTM algorithms on train and test data and then generate a training model. Generated trained model will be applied on test data to calculate prediction score, Recall, Precision and F Measure. Algorithm will learn perfectly will yield better accuracy result and that model will be selected to deploy on real system for attack detection
5. Datasets which we are using for testing are of huge size and while building model it's going to out of memory error but kdd_train.csv dataset working perfectly but to run all algorithms it will take 5 to 10 minutes. You can test remaining datasets also by reducing its size or running it on high configuration system.

5. RESULTS AND DISCUSSION SCREENSHOTS

To run project double click on 'run.bat' file to get below screen



In above screen click on 'Upload Train Dataset' button and upload dataset



In above screen we can see dataset contains 9999 records and now click on ‘Run Preprocessing TF-IDF Algorithm’ button to convert raw dataset into TF-IDF values

```
C:\Windows\system32\cmd.exe
x_test.shape before = (2000, 2978)
x_test.shape after = (2000, 2978)
y_test.shape = (2000, 13)
Model: "sequential_1"
Layer (type) Output Shape Param #
-----
dense_1 (Dense) (None, 32) 1056
dropout_1 (Dropout) (None, 32) 0
dense_2 (Dense) (None, 17) 561
Total params: 1,617
Trainable params: 1,617
Non-trainable params: 0
None
WARNING:tensorflow:From C:\Users\Admin\AppData\Local\Programs\Python\Python71\lib\site-packages\tensorflow\python\ops\tf_grad.py:1250: add_dispatch_support.<locals>.unwrap (from tensorflow.python.ops.array_ops) is deprecated and will be removed in a future version.
Instructions for updating:
Use tf.where in 2.0, which has the same broadcast rule as np.where
WARNING:tensorflow:From C:\Users\Admin\AppData\Local\Programs\Python\Python71\lib\site-packages\tensorflow\python\ops\tf_grad.py:422: The name tf.global_variables is deprecated. Please use tf.compat.v1.global_variables instead.
Epoch 1/1
192/7999 [.....] - FTB: 3.124 - loss: 0.2248 - accuracy: 0.8412
```

In above screen LSTM model is generated and its epoch running also started and its starting accuracy is 0.94. Running for entire dataset may take time so wait till LSTM and CNN training process completed. Here dataset contains 7999 records and LSTM will iterate all records to filter and build model.

```
Select C:\Windows\system32\cmd.exe
Instructions for updating:
Use tf.where in 2.0, which has the same broadcast rule as np.where
WARNING:tensorflow:From C:\Users\Admin\AppData\Local\Programs\Python\Python71\lib\site-packages\tensorflow\python\ops\tf_grad.py:422: The name tf.global_variables is deprecated. Please use tf.compat.v1.global_variables instead.
Epoch 1/1
8000/7999 [.....] - loss: 0.1448 - accuracy: 0.8412
=====
loss: 0.5412049
C:\Users\Admin\AppData\Local\Programs\Python\Python71\lib\site-packages\sklearn\metrics\_classification.py:1272: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
  zero_division, codify, msg_start, len(result))
Model: "sequential_2"
Layer (type) Output Shape Param #
-----
dense_3 (Dense) (None, 512) 152834
activation_1 (Activation) (None, 512) 0
dropout_2 (Dropout) (None, 512) 0
dense_4 (Dense) (None, 512) 162656
activation_2 (Activation) (None, 512) 0
dropout_3 (Dropout) (None, 512) 0
dense_5 (Dense) (None, 17) 8711
```



In above selected text we can see LSTM complete all iterations and in below lines we can see CNN model also starts execution

```
Epoch 10/100: 0.7200 accuracy: 0.7200 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 11/100: 0.7300 accuracy: 0.7300 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 12/100: 0.7400 accuracy: 0.7400 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 13/100: 0.7500 accuracy: 0.7500 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 14/100: 0.7600 accuracy: 0.7600 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 15/100: 0.7700 accuracy: 0.7700 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 16/100: 0.7800 accuracy: 0.7800 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 17/100: 0.7900 accuracy: 0.7900 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 18/100: 0.8000 accuracy: 0.8000 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 19/100: 0.8100 accuracy: 0.8100 val_loss: 0.0010 val_accuracy: 0.9900
Epoch 20/100: 0.8200 accuracy: 0.8200 val_loss: 0.0010 val_accuracy: 0.9900
```

In above screen CNN also starts first iteration with accuracy as 0.72 and after completing all iterations 10 we got filtered improved accuracy as 0.99 and multiply by 100 will give us 99% accuracy. So, CNN is giving better accuracy compare to LSTM and now see below GUI screen with all details



In above screen we can see Naïve Bayes algorithm output values and now click on ‘Run Decision Tree Algorithm’ to run Decision Tree Algorithm



Now click on ‘Accuracy Comparison Graph’ button to get accuracy of all algorithms

7. CONCLUSION AND FUTURE SCOPE

CONCLUSION

We have presented the AI-SIEM system in this study, which makes use of artificial neural networks and event profiles. Condensing extremely huge amounts of data into event profiles and utilizing deep learning-based detection techniques to improve cyber-threat detection capabilities are the innovative aspects of our work. By comparing long-term security data, the AI-SIEM system helps security analysts to respond quickly and effectively to important security alarms. It can also assist security analysts in quickly responding to cyber threats scattered throughout a multitude of security events by decreasing false positive warnings. We conducted a performance comparison utilizing two benchmark datasets (NSLKDD, CICIDS2017) and two real-world datasets to assess performance. Using well-known benchmark datasets, we first demonstrated how our techniques might be used as one of the learning-based models for network intrusion detection based on a comparison experiment with other approaches. Second, we demonstrated encouraging results from the evaluation using two real datasets, showing that our approach performed better in terms of accurate classifications than traditional machine learning techniques.

FUTURE SCOPE

In the future, to address the evolving problem of cyber attacks, we will focus on enhancing earlier threat predictions through the multiple deep learning approach to discovering the long-term patterns in history data. In addition, to improve the precision of labeled dataset for supervised-learning and construct good learning datasets, many SOC analysts will make efforts directly to record labels of raw security events one by one over several month

REFERENCES



- [1] S. Naseer, Y. Saleem, S. Khalid, M. K. Bashir, J. Han, M. M. Iqbal, K. Han, "Enhanced Network Anomaly Detection Based on Deep Neural Networks," *IEEE Access*, vol. 6, pp. 48231- 48246, 2018.
- [2] B. Zhang, G. Hu, Z. Zhou, Y. Zhang, P. Qiao, L. Chang, "Network Intrusion Detection Based on Directed Acyclic Graph and Belief Rule Base", *ETRI Journal*, vol. 39, no. 4, pp. 592-604, Aug. 2017
- [3] W. Wang, Y. Sheng and J. Wang, "HAST-IDS: Learning hierarchical spatial-temporal features using deep neural networks to improve intrusion detection," *IEEE Access*, vol. 6, no. 99, pp. 1792-1806,2018.
- [4] M. K. Hussein, N. Bin Zainal and A. N. Jaber, "Data security analysis for DDoS defense of cloud based networks," 2015 IEEE Student Conference on Research and Development (SCORED), Kuala Lumpur, 2015, pp. 305-310.
- [5] S. Sandeep Sekharan, K. Kandasamy, "Profiling SIEM tools and correlation engines for security analytics," In Proc. Int. Conf. Wireless Com., Signal Proce. and Net.(WiSPNET), 2017, pp. 717-721.
- [6] N. Hubballi and V. Surya Narayanan "False alarm minimization techniques in signature-based intrusion detection systems: A survey," *Comput. Commun.*, vol. 49, pp. 1-17, Aug. 2014.
- [7] A. Naser, M. A. Majid, M. F. Zolkipli and S. Anwar, "Trusting cloud computing for personal files," 2014 International Conference on Information and Communication Technology Convergence (ICTC), Busan, 2014, pp. 488-489.
- [8] Y. Shen, E. Mariconti, P. Vervier, and Gianluca Stringhini, "Tiresias: Predicting Security Events Through Deep Learning," In Proc. ACM CCS 18, Toronto, Canada, 2018, pp. 592-605.
- [9] Kyle Soska and Nicolas Christin, "Automatically detecting vulnerable websites before they turn malicious," In Proc. USENIX Security Symposium., San Diego, CA, USA, 2014, pp.625-640.