



## INTRUSION DETECTION SYSTEM USING BIG DATA FRAMEWORK

Abinesh Kamal K. U. and Shiju Sathyadevan

Amrita Center for Cyber Security Systems and Networks, Amrita School of Engineering, Amritapuri, Amrita Vishwa Vidyapeetham  
Amrita University, India

E-Mail: [abineshjerry@gmail.com](mailto:abineshjerry@gmail.com)

### ABSTRACT

In the enormous stream of network traffic, there is no way to identify which packet is benign and which is an anomaly packet. Hence, we intend to develop a new network intrusion detection model using apache-spark to improve the performance and to detect the intrusions while handling the colossal stream of network traffic in IDS. The model can detect known intrusion effectively using real-time analytics and hence identify unknown data schema compared with traditional IDS. The objective of the model addresses the following capabilities: Deep Packet Inspection (DPI) by inspecting the network traffic and examining the properties that describe the intrusion characteristics. Collaborating the vulnerability assessment with human intervention, using C.45 decision tree algorithm, optimizes pattern matching to boost detection rate. The clustered hosts are grouped based on their number of visits in a unique IP. The intrusion classifiers are developed by investigating each IP groups which reflects different properties used for prediction. The prediction model is built over Amrita Big Data Apache-Spark framework as a sequence of workflows. The above workflow is implemented in Amrita Big Data Framework (ABDF) to improve the detection time and performance, the model output provides effective results in detecting DOS attacks and port scanning attacks.

**Keywords:** apache spark, data mining, intrusion detection, network security.

### INTRODUCTION

Currently, the key research in the field of data mining has numerous shortcomings as follows: Data produced automatically integrates only small amount of intrusion characteristics, therefore rulemaking becomes inefficient. Most of the researchers go after Wenke Lee and Portnoy, because of their work upon integration of data mining algorithms such as fuzzy, genetic engineering.[1-4]. Even Though analysis in the automatic rule extraction is elementary, the data mining detection technique are not viable to identify new threats automatically. Another major factor is the evaluation criteria of the technique that includes the performance measurement and comparison of the normal behaviour of the system. Currently, accuracy and responses of any intrusion detection system (IDS) is not coherent and does not meet the demand of practical application.

Hence with the prior knowledge in Information security and data mining techniques, specific attributes have been carved and matched up with the network communications that reflect to the rules effectively. By registering apt data mining algorithm merging with appropriate anomaly detection together we deploy an efficient data-mining model based on network intrusion detection in apache-spark to improve the traditional methods.

### METHOD DESCRIPTION

A flawless IDS integrated with different detection techniques can hit up different types of intrusion. The major approach of this research paper is to distinguish legitimate and non-legitimate intrusions and their performance evaluation within the system on a criteria of Port Scanning and Deny of Service (DoS) [5-7]. Based on above methods we designed and developed a intrusion detection model deployed in apache spark to improve the

performance. The model is built by using training data which reflects the intrusion properties. To build the prediction model DARPA 1998 data set is used. This model detects new intrusions based on the classified training data. This paper dispute that the traffic source generated from various machines will have different properties with different patterns. On a server system, the traffic flow generated will be autonomous and does not send connection request in invoking. The workstation, host machine will act as an initiator of the TCP connections, the aim of the host machine is to obtain data. If a workstation machine receives huge volume of data suddenly then its behaviour is quietly abnormal and it should be monitored. To a server machine the huge number of concurrent active connection is quite normal, but if it exceeds the upper limit of the connections it should be also noticed. The decision of normal and abnormal activity is will be always unique. Each host will reflect a unique behaviour pattern as the services and streaming network traffic of hosts are different, which leads to expensive data processing. By knowing the above understandings this model verifies each and every requests based on unit time of the host in the network, then they are grouped into clusters based on their IP and visits time by Group - By features in the model framework[8]. Each IP set will reflect different behaviour pattern thus we can avoiding the individual differences' set generated helps to define how the data is distributed over the network, and the distributed training data will authorize their own classifiers by using C.45 decision tree algorithm.

### FRAMEWORK MODEL

The model framework [Figure-1] includes sequence of workflows, each workflow  $n$  takes input from the  $n-1$  workflow and pass the output to the  $n+1$ <sup>st</sup> workflow. ABDF includes several processing modules



such as Hadoop, Spark, GPGPU algorithms for large scale computing. It also have several process elements together with various file adapters such as KAFKA, flat file adapter, the data to the ABDF framework can be streaming data or an static file. Our experiments work on streaming data with the help of ABDF apache spark framework. [9]

### Data capturing by KAFKA

Network streaming data are collected from the core switch by KAFKA Pcap adapter in the framework. The network data are collected by an window size of 10. After collecting the data they are converted into Resilient Distributed Datasets (RDD). This form the input to the process flows (PF).

### Data preprocessing using ABDF preprocessor

Streaming network data include all the data that are transmitted over the network, our aim is to extract the necessary information that reflects the intrusion properties from the network packet. Once the data are converted into RDD we can select the necessary fields that are given input to the next PF.

- The ratio of deport\_same\_soip number based on destination port to all the total connection of same destination in a time window size "w".
- The ratio of deport\_diff\_soip number based on destination port to different source ip to all the total connections in a time window size "w".

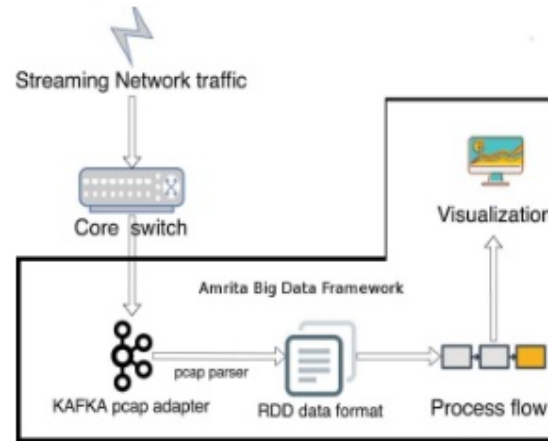


Figure- 1. Intrusion detection model workflow.

Table-1. Feature selection.

Fields selected for net PF
Date, Duration, Req_time, Service_Name, Source_Port, Destination_Port, Source_IP, Destination_IP.

### Parameters calculation

Eight different parameters are calculated to transform the selected fields in a meaningful way [10] in the References section include:

- soip\_diff\_deip : get the total count value of number of source IP connected to a different destination IP by using SQL Runner Query [Table-2] feature in the framework.
- soip\_diff\_deport number: get the total count value of number of source IP connected to a different destination hosts port by using SQL Runner Query [TABLE III] feature in the framework.
- deport\_same\_soip number: get the total count value of number of connections established by destination port to the same source IP based on the destination host window time "w".
- deport\_diff\_soip number: get the total count value of number of connections established by destination port to the different source IP based on the destination host window time "w".
- The ratio of soip\_diff\_deip number based on source IP to all the total connection of different destination IP in a time window size "w".
- The ratio of soip\_diff\_deport number based on source ip to all the total connection of different destination port in a time window size "w".

Table-2. SQL runner query for SOIP\_DIFF\_DEIP.

Fields selected for net PF
select Date, Duration, Source_Port, Destination_Port, count. Source_IP, Destination_IP, Attack_check, Attack_D escription, count.sip_diff_dip, count.sip_diff_dport from SqlRunner1 join (select Source_IP, count(Distinct Destination_IP) as sip_diff_dip, count(Distinct Destination_Port) as sip_diff_dport from SqlRunner1 Group By Source_IP) as count ON count.Source_IP=SqlRunner1.Source_IP"

Table-3. SQL RUNNER query for SOIP\_DIFF\_DEPORT.

Fields selected for net PF
select Date, Duration, Source_Port, count.Destination_Port, Source_IP, Destination_IP, Attack_check, Attack_D escription, sip_diff_dip, sip_diff_dport, count.dport_ diff_sip from SqlRunner2 join (select Destination_Port, count(Distinct Source_IP) as dport_diff_sip from SqlRunner2 group by Destination_Port) as count ON count.Destination_Port=SqlRunner2.Destination_Po rt

### Rule generation

To build the rule used DARPA 1998 dataset which have the details of intrusion properties with attack



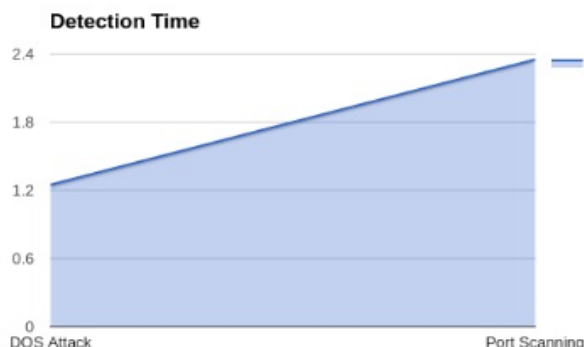


**Table-7.** Sample results obtained DoS attack- after parameter calculation.

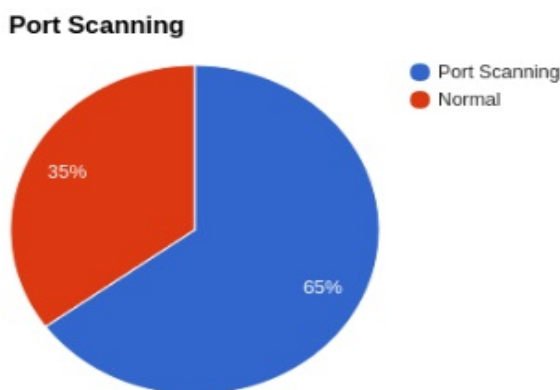
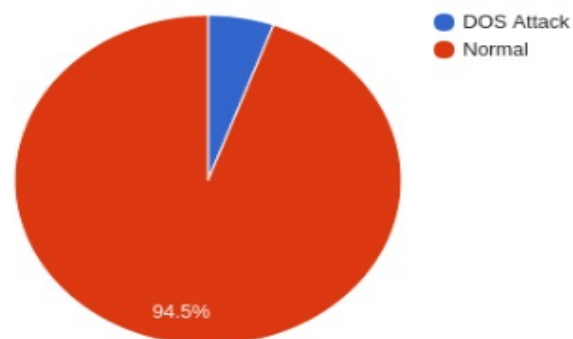
soip_diff t_deip	soip_diff t_deport	deport _same_ soip	deport_di fft_soip	Attack Flag
4.5683	0.004678	9.1367	4.5683	DoS
4.5683	0.004678	9.1367	4.5683	DoS

**Table-8.** Sample results obtained port scanning - after parameter calculation.

soip_diff t_deip	soip_diff ft_deport	deport _same_ soip	deport_di fft_soip	Attack Flag
0.002976	0.48214	9.1367	0.00297	port Scan
0.002976	0.47910	0.00297	0.00297	Port Scan

**Figure-3.** Performance analysis for prediction in dashboard.

From Figure-2 we can understand the total time taken to detect DoS attack and Port Scanning in Apache spark Framework.

**Figure-4.** Output of port scanning prediction in dashboard.**DoS Attack****Figure-5.** Output of DOS attack prediction in dashboard.

## CONCLUSIONS

Experimental results show that in data mining IDS system, previous Knowledge about information security is needed. Data mining techniques can able to extract rules from the raw data that are in specific format, but the rules derived by the rules need not to be reasonable. Moreover, with the help of our framework in the experimental setup provides High performance, detection rate and the execution time is low as the framework is built over apache spark.

## ACKNOWLEDGEMENTS

At the outset, I express my heartfelt gratitude to Shiju Sathyadevan and other faculties for their valuable guidance, timely suggestions and help in the completion of this paper. I express my heartfelt gratitude to all the staffs of the department of Amrita Cyber security systems networks, Amrita School of Engineering. I also thank the Evaluation panel for their valuable feedback, which made it possible in completing my paper. I thank all my friends who have helped me directly or indirectly, in the completion of this paper.

## REFERENCES

- [1] T. F. Lunt, R. Jagannathan, and R. Lee, "IDES: The Enhanced Prototype-A Real-Time Intrusion-Detection Expert System", Number SRI-CSL-88- 12. Computer Science Laboratory, SRI International, Menlo Park, CA, 1988.
- [2] Wenke Lee and Salvatore J. Stolfo, "Data Mining Approaches for Intrusion Detection", In Proceedings of the 7<sup>th</sup> USENIX Security Symposium, 1998.
- [3] Lee, Wenke, Salvatore J. Stolfo, and Kui W. Mok. "A data mining framework for building intrusion detection models." Security and Privacy, 1999. Proceedings of the 1999 IEEE Symposium on. IEEE, 1999.



- [4] W. Lee and S. J. Stolfo. A framework for constructing features and models for intrusion detection systems. *ACM Transaction on Information and System Security*, 3(4):227-261, Nov. 2000.
- [5] P.Porras, D. Schnackenberg, *et al.* "The Common Intrusion Detection Framework Architecture, CIDF", University of California, 1998.
- [6] Z. A. Othman, A. Bakar and I. Etubal, Improving signature detection classification model using features selection based on customized features. *Intelligent Systems Design and Applications (ISDA)*, 2010 10<sup>th</sup> International Conference on Nov. 29 2010-Dec. 1, 2010.
- [7] Z. Wanlei, Keynote: Detection of and Defense Against Distributed Denial-of-Service (DDoS) Attacks, *Proc. IEEE 11<sup>th</sup> International Conference on Trust, Security and Privacy in Computing and Communications*, 2012.
- [8] O. Al-Jarrah and A. Arafat , "Network Intrusion Detection System using attack behavior classification" , 2014 5<sup>th</sup> International Conference on Information and Communication Systems (ICICS) , pp.1 -6 .
- [9] Amrita Big Data Framework [online.available:<http://abdf.in/>].
- [10] Beniwal, Sunita, and Jitender Arora. "Classification and feature selection techniques in data mining." *International Journal of Engineering Research and Technology*. Vol. 1. No. 6 (August-2012). ESRSA Publications, 2012.
- [11] Yanjie Zhao, *et al.* "Network Intrusion Detection System Model Based on Data Mining", *Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD)*, 2016.
- [12] Padmashani, R., Shiju Sathyadevan, and Devi Dath. "BSnort IPS better Snort intrusion detection/prevention system." *Intelligent Systems Design and Applications (ISDA)*, 2012 12<sup>th</sup> International Conference on. IEEE, 2012.
- [13] 2dd [online. available: <http://2dd.it/>].