

## Study And Analysis Of Hadoop Based Network Intrusion Detection System.

Asst. Prof. Riyaz Ahmed Jamadar<sup>1</sup>, Himani Gupta<sup>2</sup>, Ankit Baghel<sup>3</sup>, Rituraj<sup>4</sup>,  
Nimish Bhandare<sup>5</sup>

<sup>1</sup>(Information Technology, AISSMS IOIT/SPPU, India)

<sup>2</sup>(Information Technology, AISSMS IOIT/SPPU, India)

<sup>3</sup>(Information Technology, AISSMS IOIT/SPPU, India)

<sup>4</sup>(Information Technology, AISSMS IOIT/SPPU, India)

<sup>5</sup>(Information Technology, AISSMS IOIT/SPPU, India)

**ABSTRACT:** Network security is a paramount concern for the organization. To secure the network, we have traditional network intrusion detection systems and firewalls but they have limitations like size of training data sets. With the inception of Hadoop technology, in industry, recently researchers have started using this new technology with traditional machine learning algorithms which generally uses pattern matching, to design and develop network intrusion detection system based on streaming of big-data using Hadoop that checks for intrusions in massive amount of data that flows in and out. In this paper, we are proposing a study and analysis of various Hadoop based network intrusion systems. Here the parameters used for comparison are detection rate and false-positive alarm rate.

**Keywords:** Big-data, Detection rate, Hadoop, Intrusion, Machine learning.

Date of Submission: 23-11-2017

Date of acceptance: 23-12-2017

### I. Introduction

Security is a major concern now-a-days as technology has reached new heights. On personal level security measures includes only installation of anti-virus and firewall. But when an organization is concerned, the solution cannot be simple. There should be a dedicated security system as the risks are many in a business organization

It is essential to have a safe and secure network for following reasons:

- To protect organization's assets
- To adhere with regulatory requirements and ethical responsibilities
- For a competitive edge

An Intrusion Detection System (IDS) is a security software that constantly monitors the network to look for suspicious or malicious activities and automatically alert the administrator. In other words, it is equivalent to a Burglar Alarm. The types of intrusion detection systems are: A host based intrusion detection and a network based intrusion detection.

Host based intrusion detection system: Host computers are installed with software and are used to analyze system activities, log files and all network-traffic received by host computers. It can discover whether an attack is successful and also record what the attacker has performed on the host.

Network based intrusion detection system: Sensors are installed throughout the network to analyze all traffic on target network in real time. These sensors describe a network-interface that receives all network-traffic and matches the defined pattern. If the pattern matches the system alerts the administrator.

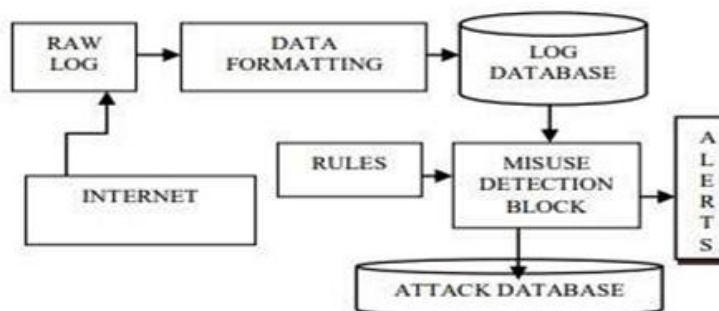


Fig 1: -Working of Network Intrusion Detection System [9]

IDS TRIGGERS: The main motive of IDS is to send an alert if an intrusion is detected, it works just like a burglar alarm.

IDS have two triggering mechanisms:

1. Misuse Detection (Signature based)
2. Anomaly Detection (profile based)

Signature based detection: This type of mechanism requires signature based files i.e. known patterns of attacks. If a pattern is matched, there is a high probability of an attack, and therefore an alarm will be triggered. It has low false-positive alarm rate as matches are based on known patterns. Signature detection fails to detect attacks which are not known or variations in the known attacks.

Anomaly based detection: This type analyzes computer and network activities and looks for an anomaly, if an anomaly is found alarm is triggered. Anomaly is abnormal behavior of network or deviation from common rule for attacks. It has high false-positive alarm rate.

The rise in use of technology has led to increase in amount of network traffic data. The traffic data is expected to be in the range of Zettabytes and is significantly increasing from past few years. This data having high Volume, Velocity and Variety is often termed as Big-data. In this generation of big-data the IDS should be good enough to process huge volume of data in real time. To overcome this challenge and the need of IDS of high accuracy and efficiency while running in parallel environment has led to introduction of IDS using Hadoop that reduces the extra overhead and does not reduce the performance. The rest of the paper has following sections. Section 2 discusses existing work done in Hadoop based IDS Section 3 provides analysis of various IDS Algorithms. Section 4 provides a conclusion of the study with comparisons.

## **II. Literature survey**

A] The work proposed by Zhiguo Shi et.al. [1] titled “An Intrusion Detection System based on Hadoop”. Here Hadoop cluster along with improved k means algorithm, as the defined K means algorithm is only used for detecting single intrusion and not for multiple intrusion detection which results in forming a distributed IDS for manipulating big-data. They proposed the use of Map reduce to select the stable centres of the k means clusters and impose it to the Hadoop clusters to form distributed IDS. The limitation of their work is that the improved K means algorithm gives accuracy of 95% which is higher than traditional K means algorithm but is still lower than decision tree which has accuracy of 96%. Also the use of Map reduce adds burden and can't be used for streaming data.

B] The Work proposed by M. Mazhar Rathore et.al. [2] titled “Hadoop based Real-time Intrusion Detection for High-Speed Networks.” includes a Hadoop based real time intrusion detection for high speed networks the main challenging problem is to identify intruders in high speed networks. So to solve this, they proposed a high speed IDS which is capable of working in big-data. This paper makes use of KDD 99 dataset from which all parameters were selected, further using Forward Selection Ranking(FSR) and backward elimination ranking(BER) only 9 best features were used for achieving faster results. Among the widely used classification algorithms, REPTree and J48 Machine Learning classifiers were selected for classifying intrusion. The proposed work makes use of Map reduce programming using single node Hadoop architecture to accurately and efficiently detect threats and to get real time efficiency Apache Spark is used. Traffic is captured using capturing device like RF-RING and TNAPI which assures that no packet remains un captured. Filtration and load balancing server(FLBS) are used for searching and comparing in intruders database. Main problem is with the data set i.e KDD99 which is not advisable by authors itself to use it, as well as architecture uses single node so it becomes a single point of failure if that node fails the system may go down.

C] The work proposed by Manish Kumar et.al. [3] titled ”Scalable intrusion detection systems log analysis using cloud computing infrastructure” implemented the IDS log analysis using the cloud architecture by taking the help of Hadoop and Map reduce. They created the distributed system in order to extract data efficiently with the help of HDFS. Using Map Reduce they merged the two alerts and store in log files. The essential information from log files produced by IDS is extracted by the log parser which is mapped to different nodes depending on the extracted Meta information. Then the information is reduced and stored in log files for the further detection. Limitation of their work is that Map Reduce increases the burden and is not adaptive and their proposed work detects only signature based intrusion.

D] The work Proposed by Sanraj Rajendra Bandre et.al. [4] titled “Design Consideration of network intrusion detection system using Hadoop and GPGPU” includes a method which deals with NIDS based on Hadoop framework and GPGPU. Big-data from organisation is applied to Hadoop framework while GPGPU is used for intrusion detection. Among the various Hadoop ecosystem like PIG, HIVE and HBASE, FLUME is used to pull collect real time streaming data. In propose system a parallel failure-less AHO-Corasik (PFAC) algorithm is used for multi string pattern matching algorithm and it removes failure transition from state

transition machine. The problem with the proposed system is to keep on updating network security with modern technologies.

E] The work proposed by Basappa Kodada et.al. [5] titled “Big-data analytics based security architecture to detect intrusions” implements a method to overcome hurdle of finding malicious activities in zettabytes of network packets. This paper evaluates the security architecture of huge amount of data that flows in and out. To process big-data, a Hadoop Map Reduce framework is used. An open source packet analyzer called Wireshark is used for network trouble shooting, as well as it allows user to put network interface which supports promiscuous mode. Big-data analytics can be used to emphasize financial transactions, log files and network traffic to detect suspicious activities as well as identify anomalies. In the proposed work results show that analysis and computation is much faster than other systems. The limitation of their work is that Map Reduce degrades the performance of the system as it is complex and time consuming in nature. Increase in the number of nodes in their proposed system can achieve better performance and results.

F] The work proposed by Gurpreet Kaur Jangla et.al. [6] titled " Development of an intrusion detection system on big-data for detecting unknown attacks" in which the researchers divides the entire IDS system into 4 parts. Firstly, the data is collected from various sources, then the data is processed and analyzed. If attack is occurred then the gets alarm. They propose the use of HDFS and Map reduce by mapping the (key, value) pair to process the data. Then the data is analyzed by using classification algorithm and prediction is done on the data. The limitation of their proposed work is although the data is processed and analyzed but Map reduce can't be used for streaming data and increases the complexity. Also, it can only predict the signature based attack and still requires advancement to find advanced threats.

G] The work proposed by Sonali Ashok Hajare [7] titled “Detection network attacks using Big-data analytics.” Includes processing of unstructured data from socio-networking sites and converting them into structured data. Map reduce is applied on this data for faster retrieval. Pre-processing of web log data is done through SVM and c means clustering. If the pattern is matched then intrusion is detected. Map Reduce extract the features from the structured data set and the key, value pair is matched with the predefined signature already available. If matched with the pattern then the intrusion is detected. If any suspicious IP packet is detected then it gets cleaned using c-means clustering. The limitation of this work is Map reduce is a huge overhead on the system as it degrades the performance of the system. More efficient machine learning algorithms other than SVM and c-means clustering can be used to improve the system further.

H] The work done by Shaik Akbar et.al. [8] titled“ A hybrid Scheme based on big-data analytics using intrusion detection system” proposed an architecture in which KDD cup dataset is used to gather data which are not similar and these data are separated into learning and detection phase . Big-data along with its 5 V's VELOCITY, VOLUME, VARIETY, VALUE and VERACITY and briefly discussed with their usage. A feature selection process is applied on KDD cup data set to select and extract main features. In learning phase known attacks like DOS, pros are detected while the attacks which are not been identified in learning phase they are detected in detection phase. In detection phase enhanced C4.5 and enhanced genetically algorithm are used. These two phases stores the data in large database forming an integrated hybrid large database. These hybrid techniques will enhance the detection rate by identifying the category of attack. The limitation of their work is due to advancement in technology the need to address real time data is essential.

**Table 1:** - Comparison of various classification Algorithms [1], [2], [10].

CLASSIFICATION ALGORITHM	ACCURACY (%)
K-MEANS	90
IMPROVED K-MEANS	95
REP Tree	72
J48	70
SUPPORT VECTOR MACHINE(SVM)	84
DECISION TREE	96

### III. Analysis

In this section, we have analyzed various IDS algorithms with respect to detection rate, false alarm rate and data type of respective algorithm.

**Table 2:** - Comparison of various IDS Algorithms with Detection Rate, False Alarm Rate and Data Type [11] [12] [13]

SR NO	IDS ALGORITHM	DETECTION RATE	FALSE ALARM RATE	DATA TYPE
1	PCA, SVM and PSO	97.75	Not Mentioned	Conventional
2	Ada-Boost Algorithm	90.88	3.42	Conventional
3	K-means and Decision tree(C4.5) Algorithms	99.6	0.1	Conventional
4	Online Ada-Boost based	90.99	0.31-1.78	Conventional
5	A hyper spherical cluster	85.47	1.48	Conventional
6	PCA based	82.86	13.3	Conventional
7	Online Adaptive PCA classifier	97.84	1.10	Conventional
8	PCA with radial SVM	97.85	0.54	Conventional
9	PCA with GMM using Spark	86.2	13	Conventional
10	Our proposed work with Decision tree using Apache Flink *expected	96-99	0.1	Streaming (Real Time)

### IV. Conclusion

It has been studied and analyzed, here that Hadoop based IDS are quite suitable for handling large training datasets. In this work, we compared various Hadoop based IDS, and also conclude that SVM, k means clustering, REP Tree and J48 classification algorithm have accuracy less than 95%. Map reduce is used to map the data based on (key, value) pair which is suitable for static data but due to advancement in technology, the need to work on streaming data is essential. The use of decision tree which has accuracy of 96% can further improve the system.

### References

- [1] Zhiguo Shi, J. A. (2015). "An intrusion detection system based on Hadoop". *IEEE*.
- [2] M. Mazhar Rathore, A. P. (2016). "Hadoop based real-time intrusion detection for high-speed networks". *IEEE*.
- [3] Manish Kumar, D. M. (2013). "Scalable intrusion detection systems log analysis using cloud computing infrastructure". *IEEE*.
- [4] Sanraj Rajendra Bandre, J. N. (2015). "Design Consideration of network intrusion detection system using Hadoop and GPGPU." *IEEE*.
- [5] Basappa Kodada, S. P. (2015). "Big-data analytics based security architecture to detect intrusions." *NJCIET*.
- [6] Gurpreet Kaur Jangla, D. .. (2015). "Development of an intrusion detection system on big-data for detecting unknown attacks". *IJARCCCE* vol.4.
- [7] Sonali Ashok Hajare(2016). "Detection network attacks using Big-data analytics." *IJRITCC* vol. 4.
- [8] Shaik Akbar, T. R. (2016). "A hybrid Scheme based on big-data analytics using intrusion detection system." *IJST* vol. 470524303359/Fig-1-Block-Diagram-of-Intrusion-Detection-System-The-figure-shows-the-block-diagram-of.ppm
- [9] [https://www.researchgate.net/profile/Premchand\\_Ambhore/publication/267390056/figure/fig1/AS:392221510651908@1470524303359/Fig-1-Block-Diagram-of-Intrusion-Detection-System-The-figure-shows-the-block-diagram-of.ppm](https://www.researchgate.net/profile/Premchand_Ambhore/publication/267390056/figure/fig1/AS:392221510651908@1470524303359/Fig-1-Block-Diagram-of-Intrusion-Detection-System-The-figure-shows-the-block-diagram-of.ppm)
- [10] Baojiang Cui, S.H.(2016). "Anomaly detection model based on Hadoop platform and Weka interface" *IEEE*.
- [11] Riyaz A Jamadar, Prof. Ms.Mousami S Vanjale (2015). "Enhanced detection rate through PCA and radial SVM in Wireless Sensor Networks." *IERJ* vol.9.
- [12] Keisuke Kato, Vitaly Klyuev. (2017) "Development of a Network Intrusion Detection System Using Apache Hadoop and Spark." *IEEE*.
- [13] Wei Hu, Weiming Hu. (2005) Network Based Intrusion Detection using Ada-Boost Algorithm. " *IEEE*

International Journal of Engineering Science Invention (IJESI) is UGC approved Journal with Sl. No. 3822, Journal no. 43302.

Asst. Prof. Riyaz Ahmed Jamadar "Study And Analysis Of Hadoop Based Network Intrusion Detection System.." International Journal of Engineering Science Invention(IJESI), vol. 6, no. 12, 2017, pp. 01-04.