

Efficient Analysis of the Big Data IDS Over Proposed Model

¹B. Priyanka, ²Indivar Shaik

ABSTRACT--In the last ten years, the Internet has grown quickly. As a result, computer and network device interconnection has become so complicated to monitor, that even security experts do not fully understand their deepest internal functions. Every year, personal computers have become very quick. It is not unusual for a very ordinary individual to connect via or faster than 20 Mbs to the Internet. The security of the network has become very important for data monitoring with this enormous network data. Big data in the intrusion detection system are a major challenge to develop. In this article, the framework for pre-processing functionality was used to create sub-sets of features related to template creation. The algorithm Random Forest was used to categorize data for the network. The knowledge benefit approach was used to improve the accuracy of the Random Forest Algorithm. To check the performance of the model proposed, the NSL-KDD standard data was used. Several assessment metrics have been suggested to assess the model proposed. The empirical results of the model proposed show that performance measures are better. The results of the proposed model and various existing algorithms are comparatively analyzed. The results show that the performance of the proposed model was higher than that of existing systems.

Keywords—Big data ids proposed model.

I. INTRODUCTION

The safety of networks has become a very important component of worldwide infrastructure, and given that personal, e-commerce, banking, and business data are shared in computer networks, safety has become one of Internet's major aspects. The Intrusion Detection System (IDS) is one of the essential areas in network security. Intrusion network management and protection is difficult, because network security teams are challenged with targeting and detection of threats. New vulnerabilities have been identified annually. It is more difficult for the security tool to automate the detection of new attacks. The intrusion detection system has become very important and helpful for preventing attacks by the computer network. Of example, most companies around the world are using firewalls to shield their hidden network data from public networks. The firewall can be used to protect the assets from the users, but it is not possible to achieve protection of the information as a whole. In addition, the network security aspects of the intrusion detection system are very important for protecting the network and detecting adversaries of network activities. The IDS tool works on the assumption that the signature of an attack activity differs from normal activity signatures. The intrusion sensing device is fitted with two ways to identify threats, whether by signature or by deviation. To order to detect connection matches, signature analyzes are used

¹ Research Scholar, Dept. of Computer Science and Engineering, SSSUTMS, Sehore, Madhya Pradesh, India.

² Research Scholar, Dept. of Computer Science and Engineering, SSSUTMS, Sehore, Madhya Pradesh, India.

to scan a list with known assault signatures. Whilst anomaly detection uses normal baseline monitoring and can issue abnormal behavior-based alerts.

II. BIG DATA IN INTRUSION DETECTION SYSTEM

This section outlines the comprehensive background and context for large-scale data detection challenges. Every day a great deal of digital data is generated as the software is through. Gross data is a term when the challenge of traditional approaches is faced by a large amount of data. In recent years different methods have been developed. So finance big data. Then technology big data. Big data is a key issue for many researchers to find a specific solution to many more security breaches such as finances, industry, medicine and other important issues due to the Network Security Use.

For any attack pattern act or unexplained packets an internal network, IDS is a security solution. Because of the complexity of network data, large data techniques are very essential for the analysis and determination of network patterns? Such as network data with many different structures and network types that are difficult to analyze using traditional methods and to collect network data from a different source. In contrast, high-dimensionality network databases face major problems. Machine learning is commonly used to help build zero-day intrusion sensing systems.

In this work of study, it focused mainly on the large data on intrusion detection through algorithms for machine learning. Intrusion from network large data is detected using the proposed method of information using random forest approaches. The inspiration behind this work is the new model introduced to help detect invasion more accurately and quickly.

In addition, for the collection of dimensionality reduction data, feature selection processes are essential. When the amount of characteristics of the data set is limited, the time and precision of a classifier is significantly reduced.

III. METHODOLOGY

The main objective of the present work is to examine the large data in the intrusion detection framework. The proposed Big Data Intrusion Detection System model will be shown in Figure 1.

An experimental study involves normal packets and abnormal packets with standard intrusion sensing data set. The data are pre-processed with the entropy method used for details. In order to build the classificatory with more precision and speed the information entropy method is used.

IV. DATA SET

For the validation of the proposed model IDS NSL-KDD is used. Data Set is a modified KDD cup'99 edition data set for NSL-KDD Standard Intrusion Detection Program. The NSL-KDD data set was developed to solve the problem that was addressed by McHugh [8] related to the KDD cup question in 1999. Can run the full set of experiments without choosing a small part randomly. The data set for NSL-KDD contains 4.898.431 entry.

The NSL-KDD data set is stored as individual network packets and is not a standard program or operating system. As such, a tag representing the category mark of the record was assigned to each record in this dataset. All

tags should be right in this data set. There are 37 forms of assault on the NSL-KDD. The simulated assaults fall into precisely one of the four categories: Server Denial, Test, Root and Regional Remote. The Table displays all types of NSL-KDD data set attacks.

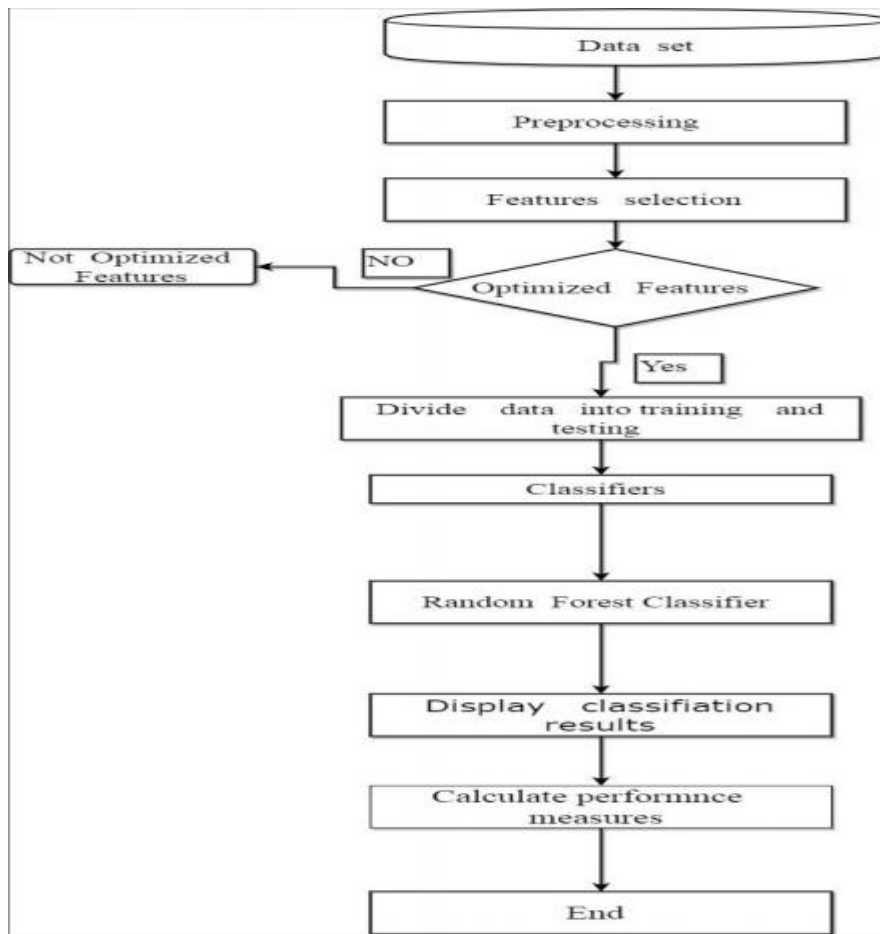


Figure 1: Proposed model

V. PREPROCESSING

Preprocessing is an essential step in managing real-world data sets in an accessible layout. Probably, in certain actions the real world samples are imperfect, noisy. For understanding trends in big data, the preprocessing step is very important. In order to improve the machine learning algorithms for the design classification, pre-processing procedures are therefore required in the big data intrusion detection framework.

VI. EXPERIMENTAL SETUP

MATLAB R2013a-64 windows 7 Ultimate with i5 core and 8 GB of RAM, the proposed model Big Data intrusion detection system is implemented, and the Weka tool has been used to compare different classification algorithms. Different performance measurements were used to test the model proposed. Hybrid model random forests and algorithms are applied in this experiment. 31 major attacks have been identified in this study.

The data only contains 18559 attacks and normal instances. The attacks in the total set of data correspond to 185559. The original set of data includes 25 MB of data. Matlab uses the hybrid model of Random Forest and

algorithms for information benefit. The results of the proposed model are shown in Table 4. This investigated that 184331 of 18559 instances were the right category of the case. In contrast, out of 185559 instances 1228 instances are misclassified.

Performance	Proposed model
Time	16.87 seconds
Correctly Classified Instances	184331
Incorrectly Classified Instances	1228
Total Number of Instances	185559

Table 4: Performance analysis of proposed model

The decision is made to work with preprocessing to improve the proposed model. Method of function selection. To strengthen the Random Forest Classifier, the method of gaining information is introduced. The process of choosing features helps to increase recognition reliability and reduce the time of template development. The number of subsets of relevant features from the original data sets was challenged when the goodness features were obtained.

Finally, the characteristics have been selected to increase the accuracy of the classification. Table 5 demonstrates the results by using a feature selection method for numerous current classifiers against the proposed model. The suggested method has been shown to outperform all current algorithms better than any other.

Classifiers	FP	TP	Accuracy	Precision
Naïve Bayes	0.003	0.949	94.9261	0.949
REP Tree	0.003	0.988	98.767	0.721
SVM	0.004	0.957	95.43	0.987
KNN	0.007	0.933	93.12	0.975
Proposed model	0.001	0.993	99.33	0.993

Table 2: Results of proposed model with different existing algorithms

VII. PERFORMANCE AND COMPARISON OF PROPOSED MODEL

The comparative method for evaluating and analyzing the current intrusion detection system design as regards intrusion detection system classifications. In order to investigate the proposal with all the various existing algorithms FP, TP, accuracy and precision performance measures are used. Table 1 summarizes the outcome obtained by using the feature selection method for the proposed and current algorithm. It is stated that greater than all the current algorithms, the effects of the proposed model are. Figure 2 the performance of the model proposed in terms of accuracy compared to various existing algorithms.

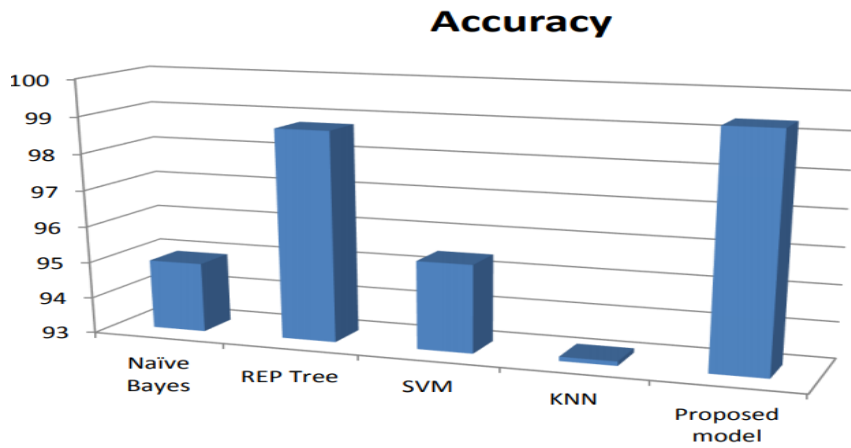


Figure 2: performance accuracy of proposed model and existing models

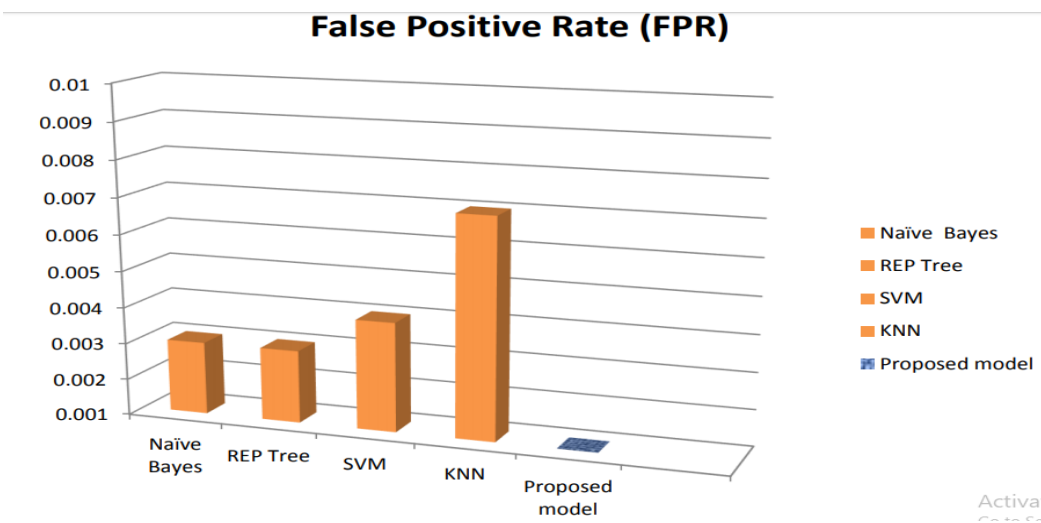


Figure 3: performance FP of proposed model and existing models

This means that the template suggested is more accurate and less time is needed to build this. Figure 3 reveals that the proposed model output and the existing classification algorithms are false positive, and it suggests that the FP of the proposed model is much lower.

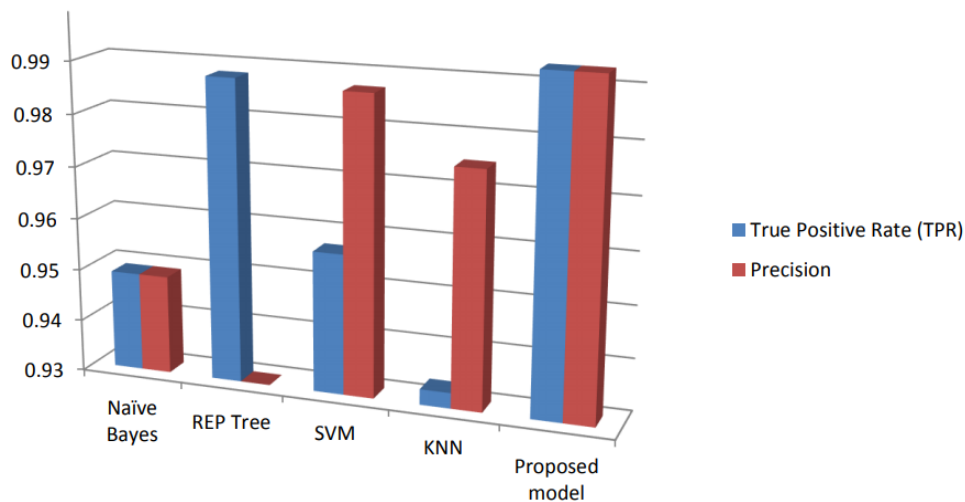


Figure 4: performance TP and precision of proposed model and existing models

Relative to another current classifier shown in figure 4, TP and accuracy measurements of the proposed model are the best. Ultimately, it is assumed that different types of attack can be observed with the most accuracy relative to other existing systems in the proposed model.

VIII. CONCLUSION

The research objective is to improve the existing intrusion detection system building algorithms. By using the proposed model, the main objective was achieved. The randomness and the accuracy of Big Data in the intrusion detection system have been improved with the proposed methods. All kinds of attacks have been used in this research. Because of the large data, the IDS system has been built. The selection method for features is however used to resolve the problem. The 13 key subset features from the original 41 software functions are created after knowledge gain. In contrast, these most significant elements were added to the proposed model. There has been an improvement in reliability and decrease in the period of production of the prototype. In evaluating the proposed model, the different performance measures have been used. The exactness of the model proposed is 99.33%. There is a descriptive study of the findings between the template suggested and the different existing ones. The new design has been found to surpass the current classification systems. The author will in future try to use soft computation through the use of different data sets.

REFERENCES

1. R. Chitrakar, and C. Huang, "Irregularity based Interruption Detection utilizing Hybrid Learning Approach of consolidating k-Medoids Clustering and Naïve Bayes Order," In *Wireless Communications, Systems administration and Mobile Computing (WiCOM)*, eighth Worldwide Conference on, pp. 1-5, IEEE, 2012.
2. M. Dhakar, and A. Tiwari, " An epic information mining based mixture interruption location structure," *Journal of Information and Computing Science*, vol 9, no. 1, pp. 037-048, 2014.
3. W. Huai-canister, Y. Hong-liang, X. U. Zhi-Jian, and Y. Zheng, "A grouping calculation use SOM and Kmeans in intrusiondetection," In *E-Business and EGovernment (ICEE)*, 2010 International Conference on, pp. 1281-1284, IEEE, 2010.
4. S. Warnars, "Mining Patterns with Attribute Oriented Acceptance," In *Proceeding of The International Meeting on Database, Data Warehouse, Data Mining and Big Information (DDDMBD2015)*, pp. 11-21, 2015.
5. V. Kachitvichyanukul, "Correlation of three developmental calculations: GA, PSO, and DE," *Mechanical Engineering andManagement Systems*,11(3), pp. 215-223, 2012.
6. R. Chitrakar, and C. Huang, "Irregularity based Interruption Detection utilizing Hybrid Learning Approach of consolidating k-MedoidsClustering and Naïve Bayes Order," In *Wireless Communications, Systems administration and Mobile Computing (WiCOM)*,8th Universal Conference on, pp. 1-5, IEEE, 2012.
7. Shi, X., Manduchi, R., 2003. A concentrate on Bayes include combination for picture order. In: *Gathering on Computer Vision and Pattern Acknowledgment Workshop, CVPRW, Madison,*

8. Wisconsin, USA, pp. 95–95.
9. <http://www.kdd.ics.uci.edu/databases/kddcup99/task.html> 7
10. Nassar M, al Bouna B, Malluhi Q (2013) Secure re-appropriating of system stream information examination. In: Big Information (BigData Congress), 2013 IEEE International Congress On. IEEE, Santa Clara, CA, USA. pp 431– 432
11. Kezunovic M, Xie L, Grijalva S (2013) The job of enormous information in improving power framework activity and assurance. In: Bulk Power System Dynamics and Control - IX Optimization, Security and Control of the Emerging Power Grid (IREP), 2013 IREP Symposium. IEEE, Rethymno, Greece. pp 1–9
12. Denning DE (1987) An interruption location model. *Softw Eng IEEE Trans SE-13(2):222–232*. doi:10.1109/TSE.1987.232894
13. Suthaharan S, Panchagnula T (2012) Relevance include choice with information cleaning for interruption location framework. In: Southeastcon, 2012 Procedures of IEEE. IEEE, Orlando, FL, USA. pp 1–6
14. Marcelo D. Holtz, Bernardo M. David and Rafael Timeote "Building Scalable Distribute Intrusion Identification System Based on the Map Reduce Structure. 2011, Intrenation diary of Revista Telecommucation pp 23-31
15. Lidong Wang*, Randy Jones "Enormous Data Analytics for System Intrusion Detection: A Survey. *Global Journal of Networks and Interchanges* 2017, 7(1): 24-31 DOI: 10.5923/j.ijnc.20170701.03
16. Jingwei Huang, Zbigniew Kalbarczyk, and David M. Nicol. "Information Discovery from Big Data for Interruption Detection Using LDA. 2014 IEEE Global Congress on Big Dat pp760-762
17. Rachana Sharma and Priyanka Sharma, Preeti Mishra and Emmanuel S. Pilli "Towards MapReduce Based Characterization approaches for Intrusion Detection". *Intrnation meeting 2016 IEEE PP-361-366*
18. Miss Gurpreet Kaur Jangla¹, Mrs Deepa.A. Amne²." Advancement of an Intrusion Detection System based on Big Data for Detecting Unknown Attacks. *Universal Journal of Advanced Research in PC and Communication Engineering* Vol. 4, Issue 12, December 2015