

Intelligent Feature Extraction through Bigram& Trigram Schemes for Detecting Network Intrusions

^[1] T. Augustine, ^[2] P. Vasudeva Reddy, ^[3] P.V.G.D. Prasad Reddy
 ^[2] Department of Engineering Mathematics, ^[2] Department of CS&SE
 ^{[1][2][3]} A.U. College of Engineering, Andhra University, Visakhapatnam, AP, INDIA

Abstract— In the present scenario of network architecture, handling long payload features is a challenge, specifically because many machine learning algorithms are not able to process these long payload features. Some of the Network Intrusion Detection Systems (NIDS) are completely avoiding these long payload features. To address this challenge, a new methodology called feature extraction through Bigram and Trigram techniques has been proposed. The long payload features are encoded though these proposed techniques and are prepared to be used in machine learning algorithms. Experiments were carried out on ISCX 2012 data set. The designed feature selection based system has shown a noticeable improvement on the performance using different metrics.

Keywords:-- Network Intrusion Detection, NID, Intrusions, Feature Extraction, Bigram, Trigram Scheme, Network Lifetime, Long Payload Features, Dataset, Dictionary Building.

I. INTRODUCTION

Generally, in real time network traffic, the payload features are long and are of different data-type. It becomes extremely difficult to any intelligent systems like machine learning systems to handle such long payload features. In the present work, lot of experiments was carried out on ISCX data set. In the literature, researchers used different IDS data sets for testing their models. However, in this paper the ISCX 2012 intrusion detection data set is used for better comparisons in the results because some of the earlier and recent works [16] used this same data set. This data set has been generated by the Information Security Centre of Excellence (ISCX) at the University of New Brunswick in 2012 [1]. The data set consists of real traces analyzed to create profiles for agents that generate real traffic for FTP, HTTP, SMTP, SSH, IMAP, POP3, etc., [1]. The generated data set contains various features that includefull packet payloads in addition to other relevant features such as total number of bytes sent and /or received. This ICSX Data set consists of different types of features like numeric, alpha numeric, date, time, categorical and strings. Usually the packet header information is represented by a combination of these above types, but the payload features are usually represented by long string values which contains very long strings that makes it very difficult for any machine learning algorithms to deal with. To address this problem, encoded schemes have been chosen to encode these features by using bigram and trigram techniques. Fig. 1illustrates the main steps of the feature extraction process that is employed to extract features using a proposed scheme. The present approach and algorithms are very similar to the procedures in [16] but in the present study, along with bigram scheme, trigram approach is also studied and experimented on the same data set. This bigram / Trigram techniques are used with payload features to investigate if the payload features contain informative features or not. It is opted to do this since most research ignores these features due to their long strings, which makes them difficult to utilize in machine learning.

II. A VIEW ON THE PREVIOUS WORK

Lot of contributions are made available in the literature by various researchers to build more efficient Network Intrusion Detection Systems (NIDSs). In the very recent past, TarfaHamed, Rozita Dara, Stefan C. Kremer [16] designed a Network intrusion detection system based on recursive feature addition and bigram technique, in which they proposed a new feature selection method called RecursiveFeature Addition (RFA) and bigram technique. In fact, this work gives motivation to this present study.In this section, some of the papers that were cited in [16] were also studied.

Apart from the above work, Studies have been conducted on applying feature selection to improve the IDS performance. In [1], the authors applied the intra-class correlation coefficient and interclass correlation coefficient to attaina class-specific subset of features. The interclass and intra-class correlation coefficients were used to measure the validity and the reliability of features respectively. The authorsShiravi A, Shiravi H, Tavallaee



M, Ghorbani AA [1] tested their model on the ISCX 2012 data set. They observed that the above combination between interclass and interclass correlation coefficients led to an increase in the detection rate and to a decrease in both execution time and false alarm rate. However, their work did not deal with the scarcity of data and interdependent features as the authors in [16] did in their work.

In other studies [2], the authors opted to build their intrusion detection system based on the normal traffic to detect unseen intrusions using the ISCX 2012 data set. The authors employed a one-class Support Vector Machine (SVM) classifier to learn http regular traffic attributes for an anomaly detection task. Their approach involved extracting appropriate attributes from normal and abnormal http traffic. The system generates an alert if it finds any deviation from the normal traffic model. The authors stated that they obtained 80% accuracy and 8.6% false alarm rate in detecting attacks on port 80. Authors in [16] stated that their work differs from the work in [2] in dealing with normal and attack data instead of dealing with normal data only. As the present work imitates the work in [16] with respect to experiments on the data set, both normal and attack data is considered in this work.

Zero-day attacks have made known to be intricate to alleviate their damage due to the lack of information [17].For this reason, there is always a need to protect against these zero-day attacks before they cause enormousdamage to networks. These attacks are also called "zero-day misuses" [3][4]. As just mentioned, Zero-day assaults have appeared to be hard to lighten their harm because of the absence of data [5][6]. Consequently, there is dependably a need to protect against these zero-day assaults before they make tremendous harm to the systems. Information mining is a method that can be utilized with interruption identification to distinguish trademark designs from the information included in that portray framework and client conduct [7][8], and preferably, cases of pernicious action. Machine learning calculations have been utilized broadly with interruption discovery to improve the precision of identification and making a safe model for the IDS against zero-day assaults or novel assaults [9][10].

To construct quick and exact IDS, it is essential to choose enlightening highlights from the information. Highlight determination has demonstrated its capacity to diminish calculation requests, over fitting, display size and increment of the exactness [8][11]. The trouble that faces

an engineer assembling these sorts of frameworks is the shortage of assault illustrations which can be utilized to prepare a learning machine to manufacture a model for identifying that specific assault. Indeed, even powerful machine learning calculations battle when there are couple of illustrations, or unequal cases and substantial quantities of highlights. The accessible useful highlights likewise influence the execution (that is the more the better). Past IDSs regularly ignored the payload highlights in spite of the fact that they contain some helpful data [12][13]. In this way, we chose to use the payload highlights and concentrate helpful data for ID purposes. Keeping in mind the end goal to enhance the identification capacity of the framework, The Bigram and Try gram strategy was utilized to encode the payload highlights into a shape that can be utilized as a part of machine learning calculations. The Bigram system is a set up procedure particularly in Deep Packet Inspection (DPI) and has been contemplated for quite a long time [14][15]. Be that as it may, in this system, another mix of utilizing highlight choice, the Bigram procedure and the application to this specific issue (interruption recognition) is exhibited. Experiments were carried out to address the issue of interruption discovery harder by concentrating on "zero-day assault" situation. With a specific end goal to reproduce this, it is deliberately fabricated a learning machine utilizing little quantities of cases and extensive quantities of highlights. The reason for that is to check on the off chance that can be even now be distinguished assaults with an informational index with the above attributes.

As outlined in this section on some of the studies in the area of NDIS, it may be concluded that in spite of decades of research in this area, handling long payload features on the network traffic still remains as a challenge.

III. PROPOSED BIGRAM & TRIGRAM SCHEMES FOR INTELLIGENTFEATURE EXTRACTION

The bigram technique is an established technique especially in Deep Packet Inspection (DPI) and has been studied for decades [14]. However, in this paper, not only a new combination of using the bigram technique and feature selection is experimented as in [16], but also a new proposed trigram technique and the application of these two schemes over the long payload features is presented.

In the proposed methodology the initial step in the feature extraction process for all payload features is construction



International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)

Vol 5, Issue 1, January 2018

of the dictionaries. Consecutively, to extract the feature vector for each payload feature, a two separate dictionaries need to be built that contain all the bigrams and trigrams respectively. Two schemes have been proposed (Bigram & Trigram), and so two separate dictionaries have been created to perform various experiments. This concept is explained in this paper by taking various examples as it is done in [16]. The feature extraction process in this approach is depicted in figure 1, which is similar to the approach in [16].

The long payload features are encoded using a Bigram / Trigram techniques. Fig. 2 illustrates the main steps of the feature extraction process that is employed to extract features using a Bigram / Trigram techniques. The bigram scheme has already been implemented in [16], in the present work the bigram approach is once again experimented on slightly different data and along with it the Trigram approach has also been experimented and results are noted and are presented for comparative analysis. These Bigram / Trigram techniques are used with payload features to investigate if the payload features include informative features or not.



The detailed steps for the dictionary construction are shown in the following algorithm 1. Even though two encoding schemes are proposed and experimented in the present work, a single algorithm is presented because of the amount of similarity. The only difference between these two is, in bigram, two adjacent words are taken from a feature where as in Trigram three are taken. Hence a common algorithm is presented below. The methodology of constructing Bigrams dictionary and Trigrams dictionary are similar. The output of the algorithm will be two different dictionaries with Bigrams and Trigrams respectively.

Algorithm 1: Dictionary Construction for Long Payload Features

- Step-1: Input the long payload features
- Step-2: Initialize the dictionary 'D', with empty
- Step-3: Take one feature from the long payload
- Step-4: Take one Bigram / one Trigram from the feature

Step-5: Check if the dictionary consist this

Bigram / Trigram already

Step 6: If dictionary D does not contain Bigram / Trigram then add feature to dictionary D

Step-7: Repeat step 4 to 6 till there is no possibility of new Bigram / Trigram from the features

Step-8: Repeat the procedure till all the payload features are encoded as Bigrams / Trigrams and added to dictionary D

Step-9: Output: Dictionary D



Fig.2 Dictionary construction stage during feature extraction for ISCX Data set using Trigram Technique

Figure 2 shows the block diagram of dictionary construction. Now, this dictionary which has been built with Bigrams and Trigrams can be used by any feature extraction algorithm more effectively even from the long payload features. The generated Bigrams and Trigrams can also be handled by any machine learning algorithms. In the present work, feature selection process was done by using one of the most widely cited machine learning algorithm, Support Vector Machine (SVM) [19] to study the results out of the proposed methodology. Weka GUI v3.8 has been used for experiments [20].

The next step is the feature vector extraction for the long payload features which is outlined in the following figure 3. The approach in [16] has been followed in the present



work for the task of feature vector extraction. The feature vector extraction step is also explained in more detail in the below Algorithm 2.

Algorithm 2: Feature Vector Extraction for Long Payload Features with Bigrams & Trigrams Step-1: Input the long payload features and constructed Bigram / Trigram Dictionary Step-2: Initialize all feature vectors to zero Step-3: Take one string at a time Step-4: Take one Bigram / Trigram Step-5: Find the index of Bigram / Trigram from the dictionary Step-6: Increment the location counter in the feature vector Step-7: Repeat till there is no possibility of finding Bigram / Trigram from the input payload feature Step-8: Finish feature vector i, and proceed to feature vector i+1 Step-9: Stop the process if all there no feature vector is left to be processed Step-10: Output the feature vectors



Fig.3 Feature vector extraction for ISCX data set using bigram technique

IV. EXPERIMENTS & RESULTS

The experiments were carried out on the available ISCX 2012 data set. The result of Algorithm 2 is the feature vectors for all the payload string features from this ISCX 2012 data set. After this step, the number of features can expand from hundreds to thousands and even to millions. This actually motivated to propose a new scheme of encoding, called 'Trigram', which somehow reduces this count. To explain the proposed scheme more clearly, a demonstration is made with small example on how the payload features are converted to bigrams according to the proposed methodology. To avoid the complexity,only three training instances have been taken as in [16] but different payload features are taken. The content of the each payload feature is different.

Generation of Bigrams:

Suppose the actual three payload features are: "B7z2Pr", "Vu3dj" and "R7z2nB" respectively then by feeding those features to the dictionary generation according to Algorithm 2, the resulting dictionary will consist of twelve (12)bigram words as follow: B7 | 7z | z2 | 2P | Pr | Vu | u3 | 3d | dj | R7 | 2n | nB. The redundant Bigrams are excluded from the list.

By applying the feature vector extraction on the given three payload features according to Algorithm 2, the resulting bigram representation for each of the above three features will be as presented in Table 1.

Table 1 : Bigram representation for the three payload features in the example

Original payload	B7	7z	z2	2P	Pr	Vu	u3	3d	dj	R7	2n	nB
B7z2Pr	1	1	1	0	0	0	0	0	0	0	0	0
Vu3dj	0	0	0	0	0	1	1	1	1	0	0	0
R7z2nB	0	1	1	0	0	0	0	0	0	1	1	1

Generation of Trigrams:

Here, the three payload features that are taken as example are same as taken for Bigram generation. The resulting trigram dictionary consists ten (10) Trigram as follows:



B7z | 7z2 | z2P | 2Pr | Vu3 | u3d | 3dj | R7z | z2n |2nB. The redundant Trigrams are excluded from the list.

By applying the feature vector extraction on the given three payload features according to Algorithm 2, the resulting Trigram representation for each of the above three features will be as depicted in Table 2

Table 2: Trigram representation for the three payload features in the example

Original payload	B7z	7z2	z2P	2Pr	Vu3	u3d	3dj	R7z	z2n	2nB
B7z2Pr	1	1	1	1	0	0	0	0	0	0
Vu3dj	0	0	0	0	1	1	1	0	0	0
R7z2nB	0	1	0	0	0	0	0	1	1	1

The three payload features in this example were converted to bigram features and Trigram features as in Table 1 and table 2 respectively. The table's header represents the standard feature vector for all the payload features in the taken example. The payload features appear in the table as per the order of their presentation to the feature extraction algorithm, that is the first row corresponds to the first payload feature discovered during dictionary generation and second row corresponds to the second payload feature discovered during dictionary generation and so on.

After all this preprocessing, in order to prepare the resulting data set for feature selection, a pre-ranking step has been conducted for the features using a quick filter selection algorithm that is the well known Correlation-based Feature Selection (CFS). The reason for this ranking is to make sure that the reducedfeatures are still informative andthus far, still contained distracting features as well. Hence a very fast feature selection method is used to pre-rank the features, which is called Correlation-based Feature Selection (CFS). By using this method, it is easy to rank the features in order to perform controlled experiments by manipulating relevant and irrelevant features.

Because the current number of features is very huge and the feature selection is really a time consuming activity, hence a small subset of 250 features from the original features are taken for experimental purpose. From the obtained ranked features list from the CFS method, a smaller subset of this ranked list is extracted in such a way that it consists of one portion of the top ranked features and nine portions of lowest ranked features. The considered portion size in this work is 25 features; therefore the subset size is obviously 250 features in total. This way of selection is intentional because to test the efficacy of the proposed system, experiments were conducted on the data set, in which 90% of the total features are bad features and only 10% are good features. The small subset of features is taken as sample from the huge set of total features to save time and computational effort.

Those resulting 250 features are used to generate data sets of size 20, 40, 100, and 400 examples respectively. In order to simulate "zero-day" attacks, the data sets have been chosen to be small in terms of number of examples. Those data sets were generated as balanced data sets (i.e., equal number of normal and attack examples). Different numbers of examples were used to monitor the behavior of feature selection with each size of data set.

To observe the effect of including the payload features to improve the detection accuracy, a very important experiment has been conducted. The well known machine learning algorithm, the SVM's classification was used to find out the accuracy and F-measure on the ISCX 2012 data sets in two cases, first without (Bigram Features / Trigram features) the payload features and second with (Bigram / Trigram features) payload features. The performance metrics have been measured before and after converting the payload features into bigram and trigram features and applying feature selection.

The main objective of this experiment is to show that these payload features include important and useful information in improving the detection accuracy. Because of the inability of handling the long payload futures, many of the previous researchers excluded payload features from the original set of features before looking for intrusions. The feature selection method was applied on the four generated data sets of ISCX 2012 and presented the maximum obtained accuracy and F-measure as shown in Table 3 and Table 4 respectively.



Data Set	Withou fea	t Bigram tures	With Bigram features			
	ACC	F-measure	ACC	F-measure		
20 examples	66.20%	66.20%	77.60%	77.00%		
40 examples	78.00%	78.00%	88.60%	88.60%		
100 examples	78.50%	78.50%	88.80%	88.70%		
400 examples	82.89%	82.89%	92.90%	92.90%		

 Table 3: Performance metrics without and with bigram
 features on the ISCX data sets

 Table 4: Performance metrics without and with Trigram
 features on the ISCX data sets

Data Set	Without feat	Trigram ures	With Trigram features			
	ACC	F- measure	ACC	F- measure		
20 examples	66.20%	66.20%	76.60%	76.00%		
40 examples	78.00%	78.00%	87.70%	88.21%		
100 examples	78.50%	78.50%	88.20%	88.20%		
400 examples	82.89%	82.89%	92.01%	91.90%		

It may be clearly observed from the above table that there is a clear improvement in the accuracy and F- Measure while considering the Bigram and Trigram features. The classifier's performance on each data set with the payload features removed and included from that data set has been studied. Columns 4 and 5 in the tables represent the maximum obtained performance from the classifier after including the payload features, expanding them using the bigram / trigram technique and applying feature selection algorithm on the resulting data set. It may be noted that there is a considerable improvement around 10% after applying classifier and feature selection algorithm on the bigram / trigram features compared to the performance of the SVMs classifier on the same data sets without bigram / trigram features.

One of the major contributions in this work is comparative analysis between the Bigram and Trigram on the same data set of long payload features. It may also be concluded after the studies on the experimental results that both Bigram and Trigram features inclusion have improved the performance of the classifier in terms of accuracy and False Measure. Even though there is a minute improvement using bigram features over trigram features in some instances, there will be a considerable reduction in the computational overhead. The below two graphs, Graph1, Graph2 will give a clear picture on the almost similar results with Bigram and Trigram with respect to Accuracy and F-measure as performance metrics.



Graph1: Accuracy achieved over Bigrams & Trigrams on the ISCX data sets



Graph2: F-measure achieved over Bigrams & Trigrams on the ISCX data sets



Hence after various observations from the experimental results it is suggested to use Trigram features to include large payload features instead of Bigram in any intelligent feature selection system to reduce the computational overhead.

V. CONCLUSION

In this paper, to handle long payload features, two encoding schemes, where one is available in the literature and other is a new one, have been studied and experiments were carried on the proposed methods on the ISCX 2012 data set. The data set has been prepared for intrusion detection by encoding the payload features which are in long strings using bigram and trigram technique. While most of the previous approaches in performing feature selection have focused only on the statistical properties of packets, an attempt was made to also try and extract useful features from the long payloads using a bigram and trigram techniques. Applying these techniques produced a high-dimensional data set with thousands of features from the given payload features. Motivated by the need to identify new threats such as zero-day attacks, and also to address the need to deal with large payload features, various experiments were carried out on large number of data sets. It is observed that the data set with large numbers of features and relatively few numbers of examples is always useful in testing the resilience of the system against over fitting. The experimental results are encouraging in drawing useful conclusions such as to recommend Trigram features to include large payload features instead of Bigram in any intelligent feature selection system to reduce the computational burden.

ACKNOWLEDGEMENTS:

Authors want to express their special thanks to Dr. Arash Habibi Lashkari, Research Associate R&D Manager, Canadian Institute for Cyber security (CIC) and other team members at CIC for the generosity in responding to the request of the authors for various data sets including ISCX 2012 [18]. The data sets provided by CIC are very useful in this current work in carrying out various experiments and draw useful conclusions.

REFERENCES

[1] Shiravi A, Shiravi H, Tavallaee M, Ghorbani AA. Toward developing a systematic approach to generate benchmark datasets for intrusion detection. ComputSecur 2012;31(3):357–74 doi. Available from: http://dx.doi.org/10.1016/j.cose.2011.12.012, http: //www. science direct.com/science/article/ pii/ S01674 048110 01672.

[2] Garcia LP, de Carvalho AC, Lorena AC. Effect of label noise in the complexity of classification problems. Neurocomputing 2015;160:108–19. Available from: http://dx.doi.org/10.1016/j.neucom.2014.10.085, http:// www.sciencedirect.com/science/article/pii/S09252312150 01241.

[3] Beigi EB, Jazi HH, Stakhanova N, Ghorbani AA. Towards effective feature selection in machine learningbased botnet detection approaches, in: 2014 IEEE Conference on Communications and Network Security; 2014, pp. 247–255. doi:10.1109/ CNS.2014.6997492.

[4] Bolon-Canedo V, Snchez-Maroo N, Alonso-Betanzos A, Bentez J, Herrera F. A review of microarray datasets and applied feature selection methods. Inf Sci (Ny) 2014;282:111–35. Available from: http:// dx. doi.org /10.1016 /j.ins.2014.05.042, http:// www. sciencedirect .com/science /article/pii/ S0020025514006021.

[5] Beniwal S, Arora J. Classification and feature selection techniques in data mining. Int J Eng Res Technol 2012;1(6):1-6.

[6] Fahad A, Tari Z, Khalil I, Habib I, Alnuweiri H. Toward an efficient and scalable feature selection approach for internet traffic classification. ComputNetw 2013;57(9):2040–57. Available from: http:// dx.doi. org/10. 1016/j.comnet. 2013.04.005, http://www.science direct.com/science/article/pii/ \$1389128613001163.

[7] Aghdam MH, Kabiri P. Feature selection for intrusion detection system using ant colony optimization. IJ NetwSecur 2016;18(3):420–32.

[8]Bolón-Canedo V, Sánchez-Maroño N, Alonso-Betanzos A. A review of feature selection methods on synthetic data. KnowlInfSyst 2013;34(3):483–519. doi:10.1007/s10115-012-0487-8.

[9] Sahu SK, Sarangi S, Jena SK. A detail analysis on intrusion detection datasets, in: 2014 IEEE International Advance Computing Conference (IACC), 2014, pp. 1348–1353. doi:10.1109/IAdCC.2014.6779523.



[10]Mancini LV, Di Pietro R. Intrusion Detection Systems. Springer; 2008

[11] Ambusaidi MA, He X, Nanda P, Tan Z. Building an intrusion detection system using a filter-based feature algorithm. IEEE selection Trans Comput 2016;65(10):2986-98. doi:10.1109/TC.2016.2519914.

[12]Abou El Kalam A., Gad El Rab M., and Deswarte Y. (2014), A model-driven approach for experimental evaluation of intrusion detection systems, Security Comm. Networks, 7, 1955–1973, doi: pages 10.1002/sec.911

[13]Mell P, Hu V, Lipmann R, Haines J, Zissman M. An overview of issues in testing intrusion detection systems. Technical Report, NIST IR 7007, National Institute of Standard and Technology, USA, 2003.

[14] NiccolòCascarano, Luigi Ciminiera, and FulvioRisso. 2010. Improving cost and accuracy of DPI traffic classifiers. In Proceedings of the 2010 ACM Symposium on Applied Computing (SAC '10). ACM, New York, NY, USA,641-646. DOI= http:// dx.doi. org/10. 1145 /1774088.1774223

eers....deretoping research [15]Laurent Bernaille, Renata Teixeira, Ismael Akodkenou, Augustin Soule, KaveSalamatian, Traffic classification on the fly, ACM SIGCOMM Computer Communication Review, v.36 n.2, April 2006 [doi>10.1145/1129582.1129589]

[16] TarfaHamed, Rozita Dara, Stefan C. Kremer, Network intrusion detection system based on recursive feature addition and bigram technique, computers & s e c u r i t y 73 (2018) 137–155

[17] Zhang M, Wang L, Jajodia S, Singhal A, Albanese M. Networkdiversity: a security metric for evaluating the resilience of networks against zero-day attacks. IEEE Trans Inf ForensicsSec 2016:11(5):1071-86. doi:10.1109/TIFS.2016.2516916.

[18] Intrusion detection evaluation dataset (ISCXIDS2012), http:// unb.c a/cic/res earch /data sets /index.html

[19] Chang C-C, Lin C-J. LIBSVM: a library for support vector machines. ACM Trans Intelligent Syst Technol 2011;2(3):27.Debar H. An introduction to intrusiondetection systems, in: Proceedings of Connect2000; 2002, pp. 1–18.

[20] Open Source, freely available and downloadable from: https: //www .cs. Waikato .ac. nz/ ml/ weka /down loading.html