

ISSUES RELATED TO SAMPLING TECHNIQUES FOR NETWORK TRAFFIC DATASET

Raman Singh¹, Harish Kumar² and R.K. Singla³

^{1,2}University Institute of Engineering and Technology, Panjab University, Chandigarh
¹raman.uiet84@gmail.com, ²harishk@pu.ac.in

³Department of Computer Science and Application, Panjab University, Chandigarh
³rksingla@pu.ac.in

ABSTRACT

Network traffic data is huge, varying and imbalanced because various classes are not equally distributed. Machine learning (ML) algorithms for traffic analysis uses the samples from this data to recommend the actions to be taken by the network administrators. Due to imbalances in dataset, machine learning algorithms may give biased or false results leading to serious degradation in performance of these algorithms. Since the network dataset is huge, during training machine learning algorithm takes more time and hence sampling should be used to reduce the training time. But using sampling may cause loss of information which should be taken care off while obtaining the samples. In this paper various sampling techniques have been analysed for loss of information and imbalances during sampling of network traffic data. Data set is collected from the Panjab University network. Various parameters like missing classes in samples, probability of sampling of the different instances have been considered for comparison.

KEYWORDS

Imbalanced learning, Sampling, Re-sampling, machine learning.

1. INTRODUCTION

Communication between different networked computers is growing day by day and so the threats to these networks. Preventive actions against any of the network anomaly can be taken by analysing this data. But to process whole of the data is not feasible due to ever growing network traffic. Hence, newer techniques like Machine learning (ML) etc. are used for the analysis. These techniques need traffic samples to analyse threats. Network traffic data is imbalanced and is performance degrader factor. Learning from this dataset may be biased because of unequal distribution of instances. Most of standard ML algorithms assume that instances are equally distributed among classes and hence the outcomes may be biased and fails to properly represent statistical properties of it. Techniques to minimize imbalances in this data are required. Rest of the paper is divided as follows. Section 2 describes experiment setup used to capture network traffic, Section 3 describes imbalanced network traffic dataset and its issues; Section 4 describes sampling of dataset, Section 5 discusses the results and Section 6 describes conclusions and future scope.

2. EXPERIMENT SETUP

There are various datasets available for the analysis. These datasets are prepared capturing traffic from different networks. In this research work, a dataset named Panjab University Campus-Wide

Area Network (PU-CAN) is prepared by capturing the traffic from Panjab University’s Campus Wide Network. To capture this dataset, a network server is configured in PU-CAN. Sub-network of PU-CAN which is used to capture dataset provides network service to three boys’ hostels covering approximately 500 users. This network is managed by team of administrators and network engineers. Internet facility is providing through Squid Server to ensure controlled access and surveillance on user’s internet usage. Proper firewall system is configured to stop malwares and unauthorized access and attacks. Dynamic Host Configuration Protocol (DHCP) server is used to dynamically distribute Internet Protocol (IP) address. The official IP addresses range which is assigned to users is 172.16.40.1 to 172.16.43.254. Some IP addresses are kept reserved for servers like 172.16.40.1 for Domain Name Server (DNS), 172.16.40.2 for Squid Server and DHCP Server. The IP address of server which is used to capture network traffic dataset is 172.16.40.11. Figure 1 shows the network diagram of sub-network of PU-CAN.

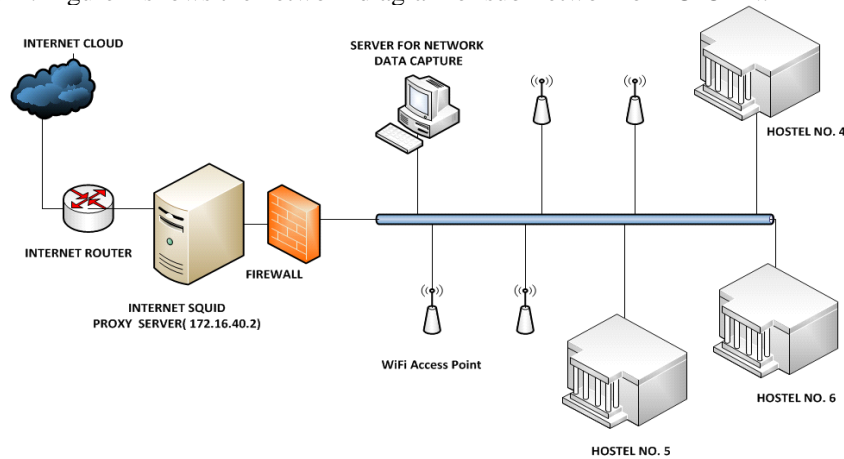


Figure 1. Network Diagram of PU-CAN Used to Capture Traffic Data

3. IMBALANCES IN NETWORK TRAFFIC DATASETS

3.1. Imbalanced Dataset

On any computer network, thousands of packets travel at any given time. Capturing each packets and then analysing network traffic dataset in order to detect malwares is very time consuming job for any Intrusion Detection System (IDS). In network traffic data the instances of malicious packets like malwares, attacks, viruses are very few in number than instances of normal packets. This problem is known as imbalanced dataset problem. This leads to serious problem of under training or biased training of the models based on machine learning and subsequently leads to miss-classification. The captured dataset is huge and need to be pre-processed efficiently before any machine learning algorithm is used to analyse it. Some classes may have hundreds and thousands of instances whereas other may have only very few number of instances [1]. In network traffic dataset, one class of packets is present in large numbers while other class has only few instances [2]. There can be un-authorized sub-networks existing in the network which can lead to performance degradation [3]. Classifier gives better performance if dataset is balanced while due to imbalanced nature of network traffic dataset, classification under performs [4].

3.2. Issues in Imbalanced Dataset

Various issues due to imbalanced nature of network traffic dataset are:

- i. ML algorithms used in network traffic classification to detect malware take it granted that dataset is balanced in nature and instances as well as network packets are equally distributed in all classes, but this is not true in case of network traffic data and results into biased classification [1].
- ii. ML technique gives poor performance on minority classes because distribution of training data may differ from testing data but these techniques are generalized and assume that both dataset have same distribution [5].
- iii. ML algorithms are designed to obtain higher accuracy rate, but this may lead to problem with IDS to detect minority attack patterns. These techniques are guided by standard accuracy rate and are biased towards majority instances, while on the other hand it is difficult to distinguish between noise and minority instances as both are few [6].

4. SAMPLING OF NETWORK TRAFFIC DATASET

Network traffic dataset requires pre-processing steps before machine learning algorithms applied to detect malwares. Pre-processing removes noisy or missing data. Huge network traffic dataset is sampled to increase performance of these algorithms. The various steps involved in pre-processing of network traffic dataset are discussed below:

a) Dataset Generation: The first step is to generate dataset for network traffic classification. Dataset is divided into training set and test set [7]. Since network traffic dataset is huge, so various sampling techniques are applied to limit size. However characteristics of dataset should not change while using sampling.

b) Feature Selection and Extraction: Network traffic dataset has various features but all those features may not contribute in classification and intrusion detection. In order to enhance accuracy and performance, important features need to be selected ignoring other redundant/irrelevant features. New features can also be derived from existing features to increase performance [7].

Commonly used sampling techniques used to sample network traffic dataset are:

a) Random Sampling: is random method to select 'p' (small letter p) instances from population 'P' (Capital letter P) packets/instances of network traffic dataset [8]. It can be done 'with/without' replacement. Probability of sampling P(s) and size of sampled dataset in percentage of total dataset can be calculated as in equation (1) and (2) below:

$$P(s) = \text{No. of favored events} / \text{Total no. of events} = p/P \quad (1)$$

$$\text{Size of sampled dataset (\%age of total dataset)} = p/P * 100 \quad (2)$$

b) Systematic Sampling: selects 'p' packets out of total population of 'P' packets by considering a packet after every regular interval starting from a point. For example if dataset size is 10000 packets and 1000 packets need to be selected then every 10th packets should be selected starting from first instance [9]. Sampling interval can be calculated as in equation (3) below:

$$\text{Sampling Interval} = \text{Total Population} / \text{Packets required} = P/p \quad (3)$$

$$\text{Starting point} = \text{1st packet}$$

c) Stratified Sampling: is capable of discovery of statistical characteristics of network traffic dataset. It is used in heterogeneous dataset where all instances are not of the same type. First, population is divided into heterogeneous sub-population called strata. Sub-population in each strata is homogenous. Then samples are selected from each strata. While selecting the samples, random/systematic/proportional-to-size sampling can be used [10]. Prior knowledge about characteristics of populations required in stratified sampling. Population 'P' is divided into 'G' groups, and each group has some different numbers of instances. Then from each group some samples are selected. Total sampled instances n(s) can be calculated by adding sampled instances from each strata as shown in equation (4).

$$n(s) = \sum_{i=0}^I P_i, \quad (4)$$

Probability of sampling for each strata can be calculated by equation (5):

$$P(s_i) = p_i/P, \text{ where } I = 1 \text{ to } G \quad (5)$$

The size of sampled dataset can be calculated by below given equation (6):

$$\text{Size of sampled dataset(In \% age)} = \left(\sum_{i=0}^I P_i / P \right), *100 \quad (6)$$

d) Re-sampling: Classified into Under-sampling and Over-sampling. In Under-sampling instances are sampled in fewer rates so that minority and majority classes have equal contribution in sampling method [11]. It is used to balance imbalances in dataset by elimination of instances of majority classes. The drawback is that sometime the loss of information may occur if some particular instances are missed and loss of useful information may occur [12][13]. The instances can be under sampled by some factor. Prior knowledge about dataset is required for deciding this factor. In over-sampling instances of minority classes are synthetically generated up to some pre-defined numbers [12]. Instances are sampled with higher sampling rate. More samples are picked in over-sampling from minority class and few picked in under-sampling from majority class. Drawback is that since the instances are generated synthetically/repeated, redundant information may leads to biased classification. Replicated instances may leads to wrong decisions [13]. Also since network traffic dataset is huge and imbalanced and further if minority classes are large in number then performing over-sampling to reduce imbalances and replicating minority classes instances further increase size of dataset and computational time [12]. SMOTE is suggested to synthetically generate the samples during under sampling and over sampling process. It also gives better performance in terms of ROC space than many other classifiers [14].

5. RESULTS AND DISCUSSION

5.1. Network Traffic Dataset Characteristics

For analysis of sampling technique on network traffic data, packets are captured and a dataset is prepared. Size of this dataset is approximately 16 GigaByte comprising of billions of instances. Out of this huge data, a small set of 60000 instances of dataset is taken for testing. Java programming language [15] is used to implement various sampling methods and to analyse results. Dataset used for analysis purpose is named as “Panjab University-Test Data Set” (PU-TDS). It covers 60000 instances of 25 different protocols. First analysis comprising of number of packets, their percentage and probability with-in sample has been carried out and the results are shown in table 1. Number of instances for each packets associated with different protocol has been calculated. Percentage of packets is also shown. P(s) is probability of packets to be picked in sampling process. Higher the P(s) value, higher the chances of packets to be picked for sampling. The analysis of test dataset shows that if this dataset is to be classified as per protocols there will be 25 classes associated with each protocol. The probability of packets for each protocol is varying proportional to numbers of packets for these protocols. Some protocols like ICMPv6 are present in higher number (12542 instances), while others like HTTP/XML, OCSP, SSL and IAPP, have very few presences. Hence this can be derived that network traffic dataset is imbalanced. This nature of data can mislead the machine learning algorithm, biased classification

and miss-classification may occur. Due to huge size of data set, pre-processing steps like sampling are required before feeding it to machine learning algorithms. But probability of certain protocol packets like HTTP/XML is minute (0.001666667), hence their chances to enter into sample are also rare. It leads to loss of information due to sampling process. Also if machine learning algorithm is used to find malware in intrusion detection system special care should be taken in pre-processing and training as malware instances are very less in number while normal packets may outnumber malicious packets. Due to these characteristics and issues specialize pre-processing/sampling and machine learning techniques should be designed for network traffic classification to detect intrusions/ attacks.

Table 1. Analysis of PU-TDS for probability of sampling and size for each protocol (classes)

Protocols	No. of Packets	%age of Packets	P(s)
DHCP	3335	5.558333333	0.005558
ARP	4501	12.50167	0.125017
ICMP	141	0.235	0.00235
HTTP	1149	1.915	0.01915
TCP	8594	14.32333	0.143233
UDP	2740	4.566667	0.045667
ICMPv6	12542	20.90333	0.209033
SSDP	5566	9.276667	0.092767
NBNS	6466	10.77667	0.107767
MDNS	158	0.263333	0.002633
LLMNR	6836	11.39333	0.113933
BROWSER	710	1.183333	0.011833
TLSv1	1788	2.98	0.0298
DB-LSP-DISC	5	0.008333	0.000083
DHCPv6	969	1.615	0.01615
DNS	351	0.585	0.00585
HTTP/XML	1	0.001667	0.00001667
IAPP	3	0.005	0.00005
IGMP	678	1.13	0.0113
IPX RIP	22	0.036667	0.000367
LLC	344	0.573333	0.005733
NBIPX	37	0.061667	0.000617
OCSP	2	0.003333	0.000033
SSL	2	0.003333	0.000033
XID	60	0.1	0.001

5.2. Analysis of various sampling techniques using PU-TDS

In order to analyse effect of various sampling techniques on network traffic dataset, various commonly used sampling methods are implemented in Java programming language and tested using PU-TDS.

5.2.1. Random Sampling

Packets are randomly selected from PU-TDS. Consider 'p' is the number of packets to be selected for sampling. Experiment is done for different value of 'p' like 500, 1000, 2000, 5000, 10000,

15000,20000 and 40000. Table 2 shows the percentage of packets selected out of total packets of different protocols for various values of 'p'.

As the network traffic dataset is imbalanced and heterogeneous, the number of packets selected is varying for different protocols. For some protocols packets the percentage of packets selected are as high as 23.40 while some packets of protocols class are missed by this sampling. It causes loss of information and machine learning will not get proper training on these sampled datasets. Misclassification may occur. In intrusion detection technique this missing samples may cause harm to network as malicious packets which are fewer in number may not be selected in sampling and Intrusion Detection System (IDS) system may fail to detect attacks. So, this sampling should not be used in sampling of network traffic data.

Table 2. Percentage of packets selected for various protocol in random sampling

p	500	1000	2000	5000	10000	15000	20000	40000
Protocol								
DHCP	5.80	4.60	6.25	5.56	5.80	5.51	5.59	5.60
ARP	13.00	11.50	12.55	12.42	12.48	12.45	12.36	12.33
ICMP	0.40	0.30	0.25	0.26	0.19	0.24	0.22	0.20
HTTP	2.40	1.20	1.55	1.80	2.00	1.97	1.88	1.91
TCP	17.60	12.90	14.95	14.50	13.38	13.86	14.30	14.31
UDP	5.80	4.50	4.30	4.46	4.66	4.75	4.52	4.62
ICMPv6	18.40	23.40	20.60	20.44	21.10	20.95	21.33	20.88
SSDP	8.00	9.80	9.05	8.94	9.12	9.40	9.20	9.30
NBNS	11.20	11.60	10.45	10.66	10.91	10.75	10.53	10.80
MDNS	0.20	0.10	0.05	0.36	0.28	0.28	0.28	0.28
LLMNR	9.60	11.50	11.20	11.64	11.79	11.59	11.56	11.41
BROWSER	0.80	0.70	0.80	1.12	1.20	1.19	1.20	1.19
TLSv1	1.80	3.50	3.60	3.36	2.97	2.93	2.95	2.90
DB-LSP-	0.00	0.00	0.00	0.00	0.00	0.01	0.02	0.01
DHCPv6	1.40	1.60	1.55	1.56	1.49	1.67	1.53	1.75
DNS	1.20	0.50	0.65	0.80	0.60	0.54	0.75	0.60
HTTP/XML	0.00	0.00	0.00	0.02	0.00	0.00	0.01	0.00
IAPP	0.00	0.00	0.00	0.02	0.00	0.00	0.01	0.00
IGMP	2.20	0.80	1.25	1.28	1.10	1.13	1.08	1.11
IPX RIP	0.00	0.10	0.10	0.08	0.03	0.03	0.03	0.04
LLC	0.20	1.30	0.75	0.58	0.69	0.58	0.56	0.61
NBIPX	0.00	0.10	0.05	0.08	0.08	0.05	0.04	0.06
OCSP	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.01
SSL	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00
XID	0.00	0.00	0.05	0.06	0.11	0.10	0.09	0.10

Figure 2 shows the number of classes missed by random sampling as per number of packets selected for sampling. From figure-2 it is clear that with increase in sampling factor, there may be decrease in number of missing classes. Since packets are randomly selected, it is not necessary that decrease should be uniform. Due to this reason, there are certain exception at $p=10000$ & $p=15000$. If sampling factor is 500 (i.e $p=500$), then 8 classes out of total 25 classes are missed in sampling and will not contribute in decision making. Further if we increase sampling factor loss of information may decrease but cannot be ruled out completely.

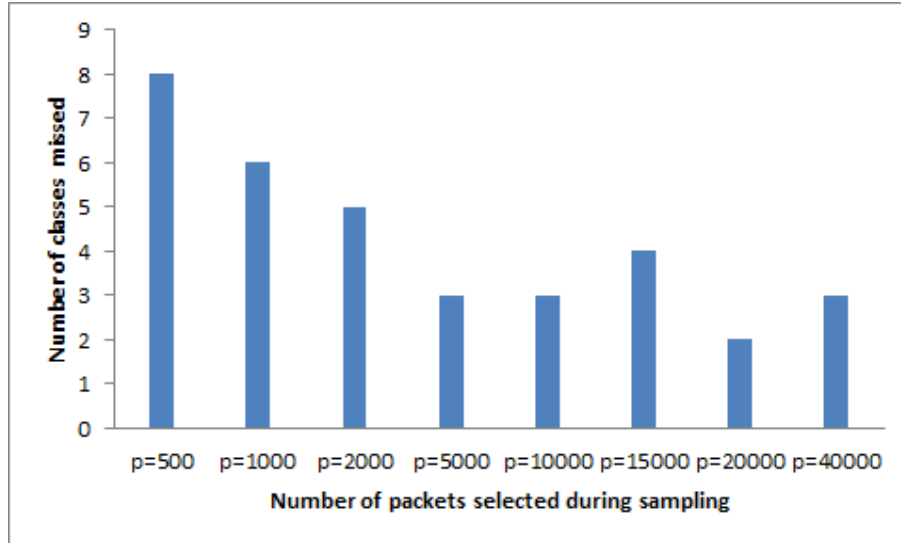


Figure 2. Loss of information in random sampling

5.2.2. Systematic Sampling

Packets are systematically selected from PU-TDS. If every I^{th} packet is selected for sampling, then 'I' is known as sampling factor. As the value of 'I' increases, number of packets selected decreases. In this experiment each 3rd, 5th, 10th, 15th, 20th, 25th and 30th packet are selected and 7 different sampled dataset are prepared. Table 3 shows the value of 'I', total no. of packets selected for each class and their percentage.

Table 3. Percentage of packets selected for various protocol in random sampling

Protocols	I=3, p=20000	I=5, p=12000	I=10, p=6000	I = 15, p=4000	I = 20, p=3000	I = 25, p=2400	I = 30, p=2000
DHCP	5.58	5.62	5.73	5.60	6.03	5.21	5.95
ARP	12.43	12.38	12.68	12.45	11.93	12.63	12.80
ICMP	0.23	0.28	0.25	0.25	0.20	0.17	0.25
HTTP	1.85	1.99	1.95	2.10	2.23	2.17	2.15
TCP	14.48	14.15	14.13	14.48	14.83	14.17	14.25
UDP	4.63	4.28	4.03	4.18	3.80	4.17	4.00
ICMPv6	21.10	21.03	20.87	20.65	20.57	21.04	19.85
SSDP	9.27	9.19	9.03	9.43	9.40	9.42	9.15
NBNS	10.82	11.32	11.65	11.38	11.40	11.38	12.10
MDNS	0.27	0.24	0.23	0.30	0.27	0.17	0.20
LLMNR	11.19	11.02	10.92	10.93	10.97	10.96	10.90
BROWSER	1.09	1.24	1.38	0.98	1.33	1.38	1.20
TLSv1	2.91	2.93	2.70	2.78	2.63	2.83	2.85
DB-LSP-DISC	0.01	0.00	0.00	0.00	0.00	0.00	0.00
DHCPv6	1.66	1.86	1.93	2.10	1.83	1.88	2.05
DNS	0.56	0.53	0.52	0.48	0.47	0.63	0.50
HTTP/XML	0.01	0.00	0.00	0.00	0.00	0.00	0.00

IAPP	0.02	0.01	0.00	0.03	0.00	0.00	0.00
IGMP	1.14	1.18	1.20	1.23	1.33	0.96	1.20
IPX RIP	0.05	0.03	0.02	0.00	0.00	0.00	0.00
LLC	0.57	0.54	0.60	0.40	0.60	0.54	0.35
NBIPX	0.06	0.05	0.07	0.05	0.07	0.13	0.05
OCSP	0.01	0.01	0.02	0.03	0.03	0.00	0.05
SSL	0.00	0.00	0.00	0.00	0.00	0.00	0.00
XID	0.13	0.15	0.08	0.23	0.07	0.21	0.15

From table 3, it can be analysed that some classes are present in majority while others are in minority. Experiment shows that this may miss some classes and can cause loss of information. Figure 3 shows analysis of loss of information in systematic sampling as per different sampling factor values. Number of classes missed is also shown.

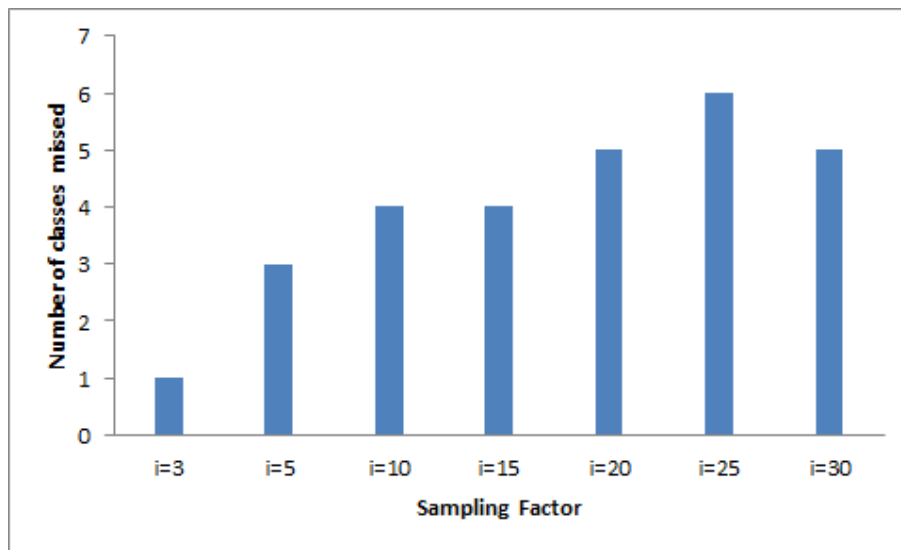


Figure 3. Loss of information in systematic sampling

5.2.3. Strata Sampling

In this sampling different heterogeneous stratas are defined which further contains homogeneous packets. In experiment, 25 different classes are considered depending on protocols. 25 different stratas are defined for each protocol or class. Then for each stratas/classes, 7 experiment performed using systematics sampling and taking values of sampling factor (I) as 4, 5, 6, 7, 8, 9 and 10. The advantage of this sampling is that each of heterogeneous packets will contribute in decision making. Out of each stratas, systematics sampling is used to select packets since some stratas have large number of packets. It is also known as two phase sampling. First each packets is divided into different 25 stratas based on protocol used and then out of each stratas some packets are selected for creating of final sampled dataset. Table 4, shows the results of strata sampling.

Table 4. Percentage of packets selected for various protocol in strata sampling

Protocols	I = 4, p=15011	I = 5, p=12011	I = 6, p=10000	I = 7, p=8586	I = 8, p=7511	I = 9, p=6679	I = 10, p=6012
DHCP	5.56	5.55	5.73	5.56	5.55	5.55	5.56
ARP	12.50	12.50	12.52	12.49	12.49	12.49	12.49
ICMP	0.24	0.24	0.14	0.24	0.24	0.24	0.25
HTTP	1.92	1.91	1.89	1.92	1.92	1.92	1.91
TCP	14.32	14.31	14.50	14.30	14.31	14.30	14.30
UDP	4.56	4.56	4.47	4.57	4.57	4.57	4.56
ICMPv6	20.89	20.89	21.07	20.87	20.88	20.87	20.87
SSDP	9.27	9.27	9.12	9.27	9.27	9.27	9.26
NBNS	10.77	10.77	11.09	10.76	10.77	10.77	10.76
MDNS	0.27	0.27	0.23	0.27	0.27	0.27	0.27
LLMNR	11.38	11.39	10.90	11.38	11.38	11.38	11.38
BROWSER	1.19	1.18	1.21	1.19	1.18	1.18	1.18
TLSv1	2.98	2.98	2.88	2.98	2.98	2.98	2.98
DB-LSP-DISC	0.01	0.01	0.01	0.01	0.01	0.01	0.02
DHCPv6	1.62	1.62	1.71	1.62	1.62	1.62	1.61
DNS	0.59	0.59	0.62	0.59	0.59	0.58	0.60
HTTP/XML	0.01	0.01	0.00	0.01	0.01	0.01	0.02
IAPP	0.01	0.01	0.02	0.01	0.01	0.01	0.02
IGMP	1.13	1.13	1.03	1.13	1.13	1.14	1.13
IPX RIP	0.04	0.04	0.06	0.05	0.04	0.04	0.05
LLC	0.57	0.57	0.60	0.58	0.57	0.58	0.58
NBIPX	0.07	0.07	0.06	0.07	0.07	0.07	0.07
OCSP	0.01	0.01	0.01	0.01	0.01	0.01	0.02
SSL	0.01	0.01	0.01	0.01	0.01	0.01	0.02
XID	0.10	0.10	0.12	0.10	0.11	0.10	0.10

From table 4, it can be analysed that since each heterogeneous packet is considered for creating strata, loss of information is minimal. However, since each stratas have homogeneous packets and number of packets in one stratas may outnumber packets in others, So imbalances may present.

5.2.4. Re-Sampling

It is used to balance the ratio of majority and minority classes. Over-sampling and under-sampling is used to correct imbalance of majority & minority class. If packets/instances in one class are very high then up to a fixed number of packets is selected. Otherwise synthetically generate/copied packets are selected. Heterogeneous packets are divided into distinguished stratas and then from each strata 't' number of packets are selected randomly. If number of packets in strata/class is greater than 't', then randomly 't' packets are selected. If it is less than 't', then 't' packets are randomly copied up to 't' numbers. Total number of packets selected is 't' multiply by number of stratas. Values of 't' taken are 100,200,300,400,500, and 700. Total number of selected packets can be calculated as in equation (7), where 's' is number of stratas.

$$\text{Total number of packets selected 'p'} = t*s \quad (7)$$

From experiments performed it is analysed that imbalances of data can be corrected by increasing instances of minority classes and decreasing that of majority classes. Then each class will have fair and equal contribution in decision making. In this experiment, since total numbers of classes are 25, during sampling 4 per-cent instances are selected from each class. In IDS, it can be used to detect malwares as these packets are few. If instance of normal packets and malicious packets made balanced then intrusions and attacks on networks can be detected easily by using ML algorithms. Disadvantage is that since imbalances are corrected in network traffic dataset by reducing majority classes and increasing minority classes, sometimes ML algorithms get wrong training by considering synthetically created minority instances as repeated attack patterns. This should be handled properly by suggesting new and improved ML techniques specifically for IDS.

6. CONCLUSIONS AND FUTURE SCOPE

Issues of imbalances in network traffic dataset are discussed in this paper. Since this dataset is very large, sampling should be used to increase performance before machine learning algorithm should be employed on dataset for intrusion detection system. Also if sampling is performed on imbalanced network traffic dataset loss of information may happen and IDS may give biased or abnormal results. Various commonly available sampling methods are discussed and experiments are performed to check worthiness of these sampling techniques. Dataset is collected using PUCAN and named as PU-TDS. In random and systematic sampling loss of information may occur which may cause wrong decision making, since malicious packets may be missed during sampling and IDS may fail to detect malwares and attacks. Strata sampling can be used to overcome problem of loss of information as in strata sampling minimum one packet of each heterogeneous packet will be selected but strata sampling does not overcome the imbalances of network traffic dataset. Re-sampling (Under-sampling & over-sampling) may be used to overcome problem of imbalances in network traffic dataset and each heterogeneous class including normal and malicious will have equal contribution of decision making and malware will be detected efficiently. Problem with re-sampling is that synthetically generated and repeated minority class packets may seem to be wrong patterns of attacks. To overcome this issue new and improved resampling approaches should be worked on for network traffic dataset to efficiently handle imbalances.

REFERENCES

- [1] Haibo He and Edwardo A. Garcia,(2009) "Learning From Imbalanced Data", IEEE Transaction On Knowledge And Data Engineering, Vol. 21, no. 9, pp 1263-1284.
- [2] Weicai Zhong, Bijan Raahemi and Jing Liu, (2009) "Learning on Class Imbalanced Data To Classify Peer-to-Peer Applications in IP Traffic using Resampling Techniques", Proceedings of International Joint Conference On Neural Networks,, June 14-19, Atlanta, Georgia, USA, pp 3548-3554.
- [3] Raman Singh, Harish Kumar and R.K. Singla, (2012)," Traffic Analysis of Campus Network for Classification of Broadcast Data", Proceedings of 47th Annual National Convention of Computer Society of India, International Conference on Intelligent Infrastructure, December 1-2, 2012, Science City, Kolkata, pp 163-166.
- [4] G. Weiss and F. Provost (2003), "Learning when training data are costly: the effect of class distribution on tree induction", Journal of Artificial Intelligence Research, Vol. 19, pp 315-354.
- [5] Mike Wasikowski and Xue-wen Chen,(2010) "Combating the Small Sample Class Imbalance Problem Using Feature Selection", IEEE Transactions on Knowledge and Data Engineering, Vol. 22, No. 10, pp 1388-1400.

- [6] Julián Luengo, Alberto Fernández, Salvador García and Francisco Herrera,(2011) “Addressing data complexity for imbalanced data sets: analysis of SMOTE-based oversampling and evolutionary undersampling”, *soft Computing* ,Vol. 15, No. 10, pp 1909-1936.
- [7] Jonathan J. Davis and Andrew J. Clark ,(2011) “Data preprocessing for anomaly based network intrusion detection: A review”, *Computers and Security, Computers & Security*, Vol. 30, No. 6–7, pp 353-375.
- [8] Guanghui He and Jennifer C. Hou , (2006) “On sampling self-similar Internet traffic”, *Computer Networks*, Vol. 50, No. 16, pp 2919-2936.
- [9] Abdun Naser Mahmood, Jiankun Hu, Zahir Tari and Christopher Leckie, (2010) “Critical infrastructure protection: Resource efficient sampling to improve detection of less frequent patterns in network traffic”, *Journal of Network and Computer Applications*, Vol. 33, No. 4, pp 491-502.
- [10] Stênio Fernandes, Carlos Kamienski, Judith Kelner, Dênio Mariz and Djamel Sadok, (2008) “ A stratified traffic sampling methodology for seeing the big picture”, *Computer Networks*, Vol. 52, No. 14, pp 2677-2689.
- [11] Liu, J. Ghosh, and C. Martin,(2007) "Generative oversampling for imbalanced datasets." *International Conference on Data Mining (DMIN)*, Las Vegas, Nevada, USA.
- [12] Sotiris Kotsiantis, Dimitris Kanellopoulos and Panayiotis Pintelas , (2006) “Handling imbalanced datasets: A review”, *GESTS International Transactions on Computer Science and Engineering*, Vol.30
- [13] Yang Liu, Xiaohui Yu, Jimmy Xiangji Huang and Aijun An, (2011) “Combining integrated sampling with SVM ensembles for learning from imbalanced datasets”, *Information Processing & Management*, Vol. 47, No. 4, pp 617-631.
- [14] Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall1 and W. Philip Kegelmeyer (2002), “MOTE: Synthetic Minority Over-sampling Technique “, *Journal of Artificial Intelligence Research*, Vol. 16, pp 321-357.
- [15] Java Programming Language, <http://www.oracle.com/technetwork/java/index.html>

Authors

Raman Singh is Research Scholar with UIET, Panjab University, Chandigarh. He completed his B.Tech. In Computer Engineering and then worked with Karman Infotech Pvt. Ltd. as Technology Specialist for two years. He completed ME (IT) from Panjab University Chandigarh. Now he is pursuing Ph.D. in Computer Science and Engineering from UIET, Panjab University.



Harish Kumar has done M.Tech. (CSE) from Department of CSEE, Punjab Agricultural University. He has been awarded Ph.D. in Computer Sci. by Panjab University. In 2002 he joined UIET in PU where he is working as Associate Prof. in Comp. Sci. & Engg. He has also worked as consultant with TCY, a leading online test portal having 10 million users prior to adopting academics as a profession . He has executed a sabbatical project at Infosys Tech. as a member of their development team in PED division. His areas of interest include MANET, Wireless Networks, Software Test Estimation, Open source software and e-learning. He has 30 publications to his credit. He has also delivered several lectures on various topics during workshops/ conference at national as well as international forums.



R.K. Singla has been Professor of Computer Sci. at Panjab University. He joined the department in 1985. He obtained the PhD in Scientific Computing-Electrohydrodynamics from Panjab Univ. and investigated the complex problem of analyzing the interaction of fluids and electric fields. His current research interests include Scientific Computing, Wireless Networks, Info. Security and Software Cost Estimation. He has around 30 Journals and 15 conference publications to his credit. He has been invited for talk at a large number of conferences/workshops. He is also a life-member of CSI. He held many IT-related positions including Chairman, DCSA; Coordinator, TIFAC-DST; Project Leader; and Programmer/Analyst. He has extensive experience with real-time integration of applications, IT administration, platform conversions (hardware and software), platforms (PC-DOS/Windows, VAX-VMS, SUN-Solaris, HP-UX and PC-Linux) and electronic communications.

