

SPID-based Method of Trojan Horse Detection

Xun-yi Ren Gui-bing Qian

College of Computer Nanjing University of Posts and Telecommunications
Email: renxy@njupt.edu.cn

Abstract

The Trojans have become a hidden threat to computer security, how to identify various Trojan efficiently and accurately is a current research focus. In this paper, based on the research of SPID's attribute meters, we proposed a method to detect and identify Trojan that based on the SPID feature optimization. The method uses SPID attribute meters to analysis with common protocol, generating a model to identify Trojan. Then, get the combination of 12 attribute meters to identify Trojan by statistical the result of the recognition. Experimental results shows that the optimized combination of attribute meters have a high efficiency to identify Trojan based on keeping SPID detection accuracy.

Keywords: Trojan horse, SPID, detection, identification

1. Introduction

Trojan [1] is a destructive program, and it is the main threat to the security of the network and host. Trojans communicate with the remote control host by pretending to be a useful or interesting program. Trojan has very large hazards, so the Trojan detection technology becomes the hotspot of the information security.

Nowadays, Trojan detection technology can be divided into several methods [2]: (1) through network

monitoring finding network traffic anomalies and then blocking it, or definite the rules that make the Trojans cannot communicate. (2) Signature technology: the method of this technology is the main technology of anti-virus software by matching features to detect the Trojan; (3) Real-time monitoring: Synchronization monitor the running process of the system (4) Behavior analysis: According to the program's dynamic behavior characteristics to identify Trojan.

These methods have their pros and cons, The Trojan detection technology of SPID feature optimized is a web-based, real-time detection technology; it is different from the previous Trojan detection technology. Using the network characteristics attributes meters to generate protocol model library and statistical-based to identify Trojans has a high recognition rate and a wide range of adaptability. In this paper, we proposed a method to detect and identify Trojan based the SPID characteristics optimized. The method analysis common protocol by SPID attribute meters and then generates a protocol model library to identify Trojan. Then, get the 12 optimized characteristic combinations of attribute meters by the statistical results.

2. SPID

SPID [3] (Statistical Protocol Identification) is a method based on statistical for protocol identification.

SPID's main purpose is to recognize network communication with which protocol the application layer is, it is not coarse-grained traffic classification (such as P2P or web), but accurate identification of which protocol it is.

Each protocol model has a range of fingerprint feature [4], which is the probability distribution of the application layer payload data or flow characteristics (size, direction, time of arrival, etc.). Each attribute meter is 2*256 arrays (see the example in Figure 1). The first line is the Counter vector, and the second line of is the Probability vector. Counter vector value is a positive integer, and each of the values expressed the value of each package to obtain in the corresponding index of this in the attribute meter.

Index	0	...	79	80	81	82	83	84	85	...	255
Counter vector	1263	...	715	935	296	919	1056	1845	643	...	1434
Probability vect.	0.006	...	0.003	0.004	0.001	0.004	0.005	0.009	0.003	...	0.007

Figure 1. Fingerprint feature of one attributemeter

Calculate the Kullback-Leibler [5](K-L) divergence of the session P with library model Q, calculate K-L value of the attribute meter value of P and Q, each attribute meter corresponds to a K-L value, last get the average K-L. For example, using attribute meter A, B, C to identify this session, calculate K-L value with default library model's attribute meter A, B, C, then get 3 K-L value, get the average K-Laver. K-L formula is as follows:

$$D_{K-L}(P_{sp} \parallel Q_{sp, prot}) = \sum_i P_{sp}(i) * \log \frac{P_{sp}(i)}{Q_{sp, prot}(i)} \quad (1)$$

Compare the K-L_{aver} with Pre-set threshold (equal 2.04), less than the threshold value, and the smallest value can be recognized P as protocol Q.

3. SPID-based feature optimization of Trojan Horse identification

3.1. Selection and generation of model library

We select 16 protocols or application as model library by observing existing product of identifying network traffic and specific Trojan identification, such as PoisonIvy 、 xPigeon 、 PcShare 、 PCanywhere 、 DameWare 、 RDP 、 Freegate and other common protocol, etc. Now, we not only select the Trojans as a model library but also the traffic without obvious characteristics and common communications applications as a background protocol.

Consider packet capture analysis carried out on the Trojans different functions (such as screen shots, to transfer files) during the Trojans crawl, The test found no big difference in the recognition effect. Training packets, generate 16 existing protocol model, see Figure 2, Sessions are multiple sessions under different environmental, Observations are packet number observed. Because of the SPID limited, it must be the source ip and destination ip corresponds to a port number during the packet training, thus distinguish packet as a single session, and then training model library.

Protocol Models		
Protocol	Sessions	Observations
BitTorrent	36	294
Dameware	2	17
eDonkey	34	333
Freegate	5	61
FTP	73	885
HTTP	121	812
Pcanywhere	1	12
pcshare	6	43
poisonIvy	4	35
POP	26	262
PPStream	4	19
QQ	2	40
RDP	1	16
SSH	54	577
SSL	81	734
Xpigeon	4	33

Figure 2. Model Library

3.2. Selection of SPID attributemeters

The SPID is based on the statistical methods and it has 34 attributemeters

available for free combination. Thus, we have to select the best combination of attributemeters for accurate identification of Trojan. The best attributemeters selection procedure is as follows:

- Crawl known protocol packets P1,P2....P34 of model library;
- Calculate K-L value between specific protocol of the model library (for example, set PcShare to P1) with the each protocol in model library, the results shown in Figure 3 (due to the large amount of data, lists only a part). Column of Meter are 34 attributemeters, column of B,C...J are the K-L values of P1 compare with model library;
- Select the attribute meter A's KL value less than 2.04 of PcShare in model library. That is the PcShare column in Figure 3. Elected large discrimination of B1 with other protocols KL value in A;
- In turn, get B2, B3.... B34 in accordance with step one, two without PcShare;
- Selected attributemeters which has common measure, Selection rules: the same attribute meter occurrences is greater than or equal to 4, the results are shown in Figure 4;
- Select the attribute meter that meets 4th step. Now, we selected 12 attributemeters, the column of Meter number is greater than or equal to 4 in Figure 4.

First, test the effect of recognition of P1 to P34, take example of P1(PcShare), the results are shown in Figure 5. Area ① are P1 packet information including the source

A	B	C	D	E	F	G	H	I	J
Basic	Checksum	Size(Torrent)	Checksum	Checksum	Checksum	Checksum	Checksum	Checksum	Checksum
AccumulatedDirectionalMeter 1	0.7561252	0.30757607	0.07592329	0.36502306	0.44744729	0.26373916	0.5900785	0.41758476	0.29516426
AccumulatedDirectionalMeter 2	0.7881292	0.31495918	0.07623149	0.36511999	0.44820691	0.26423963	0.5903097	0.41824088	0.29574039
AccumulatedDirectionalMeter 3	0.7651672	0.30747135	0.07592329	0.36502306	0.44744729	0.26373916	0.5900785	0.41758476	0.29516426
AccumulatedDirectionalMeter 4	0.7404825	0.29849489	0.07504954	0.35842137	0.43180387	0.25734898	0.57893145	0.40833163	0.29335587
AccumulatedDirectionalMeter 5	0.8021857	0.37397628	0.08271338	0.41345487	0.51282831	0.2793987	0.6897614	0.47029137	0.48389634
AccumulatedDirectionalMeter 6	0.6837495	0.33062854	0.07623149	0.36511999	0.44820691	0.26423963	0.5903097	0.41824088	0.29574039
AccumulatedDirectionalMeter 7	0.7881292	0.31495918	0.07623149	0.36511999	0.44820691	0.26423963	0.5903097	0.41824088	0.29574039
AccumulatedDirectionalMeter 8	0.6185824	0.30304819	0.07460284	0.35873083	0.44039693	0.25989229	0.57726463	0.40547716	0.29192021
AccumulatedDirectionalMeter 9	0.64162271	0.30495319	0.07504954	0.35852505	0.43258105	0.25991492	0.57893145	0.40833163	0.29335587
AccumulatedDirectionalMeter 10	0.71871515	0.30277122	0.07460284	0.35873083	0.44039693	0.25989229	0.57726463	0.40547716	0.29192021
AccumulatedDirectionalMeter 11	0.6770785	0.30047874	0.07498784	0.35836495	0.43700787	0.25927617	0.57663728	0.4072085	0.29323229
AccumulatedDirectionalMeter 12	0.6279183	0.2977837	0.07603839	0.35765267	0.44331524	0.25949197	0.57739193	0.40763983	0.29340862
AccumulatedDirectionalMeter 13	0.7881292	0.31495918	0.07623149	0.36511999	0.44820691	0.26423963	0.5903097	0.41824088	0.29574039
AccumulatedDirectionalMeter 14	0.8112222	0.36394871	0.08623571	0.42292829	0.52554009	0.28264123	0.69838487	0.4832452	0.49372388
AccumulatedDirectionalMeter 15	0.8033024	0.35211189	0.08469892	0.40851312	0.5057593	0.4328238	0.6306845	0.47457356	0.48294606
AccumulatedDirectionalMeter 16	0.9277618	0.4339888	0.10689787	0.53549651	0.61897973	0.35282262	0.58421961	0.48482705	0.4128222
AccumulatedDirectionalMeter 17	0.7767618	0.30232857	0.07603839	0.35765267	0.44331524	0.25949197	0.57739193	0.40763983	0.29340862

Figure 3. K-L output of P1 to model library

Meter No.	Basic	Checksum	Size(Torrent)	Checksum	Checksum	Checksum	Checksum	Checksum	Checksum	Checksum
1	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
7	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
8	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
2	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
3	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
9	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
10	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
11	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
12	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
13	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
14	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
15	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
16	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10
17	AccumulatedDirectionalMeter 1	AccumulatedDirectionalMeter 2	AccumulatedDirectionalMeter 3	AccumulatedDirectionalMeter 4	AccumulatedDirectionalMeter 5	AccumulatedDirectionalMeter 6	AccumulatedDirectionalMeter 7	AccumulatedDirectionalMeter 8	AccumulatedDirectionalMeter 9	AccumulatedDirectionalMeter 10

Figure 4. Comprehensive 12 meters

IP (Client IP), the destinate of the IP (Server IP), source port (C. Port.), destination port (S. Port), the observation of the number of packets (Observation) and the identified protocol (Protocol). Area ② are output K-Layer(Divergence)of the PcShare to model library as well as the discrimination compared to other protocols (Match Percentage). Regional ③ is the default model library.

3.3. Performance verification

Configuration 12 common attributemeters we selected, running SPID Algorithm Proof-of-Concept 0.4.6 [3], capture network traffic packets P include PoisonIvy, PcShare through Wireshark.

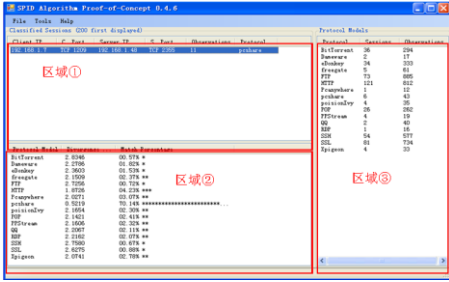


Figure 5. PcShare recognition effect

The same method as above, test the packet p which includes PoisonIvy and PcShare, we know that the port 5555 of host computer 10.10.10.226 is PcShare, port 3460 is PoisonIvy, the result is shown in Figure 6.

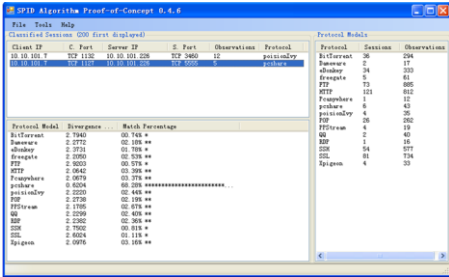


Figure 6. Effect of obfuscation packet P

We conducted a number of experiments, using 12 attributemeters can be able to maintain the original SPID detection accuracy, but detection speed significantly improved, therefore, confirmed that the proposed method is accurate and efficient.

4. Conclusion and prospect

What we proposed is a principle to early identify Trojan, generate model library by refining the traffic of Trojan as well as the remaining common protocol, it contains 16 protocols and applications. Then, using statistical observation method to extract the 12 best combinations of attribute meters can

identify the protocol of this model library. Finally, give the recognition result, so can be further verify the feasibility of this principle. However, due to the less type of model library, it can be identify the protocol which is in. In addition, the model library training is less, but theoretically, it is more accurate if it has enough training. These problems are the focus work of the future.

5. References

- [1] Ming Zhu, Sai Xu, Chunming Liu. "Analysis of Trojan Horse and Its Detection", Computer Engineering and Applications, vol.2, no 28, pp.176-178. 2003
- [2] D. Agrawal and et al, "Trojan detection using icfingerprinting", *IEEE Symp. On Security and Privacy*, pp.296-310, 2007.
- [3] Hjelmvik, E., John, W., "Breaking and Improving Protocol Obfuscation", Technical Report No: 2010-05, Department of Computer Science and Engineering, Chalmers University of Technology, Gothenburg, Sweden 2010.
- [4] E. Hjelmvik. Wiki page of SPID. <http://sourceforge.net/apps/mediawiki/spid/index.php>
- [5] S. Kullback and R. A. Leibler, "On information and sufficiency", *Annalyse of Mathematical Statistics*, vol. 22, pp. 49-86, 1951.
- [6] Xiaoqing Han, Jianfeng Wang, Wei Zhong, "Analysis and prevention of computer virus", Beijing: Publishing House of electronics industry, pp.207-209, 2006.

Acknowledgement

Supported by National Natural Science Foundation of China (61073188), China Postdoctoral Science Foundation (20100471355)