

SVM based Traffic Classification for Mitigating HTTP Attack

V. Punitha* and C. Mala

National Institute of Technology, Tiruchirappalli, India
vpunitha21@gmail.com, mala@nitt.edu

Abstract

The advancement in Internet technology brings new dimension to commercial applications, entertainment and information sharing. Consequently, many web services are launched in almost all needs of the internet users. The development of effective network infrastructure increases the usage of these services. However the convenience of using the web services are blocked by denial of service attack, which is the foremost web threat. This attack injects malicious traffics into the internet which deeply affects the availability of services. Categorizing the malicious traffic from normal traffic facilitates the elimination process. In view of eliminating the most victimized attacks which deny the services to the potential users, this paper proposes a classification method based on machine learning technique. The proposed SVM based classifier discriminates the HTTP attacks that intentionally blocks the computing resources to the legitimate users based on network flow properties. The network flow properties are selected by the proposed optimization method. The simulated results exhibit that with optimized feature set, the classification performance of the proposed classifier using RBF kernel is competently higher when compared with other kernel models.

Keywords: Denial of Service attack, Application layer attack, HTTP attack, Support Vector Machine

1 Introduction

Modern advanced technology makes the human to rely on Internet for education, entertainment and employment. Many web services are launched for commercial and educational purposes by industries and academia. Growth of bandwidth and popularity of web services increase their usage and this in turn increases the internet traffic. Most of the web services use HTTP as their application layer protocol. The ethical behaviour of the user has been changed due to intellectual challenges and economical gain [5]. This consequently affects the requesting behaviour of the users and induces to create malicious traffic in the internet. Increased HTTP traffic and advanced technology persuade the attackers to generate application layer attacks [13]. According to the study based on backbone network traffics of US and China, the occurrences of application layer attacks are higher than lower layer attacks [6]. Although various techniques are present to prevent the malicious traffic, classification of internet traffic is believed to be the perfect solution to mitigate the attacks and irregularities.

Machine learning technique is the best choice for automatic categorization of traffic. Much research work has been carried out using machine learning techniques for identification of traffic applications [4]. They apply externally observable parameters like, packet length, port number, arrival time etc., [3]. SVM (Support Vector Machine), is popularly known for its theoretical justification. It is a widely used supervised machine learning method, classifies the data using generated hyper plane. Although the machine learning techniques are applied in traffic classification to identify the applications of traffic, applying these techniques in detection of HTTP attacks needs further analysis. Hence this paper proposes a classification approach to discriminate HTTP attack from generic internet traffic.

Research Briefs on Information & Communication Technology Evolution (ReBICTE), Vol. 4, Article No. 4 (August 15, 2018)

*Corresponding author: Research Scholar, Department of of Computer Science and Engineering, National Institute of Technology, Tiruchirappalli, India

2 Related Work

This section presents few significant research works carried out on traffic classification and detection of anomalies in internet traffic.

Authors in [10], applied group testing (GT) ideologies to discriminate application layer attacks and proposed constraint-based GT system. Authors proposed three detection models using sequential and partial non-adaptive detection methods and analyzed the theoretical complexity. The methods are evaluated in terms of false positive, false negative rates and detection delay. Authors in [9] proposed a novel methodology to identify DDoS attacks using generated multidimensional access matrix from extracted network flow. The number of attributes in matrix is reduced using Principle Component Analysis (PCA) to improve the correlation among identified attributes. Then Naive Bayes and K-Nearest neighbourhood classification methodologies are applied to discriminate the traffic. Authors proved that the performance of the classification with PCA is better than the classification with access matrix. A Bio-Inspired technique is proposed in [7], to detect application layer DDoS attack. Authors proposed a bat algorithm for detection of HTTP attack at early stage to avoid the degradation of network performance. Features are extracted from absolute time interval which includes five attributes describing time frame and session. Authors proved that the proposed bat algorithm produces higher accuracy than others. Authors in [8] briefed about the new types of irregularities in HTTP/2 services and slow rate attacks. An experiential study on these attacks are presented using both plain text and encrypted messages. Here, an anomaly detection methodology is proposed using chi-square test and it classifies the attacks with highest accuracy. A new methodology to identify application layer DDoS attacks is proposed in [13]. Here, Real-time Frequency Vector (RFV) and traffic model are generated to describe the visiting status of all resources. The HTTP requests are analyzed and four types of irregularities are identified and eliminated from the traffic. Authors evaluated their proposed methodology in both simulated and real-time environment.

The above survey emphasizes that accurate assessment of HTTP request and automatic categorization of HTTP traffic in the view of discrimination of attacks need further investigation. So this paper proposes SVM based HTTP Request Classification which examines HTTP request and categorizes the HTTP attack from generic traffic. The rest of the paper is organized as follows. Section 3 introduces SVM model and the proposed SVM based HTTP Request Classification. The simulated results are analyzed in Section 4 and finally Section 5 summarizes the work

3 Proposed Model

The significance of SVM is its kernel functions and high dimensional mapping. The proposed internet traffic classification uses multiple network flow properties, so it needs complex mapping. Since SVM supports complex mapping, it is highly desired for the proposed model to map the network flow properties.

3.1 Proposed SVM based HTTP Request Classification (SHRC)

Application layer attacks are common in various computing environment like, cloud computing, ad hoc computing or mobile cloud computing. But, application layer attacks are more difficult to detect than other attacks, as the attacks are propelled only after initiating the connection. It needs intense analysis on the request structure. The proposed SVM based HTTP Request Classification (SHRC) model analyzes the structure of the request and categorizes the HTTP attack from generic using network flow features. The work flow of SHRC is depicted in Figure 1. SHRC captures the incoming traffic and extracts network flow features from the traffic. Extracted features are then optimized and transformed into training data. Finally, SHRC model is trained to classify HTTP attack.

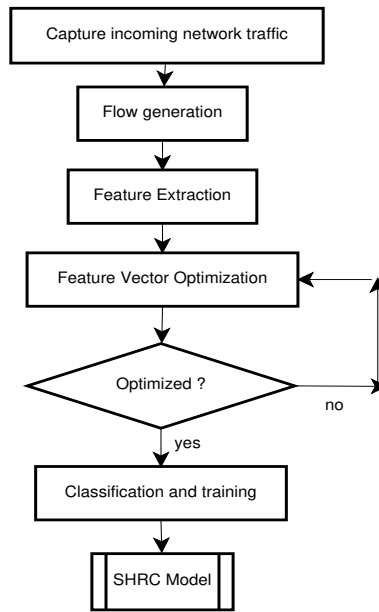


Figure 1: SVM based HTTP Request Classification

Generally, when HTTP request is prompted by the client, a TCP connection is initiated with a SYN request to the server. Then it is followed by [SYN, ACK] and ACK. After connection establishment, the HTTP-GET or POST request is communicated to the server and followed by other communications or another HTTP-GET or POST requests. Finally the client request is ended by connection termination with [FIN, ACK]. All these communications are recorded in sequence and this is called as network flow or flow. So each flow includes all the communications in both directions pertaining to a client request. Similarly, the flow of all the requests from different clients are assembled. The network flow properties such as source IP, destination IP, destination port, arrival time are collectively derived for each network flow and used as features to describe the flow. Some of the properties are directly available like, source IP, destination IP, destination port and few other properties like, service time, number of packets in a flow, average packet size etc., are computed from the available network flow properties. Around 248 network flow properties can be computed [12], but not all the properties are contributing in identification of HTTP attack. So SHRC optimizes the feature selection in an incremental way as depicted in Figure 1.

Initially, basic flow properties that are needed to identify HTTP application are chosen as feature vector, i.e., [source IP, destination IP, destination port]. Next, service time is added and in this paper this feature vector is denoted as first discriminator set (DS). As the service time of a HTTP attack is unusually longer when compared to normal request and this characteristics distinguishes the attack from normal traffic. Here, service time is computed by measuring the time between the connection establishment and termination of a HTTP request. Each request prompted by the client may include more than one HTTP-GET requests. It means that when the client is sequentially accessing more than one links/images/videos, every time distinct GET request is generated. Service time of the requests with more number of GET requests will automatically be high. So, the request cannot be declared as attack, considering only excess service time as the reason. Hence, number of GET requests in each flow is computed and added along with the above selected features, i.e., [source IP, destination IP, destination port, service time, no. of get requests], denoted as second DS. It is also noted that session time is used as a major discriminator in [7] and number of requests to the resources is used as a discriminator in [9] and [13]. Further, the service time is analyzed deeply and it is found that the duration of service time varies according to the type of information access, i.e., when the client accesses images, service time will be longer than the access time

of text pages and it is still longer when it is video content. So SHRC also considers the type of access as another feature in third DS. The type of the content is derived from each GET request in the network flow. SHRC optimizes the feature vector using the measured classification accuracy, as the accuracy is the best metric to measure the detection rate of both attack and normal traffic. The accuracy defines the ratio of correctly detected traffic to the total number of traffic. For improved detection rate of HTTP attacks, the features of the normal traffic and attacks are to be investigated deeply and it is decided to be executed in the future work. Moreover SHRC can also be applied in ad hoc computing environment or mobile cloud computing, as the constructed DS is more generalized, i.e., not specific to application environment. The features in optimized DS can be calculated from the traffic captured in ad hoc network. As the lower layer protocols differ in this network environment, other network flow features may slightly vary. For ad hoc network, transport layer protocol may be Feedback-based TCP / TCP with explicit link failure notification / TCP-buffering and sequence information / Ad hoc TCP / any other application specific transport protocol. This provides additional network flow features. As network flow features are used in this paper to predict HTTP attacks, additional flow features assist in distinguishing more attacks in ad hoc network competently. This needs further analysis and optimization on feature selection which are planned in the future work.

Detection of HTTP attacks are more difficult than flash crowd. Flash crowd can be discriminated with arrival statistics. But HTTP attacks deny the services to the legitimate users by exhausting the usage of the computing resource for long time. So, service time and the number of access requests are examined. For long service time, if the number of access is less, then it is declared as HTTP attack. In other case, if the number of access is high, then the type of access is examined. If the type is not multimedia content, then it is classified as attack, as the text link could not consume long service time. SHRC model is trained with the optimized feature vector. C and gamma values are chosen using cross validation and grid search [11], as these values influence the performance of the SVM model. Finally, SHRC model is tested with the obtained testing traffic. The probability of prediction is estimated during testing.

4 Simulation and Performance Analysis

The performance of the proposed SVM based HTTP Request Classification (SHRC) model is analyzed in a simulated environment which includes a server, clients and a dispatcher. The dispatcher is a global scheduler which is configured to capture the incoming traffic and analyzes the traffic before forwarding them to the server. The dispatcher captures the traffic in different time intervals using Wireshark, a network traffic analyzer [2]. The Wireshark filter is configured to capture HTTP requests. In this paper, three sets of traffic traces are used for analysis i.e., one for training the SHRC model and other two traffic traces are used for testing. The training traffic trace is captured for the timespan of 465ms and it consists of 560 HTTP requests. The testing traffic traces are captured in different time periods for the time span of 523ms and 652ms and they consists of 600 requests and 729 requests respectively. In this paper, the testing traffic traces are denoted as 'HTTP traffic 1' and 'HTTP traffic 2'. For smoother analysis, the attacks traffics are further merged with the captured traffic. Then the network flow features for each HTTP request is exported using Wireshark. The network flows are constructed using Java platform and feature vectors are also computed. Then using LibSVM, Java based tool [1], the SHRC model is developed. The performance measures that are applied in this paper to evaluate the proposed SHRC are accuracy, precision and recall. Precision defines the ratio of number of correctly classified data to total number of predicted data of a class. Recall defines the ratio of number of correctly classified data to number of labelled data of a class. Accuracy defines the ratio of number of data that are correctly classified for both the classes to the total number of data[4].

4.1 Cross Validation

Proper selection of C and gamma, improves the performance of SVM classification [11]. This paper implements n-fold cross validation technique to choose these values. For this reason, 560 HTTP requests (training data) are applied for n-fold cross validation, with n=10. Values of C and gamma are chosen using grid search [11]; the pair that produces higher accuracy, i.e., above 90% are measured and plotted as a contour map in Figure 2. It is inferred from Figure 2 that for a certain range of C and gamma values, the accuracy is stable. So the pair which produces highest accuracy, as well as soft margin is chosen by SHRC model and used for further classification, i.e., C=23 and gamma=21.

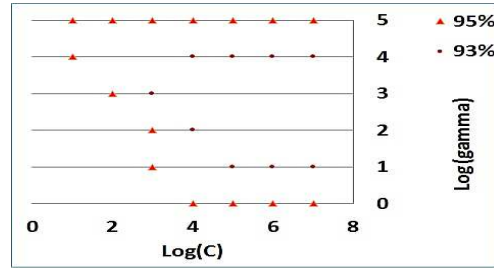


Figure 2: Kernel parameters: C Vs Gamma

4.2 Feature Selection

SHRC optimizes the selection of flow properties using an incremental method. This method applies accuracy to evaluate the selection process. Three sets of flow properties are taken for analysis as described in section 3.1. First discriminator set (DS) includes service time along with basic flow features to discriminate the HTTP attack, as the service time of attack traffic is unexpectedly high; DS1=[source IP, destination IP, destination port, service time]. Secondly, as the service time is varied due to number of access requests (GET requests), it is included in second discriminator set, i.e., DS2=[source IP, destination IP, destination port, service time, no. of get requests]. Similarly, the type of access like image, video etc., is included in third discriminator set, as it also increases the service time, so, DS3= [source

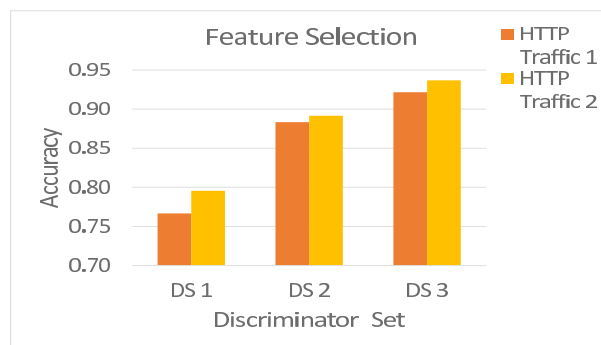


Figure 3: Discriminator set Vs Accuracy

IP, destination IP, destination port, service time, no. of get requests, type of get request]. Classification is performed with all the three DS independently and TP, TN are measured and accuracy is computed for each model and plotted in Figure 3. Accuracy = (TP+TN)/N, where N defines total number of traffic, True Positive(TP) defines number of normal traffic that are classified correctly and True Negative defines number of HTTP attacks that are predicted correctly [4]. It is inferred from Figure 3 that the third

discriminator set produces highest accuracy among the three sets. It is also noted that in terms of classification accuracy, the proposed optimized feature set produces better results than second discriminator set which includes service time and number of requests which are used as major discriminators in the works presented in the literature survey. This implies that the detection rate of attacks and normal traffic is high while using source IP, destination IP, destination port, service time, no. of get requests, type of get request as feature vector. Hence they are applied for further analysis and classification.

4.3 Performance of SHRC

In this paper two unbiased traffic traces are captured and used for analysis. SHRC discriminates HTTP attack from generic traffic using optimized feature set. When service time is long, the other two features are examined. If the lengthy service time is not because of more number of GET requests or by its type, such as multimedia content, then the request is declared as attack. The classification is performed using three kernel functions, Linear, Polynomial and RBF. Every time TP, TN and FP are computed. False Positive (FP) measures number of attacks that are not predicted correctly. Precision and recall values are computed and plotted in Figure 4. It is inferred from Figure 4 that for both traffic traces, recall is high for linear model and precision is high for RBF kernel model. Higher values of precision implies that more number of attacks are classified correctly, whereas high recall value conveys that more number of normal traffics are classified correctly. As the focus of SHRC is to identify all HTTP attacks, the precision is the chosen as best measure for detection of attacks. For further analysis, probability estimates are computed for testing traffic and plotted in Figure 5. The probability estimates illustrate the probability of classifying the attacks in different kernel models. It is inferred from Figure 4 that precision of both RBF and

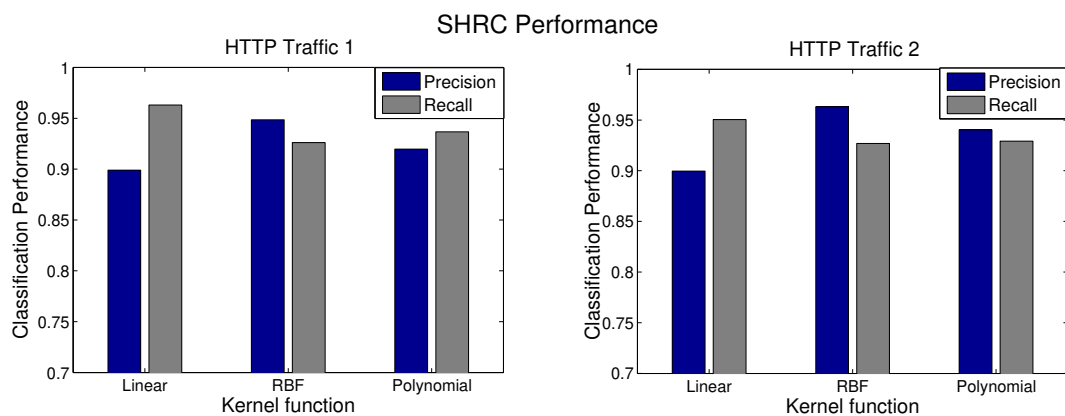


Figure 4: Kernel function Vs Classification performance

polynomial kernel models are high compared to linear kernel model. But it is observed from Figure 5 that in polynomial kernel, traffics are classified with the probability of 0.6, whereas RBF kernel classifies the traffic with the probability of 0.9. This implies that both TP and TN are classified with higher probability in RBF, in other words, both attacks and normal traffics are classified with higher probability. Linear kernel also classifies the traffic with the probability of 0.9 in traffic 2, but here the precision is minimum, shown in Figure 4. It implies that liner model identifies less number of attacks with high probability. Hence, the performance of RBF kernel model is better than other two models in identification of HTTP attacks.

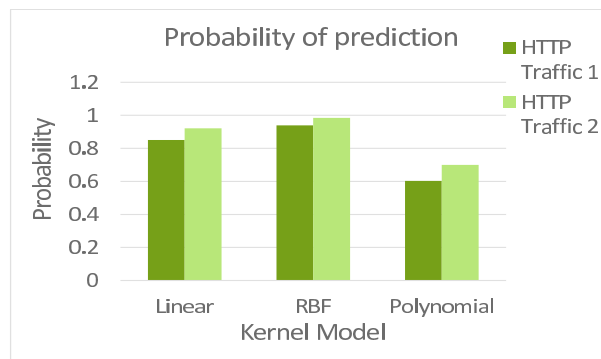


Figure 5: Kernel model Vs Probability of prediction

4.4 Classification and Attacks

Two traffic traces are captured in random period of time for identification of HTTP attacks and taken up for analysis. Traffic 1 with 600 requests and Traffic 2 with 729 requests. SHRC discriminates the attack using the optimized feature vector. HTTP attacks are intentionally created by botnets to block the usage of the resource. This paper examines the behaviour of the botnet by finding the dependency between occurrence of attack and arrival rate. If the number of attacks are high only on a specific period, say beginning of the arrival or when the arrival rate is high etc., then it is enough to execute SHRC only during that time. So, to monitor the rate of occurrence of attack during capturing period, the attack detection rate is computed for every 100 arrivals and it is plotted in Figure 6. It is observed from Figure 6 that the detection rate is not same in both traffic traces. The occurrence of attacks are high in traffic 2 during beginning of arrival, whereas in traffic 1 it is measured low. In traffic 1, maximum number of attacks are detected during middle of the arrival. So, the occurrences of attack do not depend on rate of arrival. It implies that regardless of arrival rate, the attacks spread over the entire traffic. Hence the classification of HTTP attack using SHRC is to be progressed continuously.

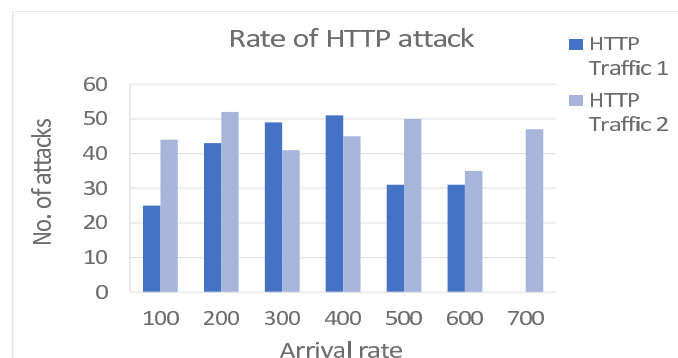


Figure 6: Arrival rate Vs Number of attacks

5 Conclusion

Classification of internet traffic discriminates most victimized attacks that deny the usage of computing resources to legitimate users. In this paper, a machine learning classification approach is proposed to detect the HTTP attack using network flow properties. An optimization technique is implemented to select the flow properties for the construction of feature set. The proposed SVM based classifier proficiently

categorizes the HTTP attacks from normal traffic with the optimized feature set. Traffics for analysis are captured using Wireshark and the proposed model is simulated using LibSVM. The performance of the classification is analyzed with precision, recall and accuracy. It is evidenced from the simulated results that using proposed optimization method, the proposed RBF kernel classifier model has outperformed other models. Deeper investigation is essential on optimization of feature selection for improved detection rate and better prediction of HTTP attacks. Further analysis on feature selection for different computing environments and optimization are planned to be implemented in the future work.

References

- [1] C.-C. Chang and C.-J. Lin. LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, April 2011.
 - [2] L. Chappell and G. Combs. *Wireshark network analysis: the official Wireshark certified network analyst study guide*. Protocol Analysis Institute, Chappell University, 2010.
 - [3] M. Finsterbusch, C. Richter, E. Rocha, J.-A. Muller, and K. Hanssgen. A survey of payload-based traffic classification approaches. *IEEE Communications Surveys & Tutorials*, 16(2):1135–1156, October 2014.
 - [4] T. T. Nguyen and G. Armitage. A survey of techniques for internet traffic classification using machine learning. *IEEE Communications Surveys & Tutorials*, 10(4):56–76, 2008.
 - [5] K. M. Prasad, A. R. M. Reddy, and K. V. Rao. Dos and ddos attacks: defense, detection and traceback mechanisms-a survey. *Global Journal of Computer Science and Technology*, 14(7), 2014.
 - [6] K. Singh, P. Singh, and K. Kumar. Application layer http-get flood ddos attacks: Research landscape and challenges. *Computers & security*, 65:344–372, March 2017.
 - [7] I. Sreeram and V. P. K. Vuppala. Http flood attack detection in application layer using machine learning metrics and bio inspired bat algorithm. *Applied Computing and Informatics*, October 2017.
 - [8] N. Tripathi and N. Hubballi. Slow rate denial of service attacks against http/2 and detection. *Computers & security*, 72:255–272, January 2018.
 - [9] S. Umarani and D. Sharmila. Predicting application layer ddos attacks using machine learning algorithms. *International Journal of Computer, control Quantum and information Engineering*, 8(10), 2014.
 - [10] Y. Xuan, I. Shin, M. T. Thai, and T. Znati. Detecting application denial-of-service attacks: A group-testing-based approach. *IEEE Transactions on parallel and distributed systems*, 21(8):1203–1216, September 2009.
 - [11] R. Yuan, Z. Li, X. Guan, and L. Xu. An svm-based machine learning method for accurate internet traffic classification. *Information Systems Frontiers*, 12(2):149–156, April 2010.
 - [12] L. Zhen and L. Qiong. A new feature selection method for internet traffic classification using ML. *Physics Procedia*, 33:1338–1345, 2012.
 - [13] W. Zhou, W. Jia, S. Wen, Y. Xiang, and W. Zhou. Detection and defense of application-layer ddos attacks in backbone web traffic. *Future Generation Computer Systems*, 38:36–46, September 2014.
-

Author Biography



V. Punitha completed a Master of Engineering (ME), Computer Science and Engineering from National Institute of Technology, Tiruchirappalli, India in 2003. Currently she is pursuing Ph.D. degree in the same Institute at Department of Computer Science and Engineering. Her research area of interest includes Parallel and Distributed Systems, Network Security and Soft Computing Techniques.



C. Mala is a Professor in the Department of Computer Science and Engineering, National Institute of Technology, Tiruchirappalli, Tamil Nadu, India – 620 015. Her research area of interest includes Data Structures & Algorithms, Computer Networks, Parallel Algorithms, Computer Architecture, Sensor Networks, Soft Computing Techniques, Image Processing, Intelligent Transportation Systems and Vehicular Adhoc Networks.