# Intrusion Detection Using Geometrical Structure

Aruna Jamdagni[1,2], Zhiyuan Tan[1], Priyadarsi Nanda[1], Xiangjian He[1,3] and Ren Liu[2]

[1]Centre for Innovation in IT Services and Applications (iNEXT)
University of Technology, Sydney
Sydney, Australia
{arunaj, thomas, pnanda, sean}@it.uts.edu.au
[2]Commonwealth Scientific and Research Organization (CSIRO)
Sydney, Australia
Ren.Liu@csiro.au
[3]Lab of Biomedical Information Technology
University of Aizu
Japan

*Abstract*—**We propose a statistical model, namely Geometrical Structure Anomaly Detection (GSAD) to detect intrusion using the packet payload in the network. GSAD takes into account the correlations among the packet payload features arranged in a geometrical structure. The representation is based on statistical analysis of Mahalanobis distances among payload features, which calculate the similarity of new data against pre-computed profile. It calculates weight factor to determine anomaly in the payload. In the 1999 DARPA intrusion detection evaluation data set, we conduct several tests for limited attacks on port 80 and port 25. Our approach establishes and identifies the correlation among packet payloads in a network.**

*Keywords-Intusion Detection; Payload; Geometrical Structure; Mahalanobis Distance; Pattern Recognition*

## I. INTRODUCTION

The growth of Internet and local area networks provide quality and convenience to human life but at the same time provides a platform for network hackers and criminals. Internet security hence becomes an important problem in near future. The concept of Intrusion Detection was introduced in 1980 by J.P. Anderson [2], and since then has become an active field of research. According to Computer Emergency Response Team (CERT) [1], 32,956 vulnerabilities were reported from many sources through 1995 until the first quarter of 2007. These vulnerabilities provide opportunities for attackers to launch attacks to computer systems and gain an access to the computers. The goal of an Intrusion Detection System (IDS) is to characterize attacks manifestations to positively identify all true attacks without falsely identifying non-attacks.

Intrusion Detection Systems are components designed to detect intrusion and also to prevent a system from being compromised. There are three major types of Intrusion Detection Systems. Anomaly detection system creates a model of normal behavior, and flags suspicious behavior or any deviation from the normal behavior. The main strength of anomaly detection is the ability to recognize novel attacks, and the major weakness is that it is susceptible to false positive alarms. Signature-based system or misuse detection system uses knowledge base to recognize directly the signatures of intrusion attempts. This technique is susceptible to a slight variation of the attack signature and also to an unknown attack. The Snort and Bro are popular examples of signature based intrusion detection system [5, 7] used commercially. Specification-based system [6] relies solely on the frequency of the input data based on system calls, or protocols such as IP, TCP and UDP. The strength of this technique is computationally light and does not require and need maintenance of many types of parameters, and or profile activities. However the weakness of this system is that it needs detail design to avoid missed attack types.

In this paper, we present a new model, called Geometrical Structure Anomaly Detection (GSAD) based on pattern recognition technique used in image processing.

The structure of the paper is as follows. Section 2 describes related work in the field of anomaly detection. Section 3 briefly describes methods used in intrusion detection. In Section 4 we discuss our proposed model. Section 5 describes the implementation of the model and Section 6 presents conclusions and future work.

## II. RELATED WORKS

The misuse detection systems or signature based systems rely on signatures of known attacks or pre-defined rules to match and identify known attacks. Presently in industry, rule based network intrusion detection systems such as Snort [5] and Bro [7] are most popular. These systems use signatures or finger prints to identify known attacks. But signature based systems are clueless in case of novel attacks. Examples of such novel attacks are Zero day attacks, Mutation attacks etc. A Zero day attack is a computer threat that tries to exploit unknown computer application vulnerability. In Mutation attack, known instances of attacks are transformed into distinct instances which have the same power of exploitation. Since attack signature is different from stored known signature due to transformations, such attacks are less likely to be detected by signature based systems.

Anomaly detection systems model the normal profile of system behaviour, and any deviation from this behaviour will be identify as a possible attack.

There are two anomaly based detection systems. One is based on specification (or a set of rules) regarded as good or normal behaviour, which depend on the human expertise, and the other one learns the behaviour of the system under normal operation automatically. Anomaly detection systems such as PAYL [14], SPADE [8], NIDES [9], PHAD [10], ALAD [11] and NATE [12] compute (statistical) models for normal network traffic and generate alarms when there is a large deviation from the normal model. Some of these systems use different algorithms to model the normal network traffic behaviour and feature extraction techniques from the available audit data. SPAD, ALAD and NIDES use source and destination IP and port addresses and TCP connection in the development of model, while PHAD uses 34 features, extracted from the packet header fields of Ethernet, IP, TCP, UDP, and ICMP packets. For these systems the detection rate of protocol based attacks is good but poor for application based attacks, as these systems ignore the payload contents.

NATE and PHAD system use first 48 bytes as a statistical features starting from IP header and can include at most 8 bytes of network packet payload. ALAD models incoming TCP request first word or token of each input line out of 1000 application payloads as a feature for HTTP and SMTP protocols.

Kruegel at al [13] describes a service-specific intrusion detection system. They use the type, length and payload distribution of the request as features to compute anomaly score of a service request and use chi-square test to calculate anomaly score of new request. They group 256 ASCII characters into six segments: 0, 1-3, 4-6, 7-11, 12-15, and 16-256, and compute one single distribution model of these six segments. Ke Wang and Salvatore J. Stolfo [14] developed full byte distribution model conditioned on the length of payloads and use Mahalanobis distance to calculate anomaly score. They also introduced the concept of automatic clustering of centroids to increase the accuracy and reduce the resource consumption. In contrast, we prepose a novel approach to develop GSAD model for packet payload. Each network connection between a pair of hosts will be viewed as an object in an image (to be recognized through image processing), and each image will be viewed as a pattern to be classified as normal or anomalous traffic class based upon the given information about the connections. This model includes the correlation between various payload features and increases the detection accuracy. We use Mahalanobis Distance Map to calculate the difference between normal and anomaly of new network traffic. We will use DARPA 1999 IDS dataset [15, 16, 21] as a benchmark to evaluate the robustness of our algorithms. This dataset is not without its critic. McHugh [17] pointed out that the DARPA/MIT Lincoln Laboratories IDS test used generated data, but MIT researchers never did any tests to show that the generated data was a representative of real data. Further more they did not conduct tests to verify that their attacks were representative of real attacks. The detail description of our model is given in Section 4.

## III. INTRUSION DETECTION METHODS

Various supervised and unsupervised algorithms used by researchers for intrusion detection with varying degree of accuracy are reviewed in [3, 4]. Some of them are summarized here in brief.

*Statistical Method*: Statistical methods are commonly used for pattern recognition. The IDS observes a set of normal behaviour and calculates one or more statistics identified by a person or some other portion of the IDS to be significant. It can provide accurate information about the malicious activities which occur over a long period of time, but it is hard to determine thresholds that balance the likelihood of false positive alarms with the likelihood of false negative alarms.

*Artificial Neural Networks*: One or more data sources are used to train the neural net to recognise normal behaviour. The neural net then identifies behaviour which does not match its training experience. It is a data clustering method based on distance measurement. This approach applies biological concepts to machines to recognise pattern. It requires minimum priory knowledge, and with enough layers and neurons can create any complex decision region.

*Data Clustering*: Data clustering is a technique for finding data in unlabelled data with many dimensions. It is an unsupervised method. It can learn from and detect intrusions in the audit data without explicit descriptions of various attack classes.

*Immune systems*: It mimics natural immunology as observed in biology. Several models exist such as negative selection, immune network model and clonal selection. Cells can sense not only the evidence for antigen presence, but also danger signals.

*Decision Tree*: This can be used to show possible consequences for particular occurrences where there are conditional probabilities for each occurrence. They perform efficiently with a large amount of data.

*Fuzzy Logic*: It is a set of rules and concepts and approaches designed to handle vagueness and imprecision. A set of rules can be created to describe a relationship between input variables and output variables, which may indicate whether an intrusion has occurred. It uses membership function to evaluate the degree of truthfulness.

However GSDA model uses statistical intrusion detection method to identify an abnormal behaviour in the network.

## IV. GEOMETRICAL STRUCTURE BASED IDS

In this section, we give a comprehensive introduction about the GSAD which employs geometrical structure into payload-based anomaly detection. This IDS is based on a statistical analysis of Mahalanobis Distances Map among characters appearing in network traffic and distinguishes abnormal traffic from normal ones with patterns. The architecture of GSAD is shown in Fig. 1.

In the following figure solid arrow indicates data flow inside the GSAD. The GSAD Architecture contains the following 5 components:

*Payload feature classifier*: This component is used in the network traffic payload classification phase. The network traffic data are grouped into various categories by using Wireshark based on four conditions including size of payload, destination address, services and direction of traffic flow. The source of the network traffic can be real network and collected tcpdump files.

*Payload feature analyst*: The payload feature analyst is first key constituents of Geometrical Structure Payload Model (GSPM). It is responsible for payload feature analysis using statistical analysis approaches and prepares raw data for the following analysis phase.

*Payload geometrical structure model*: It is the second key constituent of GSPM. The payload geometrical structure model is developed by using a statistical method for anomaly detection based on Mahalanobis Distance Map. The source data are well prepared by the payload feature analyst.

*Attack recognizer*: This part of GSAD handles the recognition of attacks from the input network traffic. It compares each incoming packet with normal and abnormal payload geometrical structure model, and then gives out the score which is the criterion to either generate alarm or not.

*Acknowledge/Communication*: In this module, the attack alarm will be generated if the score of a packet is larger than the threshold and report to the administrator. Otherwise it will consider the packet is a normal one.
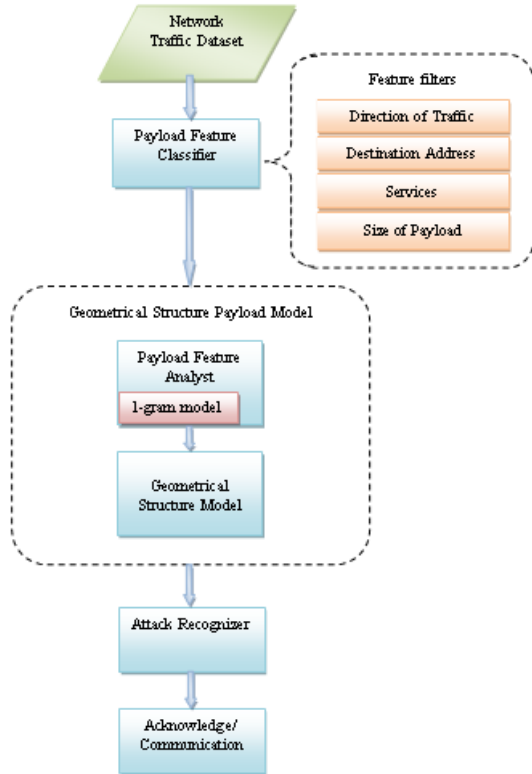


Figure 1.    GSAD architecture

## A.  GSAD Model Characteristics

The GSAD intrusion detection system uses pattern recognition techniques. They facilitate the anomaly detection ability of the system without the prior knowledge of an attack. Similar to other anomaly detection systems, GSAD models the normal behavior of the network traffic rather than the malicious ones. Moreover, the most significant contribution of GSAD is the integration of geometrical structures and payload-based anomaly detection systems, which has not been considered in other related researches. There are two models involving into our GSAD system, namely *1-gram payload model* [14] and *geometrical structure model* [13, 19].

*1) One-gram Payload Mode:* The 1-gram payload model is a payload based statistical model. The content of network packets is the analysis object of the 1-gram payload model which calculates the average frequency of each ASCII character (0-255). It does not take network packet header features into account. However, the average frequency is not the most appropriate characterizing feature for describing network behaviors because the same average frequency which can be obtained from some very different character frequencies and some steady character frequencies. Therefore, some other criteria are expected to interpret the behaviors of variant network traffic. They are the mean value and standard deviation of each byte's frequency.

In fact, these criteria are all derived from ACSII character frequency. So, when building the 1-gram payload model, feature vector is the compulsory constituent needed to be calculated first. For a payload model, the feature vector is a set of relative frequencies is the occurrences of each ASCII character to the total number of characters appearing in the payload. In general, each feature vector can be represented as the following (1).

$$X = [x_0\ x_1\ \dots\ x_{255}] \tag{1}$$

Then, given a set of feature vectors, we can compute the mean value and standard deviation of each byte's frequency. Here, we assume that there is a network traffic dataset with $n$ network packets. The mean value and standard deviation of each byte's frequency are described as (2) and (3), respectively.

$$X = [x_0\ x_1\ \dots\ x_{255}] \tag{2}$$

$$\sigma = [\sigma_0\ \sigma_1\ \dots\ \sigma_{255}] \tag{3}$$

Here,

$$x_i = \frac{1}{n}\ \sum_{k=1}^{n} x_{i,k}\ (0 \le i \le 255) \tag{4}$$

$$\sigma_i = \sqrt{\frac{1}{n}\ \sum_{k=1}^{n}(x_k - x_i)^2}\ (0 \le i \le 255) \tag{5}$$

The mean value and standard deviation vectors, $X$ and $\sigma$, are stored in a model $M$. Whereas due to the network traffic

dataset consists of traffic generated by the various network services. Therefore we need to classify network traffic based on the following features: size of payload, destination address, services and direction of traffic flow. The models are developed according to this group of features.

*2) Geometrical Structure Model:* The Geometrical Structure Model (GSM) is a pattern recognition technique used to detect similarity between the normal behavior with the new input traffic. Although this model has been adopted into the research of human detection, it is still a new concept to intrusion detection. In this subsection, we present an explanation about the practical application of geometrical structure model in payload-based anomaly detection. The model takes into account the correlations among different features (256 ASCII characters). Thus, for each network packet, there is a feature vector defined by (1). The average value of features in the 1-gram model is

$$\mu = \frac{1}{256} \sum_{i=0}^{255} x_i \qquad (6)$$

The covariance value of each feature is

$$\sum_i = (x_i - \mu)(x_i - \mu)' \ (0 \le i \le 255) \qquad (7)$$

In order to investigate the relationship among the characters, we compute the Mahalanobis distance (indicated by $d_{(i,j)}$) between every two characters.

$$d_{(i,j)} = \frac{(x_i - x_j)(x_i - x_j)'}{\sum_i + \sum_j} \ \ (0 \le i, j \le 255) \qquad (8)$$

Based to the above calculation, the Mahalanobis Distance Map (MDM) of a network packet is constructed as the following,

$$D = \begin{matrix} d_{(0,0)} & d_{(0,1)} & \cdots & d_{(0,255)} \\ d_{(1,0)} & d_{(1,1)} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ d_{(255,0)} & d_{255,1} & \cdots & d_{(255,255)} \end{matrix}$$

The above basic formulas are used in the GSM model to process a large amount of sample network traffic with normal behaviors. The distance maps of normal behaviors for each group of network traffic are calculated by (8). Simply, let us consider one group of network traffic with $m$ normal packets inside. Thus, the distance maps of normal packets are: $D_1^{nor}, \ldots, D_m^{nor}$, and the averages and variances for all elements $(i, j)$ of the distance map are computed by the following (9) and (10).

$$d_{nor(i,j)} = \frac{1}{m} \sum_{k=1}^{m} d_{nor\ i,j\ ,k} \qquad (9)$$

$$\sigma_{nor(i,j)}^2 = \frac{1}{m} \sum_{k=1}^{m} \left( d_{nor\ i,j\ ,k} - d_{nor(i,j)} \right)^2 \qquad (10)$$

Where ($0 \le i, j \le 255$) and $d_{nor\ i,j\ ,k}$ is the $(i, j)$ element of distance maps $D_k^{nor}$. The $d_{nor(i,j)}$ and $\sigma_{nor(i,j)}^2$ are all kept in a model $M_{nor}$ for further evaluation.

In the attack recognition phase, an input network packet experiences the same preprocessing procedure to construct its Mahalanobis distance map

$$D^{obj} = \left[ d_{i,j}^{obj} \right]_{256 \times 256} \qquad (11)$$

Then, a calculation is conducted to estimate the Mahalanobis distance between two distributions of $D^{obj}$ and the model $M_{nor}$.

$$w = \sum_{i,j=0}^{255,255} \frac{(d_{obj\ i,j} - d_{nor\ i,j})^2}{\sigma_{nor(i,j)}^2} \qquad (12)$$

If the weight w is larger than a threshold, we determine that the input network packet is an intrusion

## V. EXPERIMENTAL ANALYSIS AND RESULTS

We tested GSAD model on the 1999 DARPA IDS data set [16, 21], which is considered as standard data set to evaluate intrusion detection systems. In our experiment we made assumption that the number of attacks is very small in contrast to number of normal traffic. We mainly considered inbound TCP traffic only. The experiment has been done to identify crashiis attack, back attack, and mailbomb attack using 150 bytes of packet payload.

### A. Analysis and Result

The 1999 DARPA IDS data set was collected at MIT Lincoln Labs to evaluate intrusion systems. Entire network traffic was recorded in tcpdump format. The data set consists of three weeks of training of training data and two weeks of testing data. In the training data there are two weeks of attack-free data and one week of data with labelled attacks. These attacks are grouped into five classes as scan or probe, DoS, R2L, U2R and data.

In this experiment we used the inside network traffic data (week 1, week2 and week 3) which was captured between the router and the victims. We use wireshark for payload analysis and apply some filters based on payload length of 150 bytes, and for HTTP and SMTP service inbound TCP traffic.

We trained the GSAD model on the DARPA dataset using week1 and week 3 (attack free), then evaluate the model on week 2, which contains 43 instances of 15 different attacks. Test has been done on three types of attacks, crashiis, back and mailbomb. For port 80, the attacks are often malformed HTTP requests and are very different from normal requests. For instance, crashiis sends request "GET ..//..",apache2 sends request with a lot of repeated "User-Agent:sioux\r\n", back sends an HTTP request "GET //////////...." with more than 6000 slashes, which causes some versions of Apache web server to consume excessive CPU time, and for port 25, the attack mailbomb floods a user with thousands of junk emails. It is easy to identify these attacks using GSAD model and model shows a great difference in

the behaviour of these attacks with respect to the behaviour of normal network traffic for these services.

Fig. 2 (a) and (c) show the attack free and attack character relative frequencies, Fig. 2 (b) and (d) show the attack free and attack Mahalanobis Distance Map for crashiis attack. Fig. 3 (a) and (c) show the attack free and attack character relative frequencies, Fig. 3 (b) and (d) show the attack free and attack Mahalanobis Distance Map for back attack.

From Fig. 2 and Fig. 3, we can see that the character relative frequency and Mahalanobis Distance Map of the attack packets are very different from the normal packets', which can provide strong evidences to distinguish attacks from normal packets. The character relative frequencies of attack packets in both figures reveal the behaviours of

crashiis and back attack, which are different. For the crashiis attack, the "." character has the highest frequency and the other characters share even frequencies. Relatively, the statistical tendency of back attack is totally different and it is perfect match with the signature. Around 98 per cent of characters in the attack packets are "/".

Simultaneously, these experimental results illustrate the good performance of our GSAD model in detecting crashiis attack, back attack and mailbomb attack. That is clearly to be discovered from the geometrical structure models which explain the correlation among 256 ASCII characters. Both the behaviour models pairs in Figure 2 and 3 express dissimilar states between the attack free and attack packet. It can be taken as the sign to determine an intrusion.
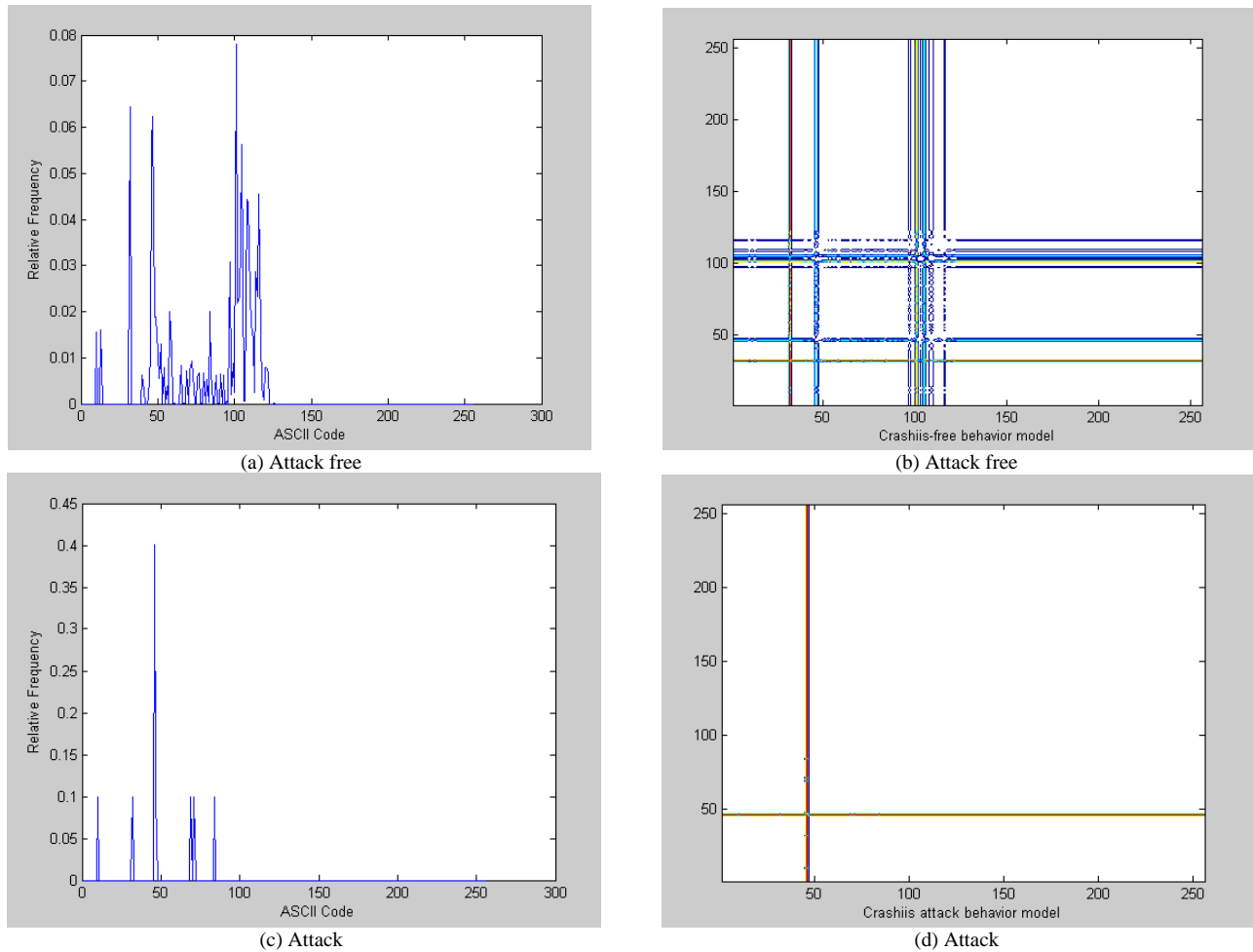


(a) Attack free      (b) Attack free

(c) Attack      (d) Attack

Figure 2.   Relative frequencies of characters (a) (c) and mahalanobis distance map (b) (d) for Crashiis attack.

(a) Attack free

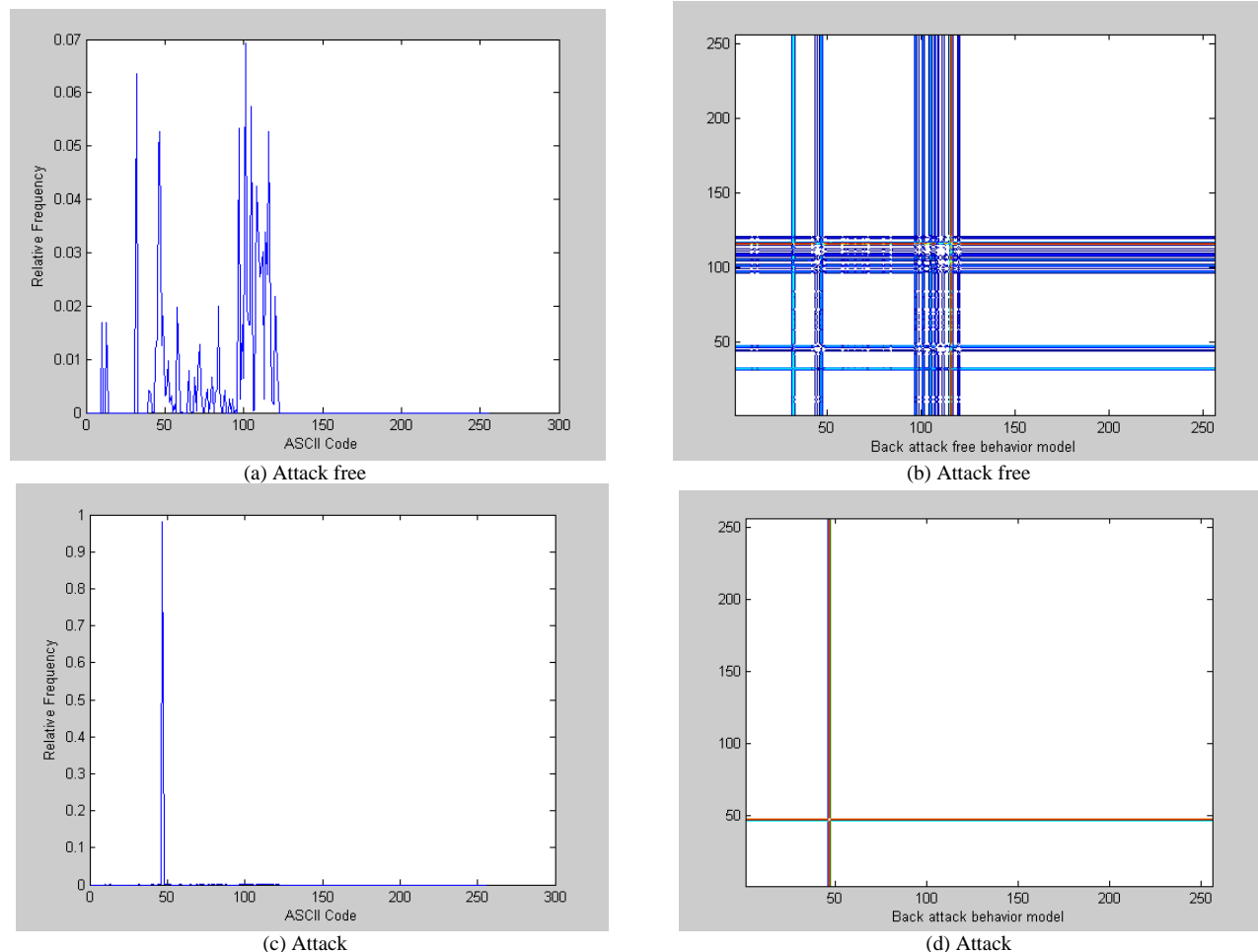(b) Attack free

(c) Attack

(d) Attack

Figure 3.   Relative frequencies of characters (a) (c) and mahalanobis distance map (b) (d) for Back attack

## VI.   CONCLUSIONS

In this paper we present an approach for network intrusion detection based on geometrical structure of anomaly payload. The key features are to compute byte distribution model and geometrical structure model for normal traffic, conditioned to service type, and payload length. The weight factor is used to compare the similarity between the new incoming packet's payload and its corresponding model using mahalanobis distance map (MDM). If the weight is greater than the threshold, the incoming packet will be considered as an attack packet. The experiments done for crahiis attack, back attack and mailbomb attack show good results.

In our future work we aim to evaluate the performance of our model and validate our results. We also plan to test this model on 1999 DARPA IDS dataset for variable length payload, protocols and services.

## REFERENCES

[1]    CERT, "CERT Statistics", http://www.cert.org/stats/#notes, 2007.

[2]    J. P. Anderson, "Computer Security Threat monitoring and surveillance", Technical report, JP Anderson Co., Ft. Washington, Pennsylvania, Apr 1980.

[3]    P. Ning and S. Jajodia, "Intrusion Detection Techniques in H. Bidgoli (Ed.)", The Internet Encyclopedia: John Wiley & Sons, 2003.

[4]    A. Patcha and J. M. Park, "An overview of anomaly detection techniques: existing solutions and latest technological trends", Computer networks, 2007.

[5]    Snort: The open source network intrusion detection system

[6]    P. Uppuluri and R. Sekar, "Experiences with Specification-Based Intrusion Detection System", In Recent Advances in Intrusion Detection: 4th International Symposium, RAID 2001 Davis, CA, USA, October 10-12, 2001, Proceedings 2001, pp. 172.

[7]    VernPaxson, "Bro: a system for detecting network intruders in real-time", Computer Networks (Amsterdam, Netherlands: 1999), 31(23-24):2435–2463, 1999V Paxson,

[8]    J.       Hoagland,       SPADE,       Silican       Defence, http://www.silicondefence.com/software/spice, 2000.

[9]    H. S. Javits and A. Valdes, "The NIDES statistical component: Description and justification", Technical report, SRI International, computer Science Laboratory, 1993.

[10]   M. Mahoney and P. Chan, "Learning non stationary models of normal network traffic for detecting novel attacks", In Proc. SIGKDD 2002, pp. 376–385, 2002.

[11] M. Mahoney, "Network traffic anomaly detection based on packet bytes", In Proc. ACM-SAC, Melbourne FL, pp. 346– 350, 2003.

[12] C. Taylor and J. Alves-foss, "NATE-Network Analysis of Anomaly Traffic Events, A Low-Cost approach", New Security Paradigms Workshop, 2001.

[13] Christopher Krgel, Thomas Toth, and Engin Kirda, "Service specific anomaly detection for network intrusion detection", In Proceedings of the 2002 ACM symposium on Applied computing, pp. 201–208, 2002.

[14] Ke Wang and S. Stolfo, "Anomalous payload-based network intrusion detection", In Recent Advances in Intrusion Detection, RAID, pages 203–222, September 2004.

[15] R. Lippmann, "The 1999 DARPA offline intrusion detection evaluation", In recent Advances in Intrusion Detection. Third International Workshop, RAID 2000, 2-4 Oct. 200, Toulouse, France (Berlin, Germany, 2000), H.Debar, L. Me, and S. Wu, Eds., Springer-Verlag, pp. 162-182.

[16] R. Lippmann, J. Haines, K. Dass, "Analysis and results Of the 1999 DARPA offline Intrusion detection evaluation", In Computer Networks, 34(4), pp. 579-595, 2000.

[17] J McHugh, "The 1998 Lincoln Laboratory IDS evaluation-a critique", In Recent Advances in Intrusion Detection, Third International Workshop, RAID 2000, 2-4 Oct. 200, Toulouse, France (Berlin, Germany, 2000), H.Debar, L. Me, and S. Wu, Eds., Springer-Verlag, pp. 145-161.

[18] TCPDUMP and LIBPCAP Project: http://www.tcpdump.org/

[19] X. He, J. Li, Y. Chen and W. Jia, "Local Binary Patterns with Mahalanobis Distance Maps for Human Detection", IEEE Congress on Image and Signal Processing, pp. 520-524, 2008

[20] Akira Utsumi and Nobuji Tetsutani, "Human Detection using Geometrical Pixel Value Structures", In Proceeding of 5[th] International Conference on Automatic Face and Gesture Recognition (FGR '02), pp. 34-39, 2002.

[21] http://www.ll.mit.edu/IST/ideval/dex.html