

Multivariate Correlation Analysis Technique Based on Euclidean Distance Map for Network Traffic Characterization

Zhiyuan Tan^{1,2}, Aruna Jamdagni^{1,2}, Xiangjian He¹, Priyadarsi Nanda¹, and
Ren Ping Liu²

¹ Research Centre for Innovation in IT Services and Applications (iNEXT)
University of Technology, Sydney, Broadway 2007, Australia

² CSIRO Marsfield, Australia

{Zhiyuan.Tan, Aruna.Jamdagni}@student.uts.edu.au,
{Xiangjian.He, Priyadarsi.Nanda}@uts.edu.au,
ren.liu@csiro.au

Abstract. The quality of feature has significant impact on the performance of detection techniques used for Denial-of-Service (DoS) attack. The features that fail to provide accurate characterization for network traffic records make the techniques suffer from low accuracy in detection. Although researches have been conducted and attempted to overcome this problem, there are some constraints in these works. In this paper, we propose a technique based on Euclidean Distance Map (EDM) for optimal feature extraction. The proposed technique runs analysis on original feature space (first-order statistics) and extracts the multivariate correlations between the first-order statistics. The extracted multivariate correlations, namely second-order statistics, preserve significant discriminative information for accurate characterizations of network traffic records, and these multivariate correlations can be the high-quality potential features for DoS attack detection. The effectiveness of the proposed technique is evaluated using KDD CUP 99 dataset and experimental analysis shows encouraging results.

Keywords: Euclidean Distance Map, Multivariate Correlations, Second-order Statistics, Characterization, Denial-of-Service Attack

1 Introduction

The growing number of network intrusive activities poses a serious threat to the reliability of network services. Businesses and individuals are suffering from these malicious interceptions. Billions of dollars loss has been recorded over the past few years [1].

As one of the major network intrusive activities, Denial-of-Service (DoS) attack receives much attention due to its continuous growth and serious impact on the Internet. A victim, such as host, router or entire network, can be overwhelmed by a DoS attack using imposed computationally intensive tasks, using

exploitation of system vulnerability or using floods with a huge amount of useless packets. The victim is then temporarily unavailable for the outside networks from a few minutes to even several days. The availability of network services is severely degraded by this type of network intrusive activities, thus effective detection mechanisms for DoS attack are highly required.

However, the work that has been done so far is still far away from being perfect. Currently, the commercially used DoS attack detection systems are mainly dominated by signature-based detection techniques [2][3]. In spite of having high detection rates to the known attacks and low false positive rates, signature-based techniques are easily evaded by new attacks and even variants of the existing attacks.

Therefore, research community has started to explore a way to achieve novelty-tolerant detection systems and developed the concept of anomaly-based detection [4][5]. The idea of anomaly-based detection is that network intrusions exhibit significantly different behaviors than the normal network activities [6], and any significant deviation from the normal network behaviors is identified as an intrusion. To implement this idea, various techniques, such as clustering [7][8], neural network [9][10], pattern recognition [11][12], support vector machine [13], nearest neighbor [14] and statistical detection techniques [15][16][17] have been used to establish anomaly-based detection systems. However, some of the techniques [7][10][14] suffer from relatively low accuracy in the task of attack detection, though they show encouraging results in other tasks.

The aforementioned discussed problem is partly raised by the low quality features which fail to provide sufficient discriminative power for correct traffic discrimination. To address the problem, this paper presents the Euclidean Distance Map (EDM) to analyze the original feature space (first-order statistics) and extracts the multivariate correlations between the first-order statistics. These multivariate correlations contain significant discriminative information and play key roles in detection accuracy. The occurrence of network intrusions cause changes to these multivariate correlations so that the changes can be used as metrics for identifying intrusive activities. Moreover, various types of traffic namely normal traffic and attack traffic can be easily and accurately characterized by using the proposed EDM-based analysis technique. By looking into the patterns, i.e. EDMs, the intrusive activities can be clearly identified and differentiated from each other.

The rest of this paper is organized as follows. Section 2 provides current work related to our research. Section 3 details the EDM-based multivariate correlation analysis technique. Section 4 evaluates the performance of the proposed approach in pattern extraction of DoS attack and makes some discussions. Finally, conclusions are drawn and future work is given in Section 5.

2 Related Work

Multivariate correlations are second-order statistics generated to reveal the relations between or among the original features, i.e. first-order statistics. This

correlative information provides important discriminative power and is proven to be more effective for object clustering and classification.

Recently, researchers have investigated to explore extracting effective multivariate correlations for DoS attack detection, and different techniques have been proposed. A team of researchers from the Hong Kong Polytechnic University [17] proposed a covariance matrix based approach to mine the multivariate correlations for sequential samples and reveal characteristics of different classes of traffic records. This idea is later adopted by Travallae et al. [18], who further used the Principal Component Analysis (PCA) to reduce the redundant information contained in the original feature space from statistical point of view. Such performance of the covariance matrix based approach can be refined by using the appropriately selected features.

Apart from the above discussed statistical-based approaches, new solutions proposed lately tend to consider the geometrical structure of the features. Jamdagni et al. [11] developed a distance measure based correlation extraction approach, in which the Mahalanobis Distance (MD) is used to measure the weighted distance between each pair of features extracted from network traffic packet payload. In addition, Tsai and Lin [19] estimated the sizes of triangle areas constructed by any signal observed data object and any two centroids of distinct clusters. New data formed by the areas of these triangles is used as a new feature space.

Although these approaches introduce some interesting concepts for multivariate correlation extraction and show their abilities in extracting discriminative power, they have some weaknesses and drawbacks. On one hand, techniques [17] [18], such as the covariance matrix, will not work under the situation where an attack linearly changes all monitored features and is vulnerable to mix-traffic containing both normal and attack traffic. On the other hand, the distance measure based techniques either suffer from high computation complexity [11] or are dependent on prior knowledge of both normal and attack traffic that causes wrong characterization of any novelty [19].

Different from the above discussed techniques, our EDM-based analysis technique is independent on prior knowledge of the features of different classes and withstands the issue of linearly changes all monitored features. It is robust to mix-traffic and does not rely on the volume of network traffic. More importantly, it is computationally efficient.

3 EDM Based Multivariate Correlation Analysis

The behavior of network traffic is reflected by its statistical properties. DoS attempts to exhaust a victims resources, and its traffic behaves differently from the normal network traffic. Therefore, the statistical properties can be used to reveal the difference. To present the statistical properties, we propose a multivariate correlation analysis approach which employs Euclidean distance for extracting correlative information (named inner correlation) from the original feature space of an observed data object. This approach shares similar properties with the one

developed in [21]. The detail of the proposed approach is given in the following section.

3.1 Multivariate Correlation Extraction

Given an arbitrary dataset $X^T = [x_1^T \ x_2^T \ \cdots \ x_n^T]$, where $x_i^T = [f_1^i \ f_2^i \ \cdots \ f_m^i]$ ($1 \leq i \leq n$) represents the i^{th} m -dimensional traffic record. The dataset can be represented in detail as

$$X = \begin{bmatrix} f_1^1 & f_2^1 & \cdots & f_m^1 \\ f_1^2 & f_2^2 & \cdots & f_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ f_1^n & f_2^n & \cdots & f_m^n \end{bmatrix}, \quad (1)$$

where f_l^i is the value of the l^{th} feature in the i^{th} traffic record, l and i are varying from 1 to m and from 1 to n respectively.

In order to further explore the inner correlations of the i^{th} traffic record on a multi-dimensional space, the record x_i^T is first transformed into a new m -by- m feature matrix x_i' by simply multiplying an m -by- m identity matrix I as shown in Equation (2).

$$x_i^T I = x_i' = \begin{bmatrix} f_1^i & 0 & \cdots & 0 \\ 0 & f_2^i & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & f_m^i \end{bmatrix}_{m \times m}. \quad (2)$$

The elements on the diagonal of the matrix x_i' are the features of the record x_i^T . Each column of the matrix x_i' is a new m -dimensional feature vector denoted by

$$F_j^{i,T} = [F_{j,1}^i \ F_{j,2}^i \ \cdots \ F_{j,m}^i], \quad (3)$$

where $F_{j,p}^i = 0$ if $j \neq p$ and $F_{j,p}^i = f_j^i$ if $j = p$. The parameters satisfy the conditions of $1 \leq i \leq n$, $1 \leq j \leq m$ and $1 \leq p \leq m$. Thus, the m -by- m feature matrix can be rewritten as

$$x_i' = [F_1^i \ F_2^i \ \cdots \ F_m^i]. \quad (4)$$

Once the transformation is finished, we can apply the Euclidean distance to extract the correlation between the feature vectors j and k in the matrix x_i' , which can be defined as

$$ED_{j,k}^i = \sqrt{(F_j^i - F_k^i)^T (F_j^i - F_k^i)}. \quad (5)$$

where $1 \leq i \leq n$, $1 \leq j \leq m$ and $1 \leq k \leq m$. Therefore, the correlations between features in the traffic record x_i^T defined by a Euclidean Distance Map (EDM) are given below.

$$EDM^i = \begin{bmatrix} ED_{1,1}^i & ED_{1,2}^i & \cdots & ED_{1,m}^i \\ ED_{2,1}^i & ED_{2,2}^i & \cdots & ED_{2,m}^i \\ \vdots & \vdots & \ddots & \vdots \\ ED_{m,1}^i & ED_{m,2}^i & \cdots & ED_{m,m}^i \end{bmatrix}. \quad (6)$$

For the dataset, its inner correlations can be represented by Equation (7).

$$EDM_X = (EDM^1 EDM^2 \dots EDM^i \dots EDM^n) \quad (7)$$

In comparison with [21], which considers only the lower triangle of the Euclidean distance map, the multivariate correlation analysis approach proposed in this paper takes the entire map into account. The EDM is then employed as a means of network traffic characterization and a visualization tool to reveal the patterns of various traffic.

3.2 Discussions

By making use of the multivariate correlations, various types of network traffic can be clearly characterized. Additionally, the distance measure facilitates our analysis to withstand the issue of linear change for all features.

The two primary advantages of the proposed EDM-based analysis technique are supported by two underlying mathematical structures. They are the transformed traffic record matrix and the Euclidean distance. These two mathematical tools help solve the dilemmas caused by the occurrence of two distinct pairs of features having the same distance on one-dimensional space and the linear change of all features.

Assume that there are three pairs of features $A(1, 2)$, $B(4, 8)$ and $C(9, 10)$ shown in Fig. 1. In this case, the ratio between the two features of A and that between the two features of B have the same value which is equal to $1/2 = 4/8 = 0.5$, and the one-dimensional linear distance between the two features of A and the distance between the two features of C have the same value which is equal to $|1 - 2| = |9 - 10| = 1$. For the conventional techniques, using ratio or one-dimensional linear distance cannot differentiate the points which coincidentally have the same value.

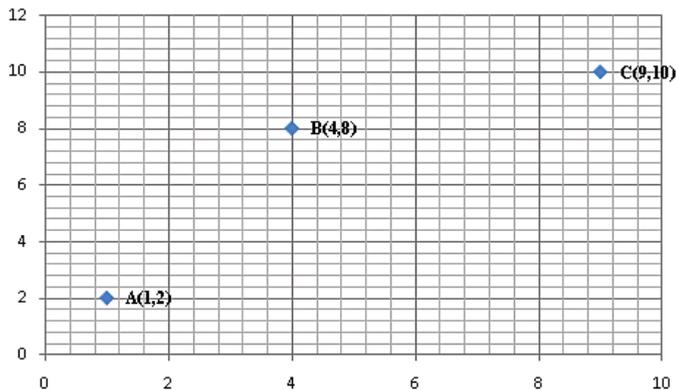


Fig. 1: Three pairs of features.

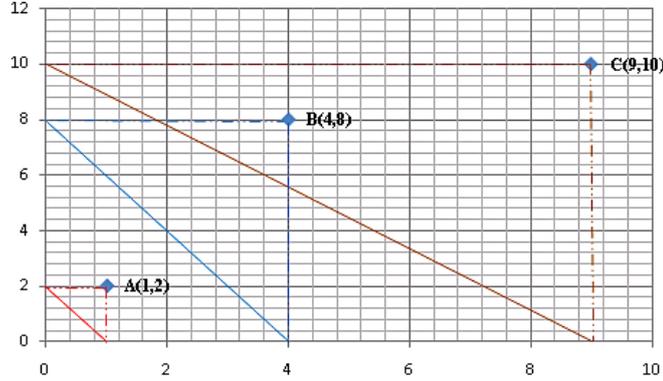


Fig. 2: Distances between the features with three pairs.

However, the proposed EDM-based analysis technique can successfully withstand these difficulties. By applying the technique, the distances between the features in the A , B and C are $\sqrt{(1-0)^2 + (0-2)^2} = \sqrt{5}$, $\sqrt{(4-0)^2 + (0-8)^2} = \sqrt{80}$ and $\sqrt{(9-0)^2 + (0-10)^2} = \sqrt{181}$ respectively. Therefore, the proposed technique can properly characterize network traffic records and is not affected by the linear change of all features.

4 Experimental Results and Analysis

The performance of our proposed EDM-based analysis technique for characterizing normal and DoS attack traffic records is evaluated using KDD CUP 99 dataset [1]. Although the dataset is not without criticism [20], it is the only public dataset with labeled attack samples. Moreover, many research works have been evaluated using this dataset as well.

4.1 Experimental Data

The 10 percent labeled dataset of the KDD CUP 99 dataset is involved in our experimentation. Six different types of DoS attacks (including Teardrop, Smurf, Pod, Neptune, Land and Back attacks) and normal network traffic from the labeled dataset are used for evaluation. The DoS attacks launch their malicious activities by exploiting three widely used network protocols respectively. Neptune, Land and Back attacks make use of TCP protocol. Teardrop sends its attack payload over UDP. The payloads of Smurf and Pod attack are carried by ICMP packets.

4.2 Results and Analysis

According to the working mechanisms discussed in Section 1, DoS attacks are expected to present different behaviors to that of the normal traffic. Thus, the

EDMs of DoS attacks should be different from the EDM of normal traffic. If significant differences can be found from these maps, the performance of the proposed technique can be proven to be good in extracting discriminative power from the first-order features of the different types of network traffic.

To demonstrate how EDM presents the correlations between the first-order features, the EDMs of normal and attack traffic records generated using 32 numerical features are given in this subsection. As shown in Fig. 3, the EDM of normal TCP traffic record is a symmetric matrix and the values of the elements along its diagonal from top left hand side to bottom right hand side are all equal to zeros. This is because the Euclidean distance measure is insensitive to the orientation of a straight line formed by any two objects in the Cartesian coordinate system, and the distance from a feature vector to itself is always zero. In other words, the distance from object D to object E is equivalent to the distance from object E to object D , and if object D and object E are the same object, then their distance is zero.

Although we can directly compare the raw EDMs to confirm the differences, it is a time-consuming task. In order to offer friendly visualization for the raw EDMs, we convert them into color images. The images of normal TCP traffic record and Neptune, Land and Back attack records are given in Fig. 4, and the

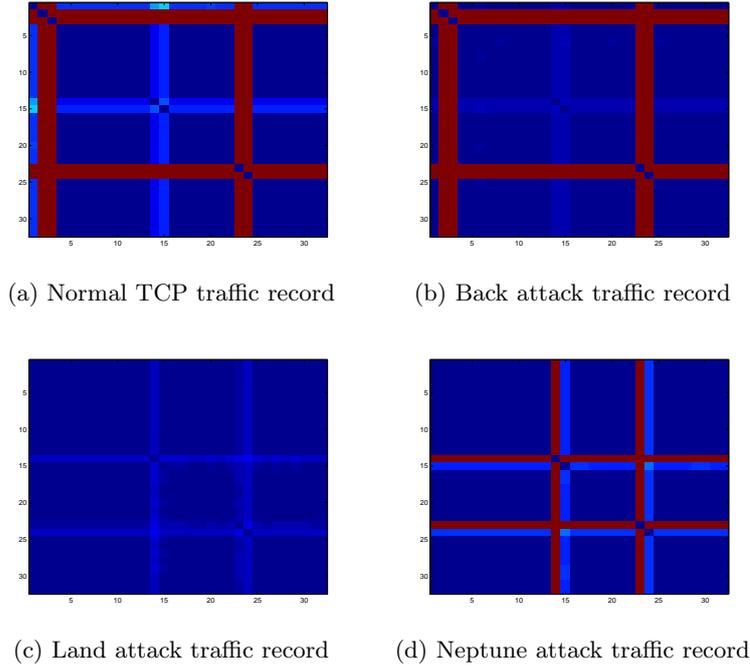


Fig. 4: Images of the EDMs of normal TCP traffic record and Back, Land and Neptune attack records.

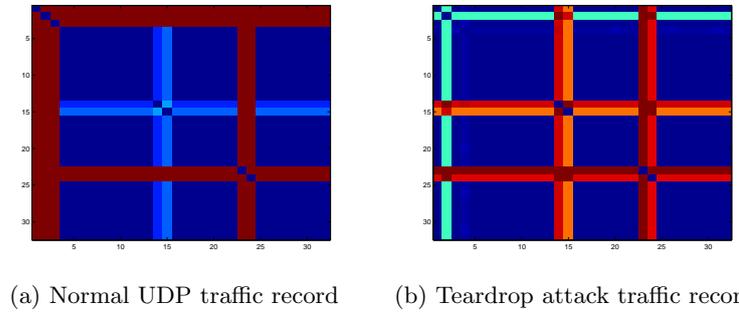


Fig. 5: Images of EDMs of UDP traffic record and Teardrop attack record.

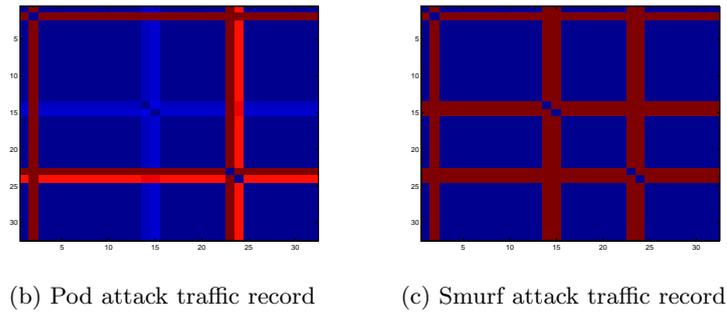
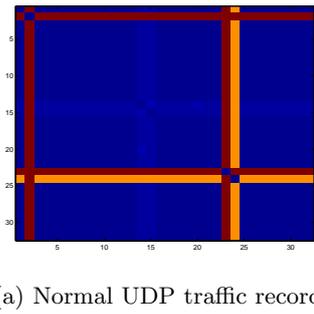


Fig. 6: Images of EDMs of ICMP traffic record and Pod and Smurf attack records.

images of normal UDP traffic record and Teardrop attack record are shown in Fig. 5. Finally, Fig. 6 presents the images of normal ICMP traffic record and Smurf and Pod attack records.

As can be seen from Fig. 4(a), the image represents the visualized pattern of the EDM of normal TCP traffic record. The color of an image point stands for the value of an element on the EDM. The lighter and warmer the color is, the greater value the element has. In other words, the darkest cold blue color areas on the image are the lowest value areas on the EDM, and conversely the lightest warm red color areas on the image are the highest value areas on the EDM. Figs. 4(b), (c) and (d) visualize the EDMs of Back, Land and Neptune attack records in the same manner respectively. The images of the attack EDMs show clear differences from the EDM of the normal TCP traffic record.

Similarly, the images of the EDMs of UDP traffic record and Teardrop attack record are exhibited in Fig. 5. The image of Teardrop shows apparent dissimilarity to the image of normal UDP traffic. In addition, the images in Fig. 6 reveal that the behaviors of ICMP-based attacks, namely Pod and Smurf attacks, are away from the normal ICMP traffic as well.

The above experimental results demonstrate that our proposed EDM-based technique achieves promising performance in characterizing various network traffic records. Our experimental results also suggest that, by taking advantage of the retained significant discriminative power, utilization of the generated multivariate correlations can improve the performance of DoS attack detection system. Moreover, by looking into the images, we can easily identify the patterns of the different traffic records. Therefore, the proposed EDM-based technique can be further applied to creating the statistical signatures of network intrusions. For detailed comparisons of effectiveness in intrusion detection with some other existing approaches, please refer to [21].

5 Conclusions and Future Work

This paper has proposed a multivariate correlation analysis approach based on Euclidean distance to extract the multivariate correlations (second-order statistics) of network traffic records. This proposed approach can better exhibit the network traffic behaviors. We have evaluated the analysis approach on the records of normal and DoS attack traffic from the KDD CUP 99 dataset. The results illustrate that these second-order statistics can clearly reveal the correlations between the first-order statistics and accurately characterize the various types of traffic records. In future, we will further evaluate the proposed technique on the task of DoS attack detection using Support Vector Data Description (SVDD) technique, which is believed to be more promising in one-class classification than SVM and NN techniques. We may also extend our research to the characterization of temporal information.

References

1. Cheng, J., Hatzis, C., Hayashi, H., Krogel, M.A., Morishita, S., Page, D., Sese, J.: KDD Cup 2001 Report. ACM SIGKDD Explorations Newsletter 3, 47-64 (2002)
2. Paxson, V.: Bro: A System for Detecting Network Intruders in Real-time. *Computer Networks* 31, 2435-2463 (1999)
3. Roesch, M.: Snort-lightweight Intrusion Detection for Networks. Proceedings of the 13th USENIX Conference on System Administration, pp. 229-238. USENIX, Seattle, Washington (1999)
4. Patcha, A., Park, J.M.: An Overview of Anomaly Detection Techniques: Existing Solutions and Latest Technological Trends. *Computer Networks* 51, 3448-3470 (2007)
5. Garcia-Teodoro, P., Diaz-Verdejo, J., Macia-Fernandez, G., Vazquez, E.: Anomaly-Based Network Intrusion Detection: Techniques, Systems and Challenges. *Computers & Security* 28, 18-28 (2009)
6. Denning, D.E.: An Intrusion-detection Model. *IEEE Transactions on Software Engineering* 222-232 (1987)
7. Jin, C., Wang, H., Shin, K.G.: Hop-count Filtering: An Effective Defense Against Spoofed DDoS Traffic. The 10th ACM Conference on Computer and Communications Security. pp. 30-41. ACM (2003)
8. Lee, K., Kim, J., Kwon, K.H., Han, Y., Kim, S.: DDoS Attack Detection Method Using Cluster Analysis. *Expert Systems with Applications* 34, 1659-1665 (2008)
9. Amini, M., Jalili, R., Shahriari, H.R.: RT-UNNID: A Practical Solution to Real-time Network-based Intrusion Detection Using Unsupervised Neural Networks. *Computers & Security* 25, 459-468 (2006)
10. Wang, G., Hao, J., Ma, J., Huang, L.: A new approach to intrusion detection using Artificial Neural Networks and fuzzy clustering. *Expert Systems with Applications* 37, 6225-6232 (2010)
11. Jamdagni, A., Tan, Z., Nanda, P., He, X., Liu, R.P.: Intrusion Detection Using GSAD Model for HTTP Traffic on Web Services. The 6th International Wireless Communications and Mobile Computing Conference. pp. 1193-1197. ACM (2010)
12. Tan, Z., Jamdagni, A., He, X., Nanda, P., Liu, R., Jia, W., Yeh, W.: A Two-tier System for Web Attack Detection Using Linear Discriminant Method. The 12th International Conference on Information and Communications Security 459-471 (2010)
13. Fugate, M., Gattiker, J.R.: Computer Intrusion Detection with Classification and Anomaly Detection Using SVMs. *International Journal of Pattern Recognition and Artificial Intelligence* 17, 441-458 (2003)
14. Lane, T., Brodley, C.E.: Temporal Sequence Learning and Data Reduction for Anomaly Detection. *ACM Transactions on Information and System Security (TISSEC)* 2, 295-331 (1999)
15. Ye, N., Emran, S.M., Chen, Q., Vilbert, S.: Multivariate Statistical Analysis of Audit Trails for Host-based Intrusion Detection. *IEEE Transactions on Computers* 810-820 (2002)
16. Manikopoulos, C., Papavassiliou, S.: Network Intrusion and Fault Detection: A Statistical Anomaly Approach. *Communications Magazine, IEEE* 40, 76-82 (2002)
17. Jin, S., Yeung, D.S., Wang, X.: Network Intrusion Detection in Covariance Feature Space. *Pattern Recognition* 40, 2185-2197 (2007)
18. Tavallaee, M., Lu, W., Iqbal, S.A., Ghorbani, A.A.: A Novel Covariance Matrix Based Approach for Detecting Network Anomalies. The Communication Networks and Services Research Conference, pp. 75-81. IEEE, (2008)

19. Tsai, C.F., Lin, C.Y.: A Triangle Area Based Nearest Neighbors Approach to Intrusion Detection. *Pattern Recognition* 43, 222-229 (2010)
20. Tavallae, M., Bagheri, E., Lu, W., Ghorbani, A.A.: A Detailed Analysis of the KDD Cup 99 Data Set. (2009)
21. Tan, Z., Jamdagni, A., He, X., Nanda, P., Liu, R.: Denial-of-Service Attack Detection Based on Multivariate Correlation Analysis. *The 18th International Conference on Neural Information Processing*. (2011) (Accpeted for Publication)