# Denial-of-Service Attack Detection Based on Multivariate Correlation Analysis

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

[1] Centre for Innovation in IT Services and Applications (iNEXT),
University of Technology, Sydney, 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 reliability and availability of network services are being threatened by the growing number of Denial-of-Service (DoS) attacks. Effective mechanisms for DoS attack detection are demanded. Therefore, we propose a multivariate correlation analysis approach to investigate and extract second-order statistics from the observed network traffic records. These second-order statistics extracted by the proposed analysis approach can provide important correlative information hiding among the features. By making use of this hidden information, the detection accuracy can be significantly enhanced. The effectiveness of the proposed multivariate correlation analysis approach is evaluated on the KDD CUP 99 dataset. The evaluation shows encouraging results with average 99.96% detection rate and 2.08% false positive rate. Comparisons also show that our multivariate correlation analysis based detection approach outperforms some other current researches in detecting DoS attacks.

**Keywords:** Denial-of-Service Attack, Euclidean Distance Map, Multivariate Correlations, Anomaly Detection.

## 1    Introduction

Network security has received public concerns with the rapid growth and the prevalence of interconnection among computer systems. It is now under spotlight due to the emergence of more sophisticated attack techniques and easy-access user-friendly attack tools which facilitates any person to easily launch network attacks with little programming knowledge. According to studies [1] and [2], there have been 10,000 new viruses or variant of existing viruses recorded in the year of 2004, and billions of dollars loss has been caused by Denial-of-Service (DoS) attacks over the past few years.

The intention of a DoS attack is to deliberately prevent a victim, such as host, router or entire network, from being accessible or being capable of receiving normal services from the Internet. The availability of network services is seriously threatened

by the continuously increasing number of DoS attacks. Thus effective mechanisms for DoS attack detection are highly demanded.

To maintain the reliability and the availability of network services, research community and industry sector have put a lot of efforts to the development of intrusion detection techniques. As one of the powerful network intrusions, DoS attack has been carefully studied in the intrusion detection research over the last decade. Generally, network intrusion detection can be grouped into two main categories, namely signature-based detection [3] and anomaly-based detection [4]. Benefiting from the principal of detection, which monitors and flags any network activity presenting significant deviation from their normal profiles as a suspicious, anomaly-based detection techniques show more advanced in detecting zero-day intrusions [5].

Therefore, recent works in DoS attack detection mainly focus on anomaly-based techniques, and various detection techniques have been proposed. For example, clustering [6] [7], neural network [8] [9], pattern recognition [10], support vector machine [11], nearest neighbor [12] and statistical detection techniques [13] [14] [15].

However, some of these proposed techniques often suffer high false positive rate since the dependencies and correlations of the features are intrinsically neglected [16]. The other techniques are either invalid to flooding-based DoS attacks [10] or incapable of identifying individual attack packets from a group of samples [15].

To address the aforementioned problems, a Euclidean Distance Map (EDM) based multivariate correlation (second-order statistics) analysis approach is proposed in this paper to discover the relations among features within the observed data objects. Significant changes of these relations indicate occurrences of intrusions. Owing to the computational simplify of Euclidean distance and the extracted valuable correlative information, application of the multivariate correlation analysis makes the DoS attack detection more effective and efficient. It achieves high detection accuracy while retaining a low false positive rate. Moreover, benefiting from the principal of anomaly detection, our DoS attack detection approach is independent on prior knowledge of attack and is capable of detecting both known and unknown DoS attacks.

The rest of this paper is organized as follows. Section 2 provides current work related to our research. Section 3 proposes a novel multivariate correlation analysis approach. Section 4 presents a detailed description on the applications of the proposed multivariate correlation analysis approach in DoS attack detection. Section 5 shows the evaluation results of the proposed approach on KDD CUP 99 dataset and makes some analysis. Finally, conclusions and future work are given in Section 6.


## 2    Related Work

Owing to the advantage in detecting unknown attacks, anomaly intrusion detection mechanism has captured the major attention from research community. Researchers focused on studies of sequential change-point detection based statistical DoS attack detection approaches [17] [18] in the early 2000s. The approaches make use of the abrupt change occurring in the observed sequential data, such as Management Information Base (MIB) variables and statistics between the number of SYN packets and the number of FIN or SYN/ACK packets. They have been proven effective in

detecting any abrupt change in network traffic. However, on one hand, in Thottan and Ji's approach [17] the operator matrix may need to frequently update its feature set to cover the emergent attacks, and the approach suffers time granularity problem. On the other hand, the approach proposed by Wang et al. [18] only targets on SYN flooding attacks, and its performance is affected by the passive RST packets. Moreover, the two approaches consider only the first-order statistics and ignore the correlative information which is important to detection accuracy.

Recently, the intrusion detection research community started recognizing the importance of the second-order statistics of monitored-network features. Several researches have been conducted to explore the use of the second-order statistics in DoS attack detection. A team of researchers from the Hong Kong Polytechnic University [15] employed the covariance matrices of the sequential samples and proposed a threshold based detection approach to detecting various types of DoS attacks. Travallaee et al. [19] applied Covariance Matrix Sign (CMS) for DoS attack detection. These approaches achieve encouraging detection rates. However, they still suffer from comparative high false positive rates and do not work under the situation where an attack linearly changes all monitored features. In addition, the approaches can only label a group of observed samples as normal or attacks, and cannot identify individual attack packets from the crowd.

Our work also makes use of the idea of change and the second-order statistics. We investigate the change of the correlations between features which are the second-order statistics of the features in a single observed data object. This makes our approach more advanced in detection accuracy and the ability of labeling individual attack packets.

## 3 Multivariate Correlation Analysis

The behavior of network traffic is reflected by its statistical properties. DoS attack is a type of intrusions attempting to exhaust a victim's resource, and its traffic behaves different from the normal network traffic. Therefore, the statistical properties can be used to reveal the difference. To well present the statistical properties, we propose a novel multivariate correlation analysis approach which employs Euclidean distance for extracting correlative information (named inner correlation) among the features within an observed data object. The detail is given in the following.

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

$$X = \begin{bmatrix} f_1^1 & f_2^1 & \cdots & f_m^1 \\ f_1^2 & f_2^2 & \cdots & f_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ f_1^m & f_2^m & \cdots & f_m^m \end{bmatrix} \tag{1}$$

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 (2).

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

The elements on the diagonal of the matrix $x_i'$ of 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

$$F_j^{i\,T} = [F'^i_{j,1} \ F'^i_{j,2} \cdots F'^i_{j,m}] \ . \tag{3}$$

where $F'^i_{j,p} = 0$ if $j \neq p$ and $F'^i_{j,p} = f_j^i$ 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 $x_i'$ can be rewritten as (4).

$$x_i' = [F_1^i \ F_2^i \cdots F_m^i] \ . \tag{4}$$

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

$$ED_{j,k}^i = \sqrt{\left(F_j^i - F_k^i\right)^T \left(F_j^i - F_k^i\right)}, \tag{5}$$

where $1 \leq i \leq n, 1 \leq j \leq m$ and $1 \leq k \leq m$. However in practice, (5) can be simplified and rewritten as (6) in order to reduce computational complexity.

$$ED_{j,k}^i = \begin{cases} \sqrt{\left(f_j^i - 0\right)^2 + \left(0 - f_k^i\right)^2}, & j \neq k \\ 0, & j = k \end{cases} \tag{6}$$

Therefore, the correlations between features in the traffic record $x_i{}^T$ are defined by a Euclidean Distance Map (EDM) given below.

$$EDM^i = \begin{bmatrix} ED_{1,1}^i & ED_{1,2}^i & \cdots & ED_{1,m}^i \\ ED_{2,1}^i & ED_{2,2}^i & \cdots & ED_{2,m}^i \\ \vdots & \vdots & \ddots & \vdots \\ ED_{m,1}^i & ED_{m,2}^i & \cdots & ED_{m,m}^i \end{bmatrix}. \tag{7}$$

Since the EDM is a symmetric matrix (in which $ED_{j,k}^i = ED_{k,j}^i$) and there is a zero distance from a feature vector to itself ($ED_{j,k}^i = 0, if \ j = k$), the upper or the lower triangle of the matrix is sufficient to reveal the inner correlations. Hence, the EDM

can be simplified and converted into a new inner correlation vector containing only the lower triangle of the EDM as (8).

$$EDM_{lower}^{i}{}^{T} = \left[ED_{2,1}^{i}\ ED_{3,1}^{i}\ \cdots\ ED_{m,1}^{i}\ ED_{3,2}^{i}\ ED_{4,2}^{i}\ \cdots\ ED_{m,2}^{i}\ \cdots\ ED_{m,m-1}^{i}\right]. \tag{8}$$

For the dataset $X$, its inner correlations can be represented by (9).

$$EDM_{lower}{}^{T} = [EDM_{lower}^{1}{}^{T}\ EDM_{lower}^{2}{}^{T}\ \cdots\ EDM_{lower}^{i}{}^{T}\ \cdots\ EDM_{lower}^{n}{}^{T}]. \tag{9}$$

By making use of the inner correlations, the changes of network behavior caused by DoS attack can be clearly revealed. Additionally, the distance measure facilitates our analysis approach to withstand the issue of linear change of all features.

## 4 Multivariate Correlation Analysis Based Detection Approach

The objective of this paper is to develop a detection approach that is effective in detecting any known and unknown DoS attacks. Thus, the concept of anomaly-based IDS, which attempts to identify network intrusions by detecting any significant deviations from a profile generated using only normal traffic records in training phase, is the best fit to our problem.

### 4.1 Norm Profile Generation

In this work, the norm profile is first built through the density estimation of the Mahalanobis Distances (MDs) between observed normal traffic records and the expectation of the normal traffic record. To obtain the distribution of the MDs, two parameters are required to be determined. They are the mean $\mu$ and the standard deviation $\sigma$ of the distances.

Assume that there is a set of $g$ normal training traffic records, which is denoted by $EDM_{lower}^{normal}{}^{T} = [EDM_{lower}^{norml,1}{}^{T}\ EDM_{lower}^{normal,2}{}^{T}\ \cdots\ EDM_{lower}^{normal,g}{}^{T}]$, the parameters can be determined by using the equations shown below. The mean $\mu$ is defined as

$$\mu = \frac{1}{g}\sum_{i=1}^{g} MD^{normal,i}, \tag{10}$$

$$MD^{normal,i} = \sqrt{\frac{(EDM_{lower}^{normal,i} - \overline{EDM_{lower}^{normal}})^{T}(EDM_{lower}^{normal,i} - \overline{EDM_{lower}^{normal}})}{Cov}}. \tag{11}$$

The expectation of the lower triangles of normal EDMs ($\overline{EDM_{lower}^{normal}}$) and the covariance matrix ($Cov$) are given in (12) and (13).

$$\overline{EDM_{lower}^{normal}} = \frac{1}{g}\sum_{i=1}^{g} EDM_{lower}^{normal,i} \tag{12}$$

$$Cov = \begin{bmatrix} \sigma_{ED_{2,1}^{normal}ED_{2,1}^{normal}} & \sigma_{ED_{2,1}^{normal}ED_{3,1}^{normal}} & \cdots & \sigma_{ED_{2,1}^{normal}ED_{m,m-1}^{normal}} \\ \sigma_{ED_{3,1}^{normal}ED_{2,1}^{normal}} & \sigma_{ED_{3,1}^{normal}ED_{3,1}^{normal}} & \cdots & \sigma_{ED_{3,1}^{normal}ED_{m,m-1}^{normal}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{ED_{m,m-1}^{normal}ED_{2,1}^{normal}} & \sigma_{ED_{m,m-1}^{normal}ED_{3,1}^{normal}} & \cdots & \sigma_{ED_{m,m-1}^{normal}ED_{m,m-1}^{normal}} \end{bmatrix} \tag{13}$$

The covariance between two arbitrary second-order statistics in the lower triangle of normal EDM is as given in (14)

$$\sigma_{ED_{j,k}^{normal}ED_{l,v}^{normal}} = \frac{1}{g-1}\sum_{i=1}^{g}(ED_{j,k}^{normal,i} - \mu_{ED_{j,k}^{normal}})(ED_{l,v}^{normal,i} - \mu_{ED_{l,v}^{normal}}) \ , \tag{14}$$

where $\mu_{ED_{j,k}^{normal}} = \frac{1}{g}\sum_{i=1}^{g} ED_{j,k}^{normal,i}$. The standard deviation $\sigma$ can be obtained by using (15).

$$\sigma = \sqrt{\frac{\sum_{i=1}^{g}(MD^{normal,i} - \overline{MD^{normal}})^2}{g-1}} \tag{15}$$

The obtained distribution of the normal training traffic records is stored in the norm profile for attack detection.

## 4.2   Attack Detection

According to the definition of normal distribution, roughly 99.7% of the values are within 3 standard deviations from the mean. Therefore, the decision can be made by comparing the distance of an observed object to the mean of the distribution. If the distance is greater than 3 standard deviations from the mean, it is flagged as an attack with 99.7% confidence.

To make the comparison, the EDM of the observed traffic record ($EDM_{lower}^{observed}$) needs to be generated using the proposed multivariate correlation analysis approach. Then, the MD between the $EDM_{lower}^{observed}$ and the expectation ($\overline{EDM_{lower}^{normal}}$) of the lower triangles of normal EDMs stored in the normal profile is computed using (16).

$$MD^{observed} = \sqrt{\frac{(EDM_{lower}^{observed} - \overline{EDM_{lower}^{normal}})^T (EDM_{lower}^{observed} - \overline{EDM_{lower}^{normal}})}{Cov}} \tag{16}$$

Then, the values are compared with the pre-defined threshold given in (17).

$$Threshold = \mu + \sigma * n \tag{17}$$

For a normal distribution, $n$ is usually ranged from 1 to 3. This means that we would like to make a detection decision with a certain level of confidence varying from 68% to 99.7% in associate with the selection of different values of $n$. Therefore, if the observed MD is greater than the threshold, it will be considered as an attack.

## 5    Evaluation and Analysis

This section describes and analyzes the results obtained from the evaluation of the performance of the proposed multivariate correlation analysis based approach on DoS attack detection. The evaluation is conducted on KDD CUP 99 dataset [20].

The 10 percent labelled dataset of KDD CUP 99 dataset is applied for evaluation purpose. Although the dataset is not without criticism [21], it is the only public dataset with labelled attack samples. Moreover, for the comparison, many research works were evaluated using this dataset. There are six different types of DoS attacks available from the 10 percent labelled dataset. They are Teardrop, Smurf, Pod, Neptune, Land and Back attacks.

To evaluate an IDS, detection rate and false positive rate are two important metrics. We aim to achieve a high detection rate while retaining a low false positive rate. To visually reveal the performance of the IDS, Receiver Operating Characteristic (ROC) curve is employed to show the relations between these two metrics.

### 5.1    Preprocessing

We first filter all records with the labels of Normal, Teardrop, Smurf, Pod, Neptune, Land and Back from the 10 percent labelled dataset. Then, further classify them into different clusters according to their labels. The description of the filtered data is presented in Table 1.

**Table 1.** Number of records of normal and DoS attack records

| Normal | Teardrop | Smurf | Pod | Neptune | Land | Back |
|--------|----------|-----------|-------|-----------|------|--------|
| 97,260 | 9,790 | 2,807,900 | 2,640 | 1,072,010 | 210 | 22,030 |

### 5.2    Results and Analysis

To evaluate the detection performance of the proposed approach, we conduct 10 fold cross-validations and 32 continuous features are used. Norm profiles are built with respect to different types of traffic, namely TCP, UDP and ICMP traffic. In the training phase, we only employ the Normal records, while Normal records and the attack records are all involved in the test phase. In the test phase, we vary the parameter $n$ given in the (17) from 1 to 3 with an increment of 0.5, in order to compare the detection accuracy with the change of the threshold. The results are shown in Table 2.

As can be seen from Table 2, our proposed multivariate correlation analysis based detection approach performs very well in most of the cases. The detection rates of Normal records rise from 97.92% to 99.13% along with the increase of the threshold. The Smurf and Pod attack records are completely detected without being affected by the change of the threshold. For Teardrop and Neptune attacks, our approach achieves approximately 100% almost in all cases. The detection approach suffers from mirror degeneration in the case of Land attack when the threshold set greater than 2σ, but it still manages to detect around 74.76% of the attack records. However, when detecting Back attacks, the detection accuracy suffers from a sharp decrease from 92.37% to 45.71% and finally down to 6.42% when the threshold keeps increasing. This problem may be caused by the non-normalized data in which some features dominate the detection performance during their comparatively large values.

**Table 2.** Detection rates for normal and DoS attack records against different thresholds

| Type of records | Threshold | | | | |
|---|---|---|---|---|---|
| | 1σ | 1.5 σ | 2 σ | 2.5 σ | 3 σ |
| Normal | 97.92% | 98.47% | 98.75% | 98.99% | 99.13% |
| Teardrop | 100% | 100% | 100% | 99.99% | 99.98% |
| Smurf | 100% | 100% | 100% | 100% | 100% |
| Pod | 100% | 100% | 100% | 100% | 100% |
| Neptune | 100% | 100% | 100% | 99.99% | 99.99% |
| Land | 100% | 100% | 96.19% | 87.62% | 74.76% |
| Back | 92.37% | 45.71% | 9.93% | 9.16% | 6.42% |

To better understand the performance of our proposed multivariate correlation analysis based detection approach, ROC curves are given in Fig. 1.
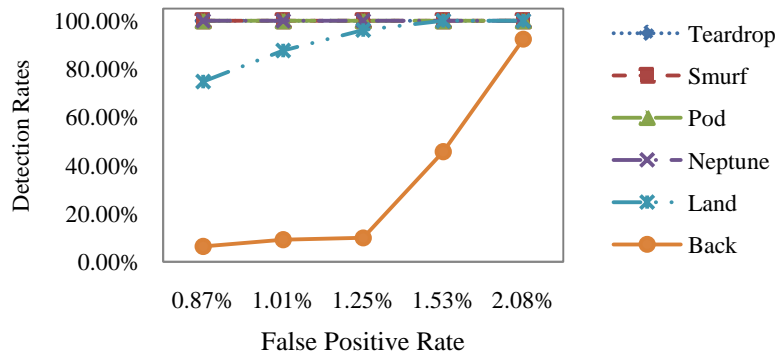


**Fig. 1.** ROC curves for the detection of DoS attacks

The ROC curves show clear tradeoff between the detection rate and false positive rate. There is a common trend that detection rates of all attacks increase when larger numbers of false positive are tolerated. If high detection rates are required, we have to endure a comparatively high false positive rate of 2.08%. However, 2.08% false

positive rate is still in the acceptable range and even better than other detection approaches.

To further evaluate the performance of the proposed approach, a comparison with covariance feature space based network intrusion detection [15] and HCVM [22] is given in Table 3.

**Table 3.** Performance comparison of different detection approaches

|  | Multivariate correlation analysis based detection approach (Threshold = 1δ) | Covariance feature space based network intrusion detection (Threshold approach with 3D principle) [15] | HCVM[22] |
|---|---|---|---|
| Detection rate | 99.96% | 99.95% | 93.2% |
| False positive rate | 2.08% | 10.33% | 5.4% |

In general, our multivariate correlation analysis has been proven by the evaluation results and it can correctly extract the statistical properties to exhibit the behaviour of network traffic. The application of multivariate correlation analysis in DoS detection gives promising outcomes.

## 6    Conclusions and Future Work

This paper has proposed a Euclidean distance based multivariate correlation analysis approach to extract the inner correlations (second-order statistics) of network traffic records, which can better exhibit the network traffic behaviours. We have evaluated the analysis approach using the KDD CUP 99 dataset. The results show that these second-order statistics can clearly reveal the changes of network behavior caused by DoS attack. The multivariate correlation analysis based DoS attack detection approach achieves 99.96% detection rate and 2.08% false positive rate. The detection accuracy is improved by involving the second-order statistics instead of the original first-order statistics into the classification.

However, our approach still suffers from a high false negative rate in detecting Back attack. This may be caused by the non-normalized data or the redundant features in the dataset. Therefore, we will employ data normalization methods and optimal feature selection in our future work in order to improve the detection accuracy. Also, temporal information will be considered in the successive research.

## References

1. Kay, J.: Low Volume Viruses: New Tools for Criminals. Network Security 2005 (2005) 16-18
2. Gordon, L.A., Loeb, M.P., Lucyshyn, W., Richardson, R.: 2005 CSI/FBI Computer Crime and Security Survey. Computer Security Journal 21 (2005) 1

3. Roesch, M.: Snort-lightweight Intrusion Detection for Networks. The 13th USENIX Conference on System Administration. USENIX, Seattle, Washington (1999) 229–238
4. Patcha, A., Park, J.M.: An Overview of Anomaly Detection Techniques: Existing Solutions and Latest Technological Trends. Computer Networks 51 (2007) 3448-3470
5. Denning, D.E.: An Intrusion-detection Model. IEEE Transactions on Software Engineering (1987) 222-232
6. 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. ACM (2003) 30-41
7. Lee, K., Kim, J., Kwon, K.H., Han, Y., Kim, S.: DDoS Attack Detection Method Using Cluster Analysis. Expert Systems with Applications 34 (2008) 1659-1665
8. 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 (2006) 459-468
9. 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 (2010) 6225-6232
10. 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 ACM (2010) 1193-1197
11. Fugate, M., Gattiker, J.R.: Computer Intrusion Detection with Classification and Anomaly Detection Using SVMs. International Journal of Pattern Recognition and Artificial Intelligence 17 (2003) 441-458
12. Lane, T., Brodley, C.E.: Temporal Sequence Learning and Data Reduction for Anomaly Detection. ACM Transactions on Information and System Security (TISSEC) 2 (1999) 295-331
13. Ye, N., Emran, S.M., Chen, Q., Vilbert, S.: Multivariate Statistical Analysis of Audit Trails for Host-based Intrusion Detection. IEEE Transactions on Computers (2002) 810-820
14. Manikopoulos, C., Papavassiliou, S.: Network Intrusion and Fault Detection: A Statistical Anomaly Approach. Communications Magazine, IEEE 40 (2002) 76-82
15. Jin, S., Yeung, D.S., Wang, X.: Network Intrusion Detection in Covariance Feature Space. Pattern Recognition 40 (2007) 2185-2197
16. Sarasamma, S.T., Zhu, Q.A., Huff, J.: Hierarchical Kohonen Net for Anomaly Detection in Network Security. Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on 35 (2005) 302-312
17. Thottan, M., Ji, C.: Anomaly Detection in IP Networks. Signal Processing, IEEE Transactions on 51 (2003) 2191-2204
18. Wang, H., Zhang, D., Shin, K.G.: Change-point Monitoring for the Detection of DoS Attacks. IEEE Transactions on Dependable and Secure Computing (2004) 193-208
19. 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. IEEE (2008) 75-81
20. Cheng, J., Hatzis, C., Hayashi, H., Krogel, M.A., Morishita, S., Page, D., Sese, J.: KDD Cup 2001 report. ACM SIGKDD Explorations Newsletter 3 (2002) 47-64
21. Tavallaee, M., Bagheri, E., Lu, W., Ghorbani, A.A.: A Detailed Analysis of the KDD Cup 99 Data Set. The Second IEEE International Conference on Computational Intelligence for Security and Defense Applications (2009)
22. Chen, Y., Pang, S., Kasabov, N., Ban, T., Kadobayashi, Y.: Hierarchical Core Vector Machines for Network Intrusion Detection. The 16th International Conference on Neural Information. Springer (2009) 520-529