

Cyber Attack Classification Based on Parallel Support Vector Machine

Mital Patel

MTech Scholar,
SIRT, Bhopal,

Yogdhar Pandey

Assistant Professor, CSE Dept,
SIRT, Bhopal

ABSTRACT

Cyber attack is becoming a critical issue of organizational information systems. A number of cyber attack detection methods have been introduced with different levels of success that is used as a countermeasure to preserve data integrity and system availability from attacks. The classification of attacks against computer network is becoming a harder problem to solve in the field of network security. This paper describes a Subset Selection Decision Fusion method to choose features (attributes) of KDDCUP 1999 intrusion detection dataset. Selection algorithm for distributed cyber attack detection and classification is proposed. Different types of attacks together with the normal condition of the network are modeled as different classes of the network data. We proposed Parallel Support Vector Machine (PSVM) algorithm for detection and classification of cyber attack dataset. Support Vector Machines (SVM) are the classifiers which were originally designed for binary classification. The classification applications can solve multi-class problems. Result shows that PSVM gives more detection accuracy for classes and comparable to false alarm rate.

Keywords:

Distributed cyber attack detection and classification, Subset selection decision fusion, Parallel Support Vector Machine, KDDCUP’99 and Confusion Matrix.

1. INTRODUCTION

The rapid increase in connectivity and accessibility of computer system has resulted frequent chances for cyber attacks. Attack on the computer infrastructures are becoming an increasingly Serious problem. Basically the cyber attack detection is a classification problem, in which we classify the normal pattern from the abnormal pattern (attack) of the system. Subset selection decision fusion method plays a key role in cyber attack detection. It has been shown that redundant and/or irrelevant features may severely affect the

accuracy of learning algorithms. The SDF is very powerful and popular data mining algorithm for decision-making and classification problems. It has been using in many real life applications like medical diagnosis, radar signal classification, weather prediction, credit approval, and fraud detection etc.

In this paper we proposed Parallel Support Vector Machine (PSVM) algorithm for detection and classification of cyber attack dataset. As we know that the performance of support vector machine is greatly depend on the kernel function used by SVM. Therefore, we modified the Gaussian kernel function in data dependent way in order to improve the efficiency of the classifiers. The relative results of the both the classifiers are also obtained to ascertain the theoretical aspects. The analysis is also taken up to show that PSVM performs better than SDF. The classification accuracy of PSVM remarkably improve (accuracy for Normal class as well as DOS class is

almost 100%) and comparable to false alarm rate and training, testing times.

The remainder of the paper is organized as follows. In Section II, we present KDDCUP’99 dataset. The Preliminary work of distributed cyber attack detection and classification is formulated in Section III. In section IV PSVM is proposed. The proposed Parallel Support Vector Machine algorithm is evaluated using KDD1999 intrusion detection datasets. The performance is analyzed by comparing to the feature subset selection and parallel support vector algorithm. Conclusions are provided in Section V.

2. KDD CUP ‘99 DATA SET DESCRIPTION

To check performance of the proposed algorithm for distributed cyber attack detection and classification, we can evaluate it practically using KDD’99 intrusion detection datasets [1]. In KDD99 dataset these four attack classes (DoS, U2R,R2L, and probe) are divided into 22 different attack classes that tabulated in Table I. The 1999 KDD datasets are divided into two parts: the training dataset and the testing dataset. The testing dataset contains not only known attacks from the training data but also unknown attacks. Since 1999, KDD’99 has been the most widely used data set for the evaluation of anomaly detection methods.

This data set is prepared by Stolfo et al. [11] and is built based on the data captured in DARPA’98 IDS evaluation program [12]. DARPA’98 is about 4 gigabytes of compressed raw (binary) tcpdump data of 7 weeks of network traffic, which can be processed into about 5 million connection records, each with about 100 bytes. For each TCP/IP connection, 41 various quantitative (continuous data type) and qualitative (discrete data type) features were extracted among the 41 features, 34 features (numeric) and 7 features (symbolic).

TABLE I.

DIFFERENT TYPES OF ATTACKS IN KDD99 DATASET

4 Main Attack Classes	22 Attack Classes
Denial of Service (DoS)	back, land, neptune, pod, smurt, teardrop
Remote to User (R2L)	ftp_write, guess_passwd, imap, multihop, phf,spy, warezclient, warezmaster
User to Root (U2R)	buffer_overflow, perl, loadmodule, rootkit

Probing(Information Gathering)	ipsweep, nmap, portsweep, satan
--------------------------------	---------------------------------

To analysis the different results, there are standard metrics that have been developed for evaluating network intrusion detections. Detection Rate (DR) and false alarm rate are the two most famous metrics that have already been used. DR is computed as the ratio between the number of correctly detected attacks and the total number of attacks, while false alarm (false positive) rate is computed as the ratio between the number of normal connections that is incorrectly misclassified as attacks and the total number of normal connections [87]. In the KDD Cup 99, the criteria used for evaluation of the participant entries is the Cost Per Test (CPT) computed using the confusion matrix and a given cost matrix [21]. A Confusion Matrix (CM) is a square matrix in which each column corresponds to the predicted class, while rows correspond to the actual classes. An entry at row *i* and column *j*, CM (*i*, *j*), represents the number of misclassified instances that originally belong to class *i*, although incorrectly identified as a member of class *j*. The entries of the primary diagonal, CM (*i*, *i*), stand for

the number of properly detected instances. Cost matrix is similarly defined, as well, and entry *C* (*i*, *j*) represents the cost penalty for misclassifying an instance belonging to class *i* into class *j*. Cost matrix values [21] employed for the KDD Cup 99 classifier learning contest are shown in Table 2. A Cost Per Test (CPT) is calculated by using the following formula:

$$PT = 1/N \sum_{i=1}^m \sum_{j=1}^m CM(i, j) * C(i, j)$$

Where CM and C is confusion matrix and cost matrix, respectively, and N represents the total number of test instances, m is the number of the classes in classification. The accuracy is based on the Percentage of Successful Prediction (PSP) on the test data set.

$$PSP = \frac{\text{number of successful instance classification}}{\text{number of instance in the test set}}$$

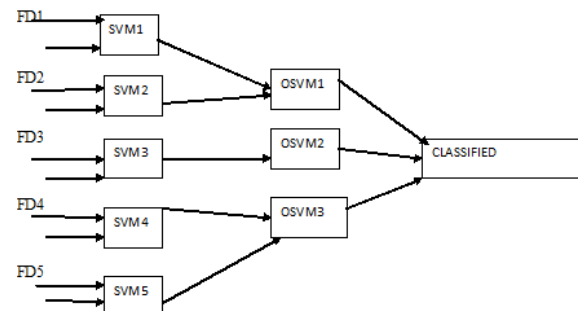
3. RELATED WORK

Support Vector Machine is a powerful tool to classify cyber attacks. But still it has some drawback. The first drawback is that SVM is very sensitive for attacks [17].The second, SVM designed for the two class problems it has to be extended for multiclass problem by choosing suitable kernel function. The performance of the SVM depends upon the kernel function. Some methods to improve the performance of SVM were proposed. Fuzzy SVM [13] is one of the improvements made on the traditional SVM. Several machine learning paradigms including Artificial Neural Network [14], Linear Genetic Programming (LGP) [15], Data Mining [16], etc. have been investigated for the classification of cyber attack. Also the machine learning techniques are sensitive to the noise in the training samples. The presence of mislabeled data if any can result in highly nonlinear decision surface and over fitting of the training set. This leads to poor generalization ability and classification accuracy. Decision-tree-based support vector machine which combines support vector machines and decision tree can be an effective way for solving multi-class problems. This method can decrease the training and testing time, increasing the efficiency of the system [2]. Improved Support Vector Machine (iSVM) algorithm for classification of cyber attack dataset which gives 100% detection accuracy

for Normal and Denial of Service (DOS) classes and comparable to false alarm rate, training, and testing times [8]. A new feature selection algorithm for distributed cyber attack detection and classification is proposed. Different types of attacks together with the normal condition of the network are modeled as different classes of the network data. Binary classifiers are used at local sensors to distinguish each class from the rest [9].

4. PROPOSED WORK

We proposed a new method for cyber attack classification based on parallel support vector machine based on distant feature set of attack attribute. All of the features are ranked based on their KullbackLeibler (K-L) distances, which is an alternative way to measure the importance of a feature in discriminating two classes. The features discriminating based on the euclidean distance formula for finding a similarity of features based on attack category. After calculation of discriminate we apply parallel support vector machine. SVM which was developed by Vapnik is one of the methods that is receiving increasing attention with remarkable results. SVM implements the principle of Structural Risk Minimization by constructing an optimal separating hyper plane in the hidden feature space, using quadratic programming to find a unique solution. Originally SVM was developed for pattern recognition problems. Recently, a regression version of SVM has emerged as an alternative and powerful technique to solve regression problems by introducing an alternative loss function. Although SVM has been successfully applied in many fields, there is a conspicuous problem appeared in the practical application of SVM. In parallel SVM machine first we reduced non-classified features data by distance matrix of binary pattern. From this concept, the cascade structure is developed by initializing the problem with a number of independent smaller optimizations and the partial results are combined in later stages in a hierarchical way, as shown in figure 1, supposing the training data subsets and are independent among each other.



(fig. 1 Cascaded SVM)

This figure shows that cascaded support vector machine, in this machine we passed five stage of features discernment and all these passes to optimized support vector machine for the processing of classification.

4.1 Step for data preprocessing.

- Transform data to the format of an SVM
- Conduct scaling on the data
- Consider the RBF kernel $K(x; y)$
- Use cross-validation to 2nd the best parameter C and

- Use the best parameter C and to train the whole training set
- Generate formatted data.

4.2 Step of cyber data classification.

- Read preprocessing data
 - For all the classes are represented
- ```
BEGIN
Find class with no attribute
Find class at Max cross product rate
Find the class at half cross product
REPEAT
Pointer= False
Find the intervals of hyper plane
If the end condition is met
Pointer = True
If the first interval has better results we should
Use this, otherwise the other
Find the class evaluation after cross product class
Instances middle times
UNTIL pointer= False
END
```
- Multiply all the classes with the best factor obtained;
  - Data classified.

## 5. CONCLUSION

In this paper we presented cyber attack detection and classification system to classify cyber attacks. A parallel support vector machine method for distributed cyber attack detection and classification is proposed. The new PSVM is shown more efficient for detection and classification of different types of cyber attacks compared to SDF. In both the method we use KDD'99 dataset to evaluate and compare the performance.

## 6. REFERENCE

- [1] KDD Cup 1999. Available on: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, October 2007.
- [2] Snehal A. Mulay, P.R. Devale, G.v. Garje, "Intrusion Detection System using Support Vector Machine and Decision Tree", International Journal of Computer Applications (0975 - 8887) Volume 3 - No.3, June 2010
- [3] Latifur Khan, Mamoun Awad, Bhavani Thuraisingham, "A new intrusion detection system using support vector machines and hierarchical clustering", The VLDB Journal DOI 10.1007/s00778-006-0002. 2007.
- [4] YMahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A. Ghorbani. "A detailed analysis of KDD CUP'99 data set", IEEE-2009.
- [5] <http://kdd.ics.uci.eduidatabases/kddcup99/kddcup99.html>
- [6] G.MeeraGandhi, Kumaravel Appavoo, S.K. Srivatsa, "Effective Network Intrusion Detection using Classifiers Decision Trees and Decision rules" Int. J. Advanced Networking and Applications Volume: 02, Issue: 03, Pages: 686-692 ,2010
- [7] P. Srinivasulu, R. Satya Prasad and I. Ramesh Babu, "Intelligent Network Intrusion Detection Using DT and BN Classification Techniques" Int. J. Advance. Soft Comput. Appl., Vol. 2, No. 1, March 2010 ISSN 2074-8523; Copyright © ICSRS Publication, 2010, www.i-csrs.org
- [8] Shailendra Singh Member, IEEE, IAENG, Sanjay Agrawal, Murtaza, A. Rizvi and Ramjeevan Singh Thakur, "Improved Support Vector Machine for Cyber Attack Detection", Proceedings of The World Congress on Engineering and Computer Science 2011 Vol I WCECS 2011, October 19-21, 2011, San Francisco, USA
- [9] Hoa Dinh Nguyen, Qi Cheng, "An Efficient Feature Selection Method For Distributed Cyber Attack Detection and Classification", 978-1-4244-9848-2/11 \$26.00©2011 IEEE
- [10] Dr. Adnan Mohsin Abdulazeez Brifcani & Adel Sabry Issa, "Intrusion Detection and Attack No.2, 2011
- [11] S. J. Stolfo, W. Fan, W. Lee, A. Prodromidis, and P. K. Chan, "Costbased modeling for fraud and intrusion detection: Results from the jam project," discex, vol. 02, p. 1130, 2000.
- [12] R. P. Lippmann, D. J. Fried, I. Graf, J. W. Haines, K. R. Kendall, D. McClung, D. Weber, S. E. Webster, D. Wyszogrod, R. K. Cunningham, and M. A. Zissman "Evaluating intrusion Detection systems: The 1998 darpa off-line intrusion detection evaluation," discex, vol. 02, p. 1012, 2000.
- [13] Xiong, Sheng-Wu, Liu Hong-bing, Niu Xiao-xiao, Fuzzy support vector machines based on FCM clustering. Proceedings of the fourth international conference on Machine Learning and Cybernetics, Guangzhou, China, Aug 18-21, IEEE, p.2608-2613, 2005.
- [14] A. K. Ghosh and A. Schwartzbard. "A study in Using Neural Networks for Anomaly and Misuse detection" Proceeding of the 8th USENIX Security Symposium, pp. 23-36. Washington, D.C. US. 1999
- [15] Mukkamala S., Sung AH, Abraham A. Modeling Intrusion Detection Systems Using linear genetic programming approach, The 17th international conference on industrial & engineering applications of artificial intelligence and expert systems, innovation in applied artificial intelligence.
- [16] W.Lee, S.J.Stolfo and K. Mok. Data mining in work flow environments: Experience in intrusion detection, Proceedings of the Conference on Knowledge Discovery and Data Mining (KDD-99), 1999.
- [17] Liu Yi-hung, Chen Yen-ting, face recognition using total margin based adaptive fuzzy support vector machines. IEEE Transactions on Neural Networks, 18(1): 178-192, 2007.