

Intrusion Detection System Using Data Mining Technique: Support Vector Machine

Yogita B. Bhavsar¹, Kalyani C. Waghmare²

¹Post Graduate Student, ²Assistant Professor, Pune Institute of Computer Technology, Pune, Maharashtra, India

Abstract— Security and privacy of a system is compromised, when an intrusion happens. Intrusion Detection System (IDS) plays vital role in network security as it detects various types of attacks in network. So here, we are going to propose Intrusion Detection System using data mining technique: SVM (Support Vector Machine). Here, Classification will be done by using SVM and verification regarding the effectiveness of the proposed system will be done by conducting some experiments using NSL-KDD Cup'99 dataset which is improved version of KDD Cup'99 data set. The SVM is one of the most prominent classification algorithms in the data mining area, but its drawback is its extensive training time. In this proposed system, we have carried out some experiments using NSL-KDD Cup'99 data set. The experimental results show that we can reduce extensive time required to build SVM model by performing proper data set pre-processing. Also when we do proper selection of SVM kernel function such as Gaussian Radial Basis Function, attack detection rate of SVM is increased and False Positive Rate (FPR) is decrease.

Keywords— Classification, Intrusion Detection System (IDS), Kernel Function, NSL- KDD, Pre-processing, Support Vector Machine (SVM)

I. INTRODUCTION

As network-based computer systems have important roles in modern society, they have become the targets of intruders. Therefore, we need to find the best possible ways to protect our systems. The security of a computer system is compromised when an intrusion takes place. An intrusion can be defined as any action done to hamper the integrity, confidentiality or availability of the system. There are some intrusion prevention techniques which can be used to protect computer systems as a first line of defense. But only intrusion prevention is not enough. As systems become more complex, there are always exploitable weaknesses in the systems due to design and programming errors, or various penetration techniques. Therefore Intrusion detection is required as another measure to protect our computer systems from such type of vulnerabilities.

II. RELATED WORK

In 1980, the concept of intrusion detection began with Anderson's seminal paper [1]; he introduced a threat classification model that develops a security monitoring surveillance system based on detecting anomalies in user behavior.

In 2003, Kaining Lu Zehua Chen Zhigang Jin Jichang Guo, [4] has presented one collaborate IDS module to make a real-time detection and block intrusions before occurrences based on HIDS using sequences of system call anomaly detection. In 2009, Chunhua Gu and Xueqin Zhang,[6] proposed a system using rough set for attribution reduction and support vector machine for intrusion detection classification. In 2009, Yong-Xiang Xia Zhi-Cai Shi and Zhi-Hua Hu, [5] proposed a method of detecting intrusion using incremental SVM based on key feature selection. Again in the same year, Rung-Ching Chen, Kai-Fan Cheng and Chia-Fen Hsieh, [7] used RST (Rough Set Theory) and SVM (Support Vector Machine) to detect intrusions. First, RST is used to preprocess the data and reduce the dimensions. Next, the features were selected by RST will be sent to SVM model to learn and test respectively.

In 2010, Heba F. Eid [8] effectively introduced intrusion detection system by using Principal Component Analysis (PCA) with Support Vector Machines (SVMs) as an approach to select the optimum feature subset. In 2011, Shingo Mabu, Member, *IEEE*, [9] has described a novel fuzzy class-association rule mining method based on genetic network programming (GNP) for detecting network intrusions. Again in the same year, Carol J Fung and Jie Zhang, [10] have proposed Dirichlet-based trust management to measure the level of trust among IDSes according to their mutual experience.

Recently in 2012, [11] has described an adaptive network intrusion detection system which uses a two stage architecture. In the first stage a probabilistic classifier is used to detect potential anomalies in the traffic. In the second stage a HMM based traffic model is used to narrow down the potential attack IP addresses. Again in 2012, V. Jaiganesh, [15] proposed Kernelized Support Vector Machine with Levenberg-Marquardt (LM) Learning. Again In 2012, Gholam Reza Zargar, Tania Baghaie, [13] proposed a category-based selection of effective parameters for intrusion detection using Principal Components Analysis (PCA).

III. DATA SET COLLECTION AND PRE-PROCESSING

A. Data Set Collection

To verify the effectiveness and the feasibility of the proposed IDS system, we have used NSL-KDD dataset.

Here, we have determined two target classes: class ‘zero’ for normal instance and class ‘one’ for attack or intrusion. Then we have to save target class and feature values of each instance in libSVM format. LibSVM format is:

[Label][index1]:[value1][index2]:[value 2].....

Where,

‘Label’ is target classes of classification. Usually we put integers for the class value. We can have [0, 1] for target class or [-1, +1] for target class.

Here, we have used [0, 1] target class. Where class ‘0’ indicates ‘normal’ and class ‘1’ indicates ‘attack’.

‘Index’ is the ordered index. Usually continuous integer.

‘Value’ is the input data for training. Usually lots of real numbers.

Input dataset to the problem we are trying to solve involves 41 of ‘features’, so the input will be a set of these 41 features.

After this conversion, we have to perform linear scaling of libSVM format datasets and store these scaled datasets for further use. Linear scaling of datasets is done to improve the performance of classification using SVM.

VI. SUPPORT VECTOR MACHINE

The Support Vector Machine is one of the most successful classification algorithms in the data mining area. SVM uses a high dimension space to find a hyper-plane to perform binary classification. SVM approach is a classification technique based on Statistical Learning Theory (SLT). It is based on the idea of hyper plane classifier. The goal of SVM is to find a linear optimal hyper plane so that the margin of separation between the two classes is maximized.

The SVM uses a portion of the data to train the system. It finds several support vectors that represent the training data. These support vectors will form a SVM model. According to this model, the SVM will classify a given unknown dataset into target classes.

VII. INTRUSION DETECTION USING SVM

In the proposed system, we have constructed a SVM model for classification. While intrusion behaviors happen, SVM will detect the intrusion. A classification task involves training set and testing set which consist of instances. Each instance in the training set contains one “target value” (class labels: Normal or Attack) and several “attributes” (features). The goal of SVM is to produce a model which predicts target value of data instance in the testing set which is given only attributes.

To achieve this goal, we have used kernel functions available with SVM. There are 3 major SVM kernel functions:

- (i) Gaussian Kernel (Radial Basis Function)
- (ii) Polynomial kernel
- (iii) Sigmoid kernel

(i) *Gaussian Kernel Function:* The Gaussian kernel is an example of radial basis function kernel.

$$K(X_i, X_j) = \exp \left\{ -\frac{\|X_i - X_j\|^2}{2\sigma^2} \right\}$$

Where, σ stands for window width.

(ii) *Polynomial Kernel Function:* This Polynomial kernel is a non-stationary kernel. Polynomial kernels are well suited for problems where all the training data is normalized. Adjustable parameters are the constant term c and the polynomial degree d .

$$K(X_p, X_j) = (X_p, X_j)^d + c$$

(iii) *Sigmoid Kernel Function:* Sigmoid Kernel is also called as Hyperbolic Tangent Kernel and the Multilayer Perceptron (MLP) kernel. The Sigmoid Kernel comes from the Neural Networks field, where the bipolar sigmoid function is often used as an activation function for artificial neurons

$$K(X_i, X_j) = \tanh(k(X_i, X_j) + r)$$

In classification phase, SVM training model is build and SVM kernel function is selected to generate classification results. The system design for IDS using SVM is shown in figure 3.

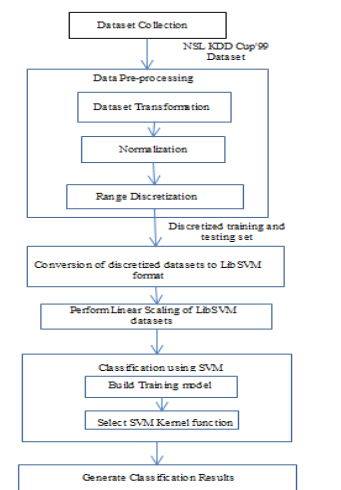


Figure 3: System Design for IDS using SVM

VIII. EXPERIMENTAL RESULTS

The proposed Intrusion Detection System is experimented using the Waikato Environment for Knowledge Analysis (WEKA 3.7) and LibSVM 1.5.

WEKA is a complete set of Java class libraries that execute several state-of-the-art machine learning and data mining approaches [20] and LibSVM is a library for support vector machines [19]. Dataset used for experiment purpose is NSL-KDD dataset which is a new version of KDDcup99 dataset and has some advantages over KDDcup99 dataset.

Time required to build model using different SVM kernel functions and 10 fold cross validation as well as 10 fold cross validation with re-evaluation using supplied test set is shown in table 4 and respective graph is shown in figure 4.

Table 5 shows, accuracy achieved using different SVM kernel functions and classification using 10 fold cross validation as well as 10 fold cross validation with re-evaluation using supplied test set. The graph for accuracy using different SVM kernel functions at the time of classification is shown in figure 5.

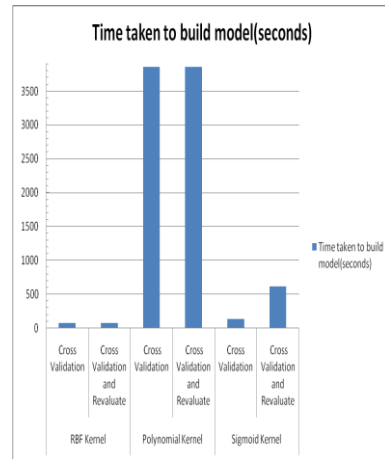


Figure 4: Model Building Time with different SVM kernel functions

**TABLE IV
TIME TAKEN TO BUILD MODEL USING SVM**

Kernel Type	Classification Type	Time taken to build model(seconds)
RBF Kernel	Cross Validation	77.01
	Cross Validation and Reevaluate	77.01
Polynomial Kernel	Cross Validation	3859.57
	Cross Validation and Reevaluate	3859.57
Sigmoid Kernel	Cross Validation	134.64
	Cross Validation and Reevaluate	615.75

**TABLE V
ACCURACY OF DIFFERENT SVM KERNEL FUNCTIONS**

Kernel Type	Classification Type	Accuracy
RBF Kernel	Cross Validation	94.1857 %
	Cross Validation and Reevaluate	98.5749 %
Polynomial Kernel	Cross Validation	98.4281 %
	Cross Validation and Reevaluate	98.4281 %
Sigmoid Kernel	Cross Validation	98.5749 %
	Cross Validation and Reevaluate	73.0886 %

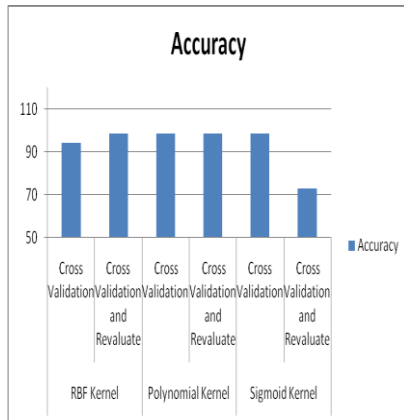


Figure 5: Accuracy of different SVM Kernel functions

IX. CONCLUSION

Now a days, intrusion which affects the security and privacy of the system, has become major concern for many organizations. Hence, there is a need of strong IDS which can detect novel attack with high attack detection accuracy. In this paper, we have proposed a method of intrusion detection using SVM which can reduce the time required to build model for classification and increase the intrusion detection accuracy when Gaussian RBF kernel is used.

The experimental results show that, when data sets are properly processed and proper SVM kernel is selected i.e. Radial Basis Function (RBF), it can overcome the drawback of SVM i.e. extensive time required to build model.

When we have conducted experiment with 10 fold cross validation and Gaussian RBF kernel of SVM, the time required to build model was **77.07 seconds** and attack detection accuracy achieved was **94.1857 %**. This attack detection accuracy was increased to **98.5749 %**, when we have changed classification to 10 fold cross validation and re-evaluation using supplied test set with same RBF SVM kernel function.

REFERENCES

- [1] James P. Anderson, "Computer Security Threat Monitoring and Surveillance," Technical report, James P. Anderson Co., Fort Washington, Pennsylvania. April 1980.
- [2] Tomas Abraham, "IDDM: INTRUSION Detection using Data Mining Techniques", Technical report DTSO electronics and surveillance research laboratory, Salisbury.
- [3] Wenke Lee and Salvatore J. Stolfo, "A Framework for constructing features and models for intrusion detection systems", ACM transactions on Information and system security (TISSEC), vol.3, Issue 4, Nov 2000.
- [4] Kaining Lu Zehua Chen Zhigang Jin Jichang Guo." An Adaptive Real-Time Intrusion Detection System Using Sequences of System Call", CCECE 2003.
- [5] Yong-Xiang Xia, Zhi-Cai Shi, Zhi-Hua Hu," An Incremental SVM for Intrusion Detection Based on Key Feature Selection" 2009 Third International Symposium on Intelligent Information Technology Application.

- [6] Chunhua Gu and Xueqin Zhang," A Rough Set and SVM Based Intrusion Detection Classifier", Second International Workshop on Computer Science and Engineering, 2009.
- [7] Rung-Ching Chen, Kai-Fan Cheng and Chia-Fen Hsieh ,"Using Rough Set And Support Vector Machine For Network Intrusion Detection" International Journal of Network Security & Its Applications (IJNSA), Vol 1, No 1, 2009.
- [8] Heba F. Eid, Ashraf Darwish, Aboul Ella Hassanien, and Ajith Abraham" Principle Components Analysis and Support Vector Machine" based Intrusion Detection System", IEEE 2010.
- [9] Shingo Mabu, Ci Chen, Nannan Lu, Kaoru Shimada, and Kotaro Hirasawa," An Intrusion-Detection Model Based on Fuzzy Class-Association-Rule Mining Using Genetic Network Programming", IEEE Transactions On Systems, Man, And Cybernetics—Part C: Applications And Reviews, Vol. 41, No. 1, January 2011.
- [10] Carol J Fung and Jie Zhang," Dirichlet-Based Trust Management for Effective Collaborative Intrusion Detection Networks" ,IEEE Transactions on Network And Service Management, Vol. 8, No 2, June 2011.
- [11] R Rangadurai Karthick, Vipul P. Hattiwale, Balaraman Ravindran," Adaptive Network Intrusion Detection System using a Hybrid Approach ", IEEE 2012.
- [12] Vincent F. Mancuso, Dev Minoira, Nicklaus Giacobbe, Michael McNeese and Michael Tyworth " *idsNETS*: An Experimental Platform to Study Situation Awareness for Intrusion Detection Analysts", IEEE International Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support, New Orleans, LA, 2012.
- [13] Gholam Reza Zargar, Tania Baghaie, "Category-Based Intrusion Detection Using PCA", Journal of Information Security, 2012.
- [14] Neethu B, "Classification of Intrusion Detection Dataset using machine learning Approaches", International Journal of Electronics and Computer Science Engineering, 2012.
- [15] V. Jaiganesh, "Intrusion Detection Using Kernelized Support Vector Machine With Levenbergmarquardt Learning", International Journal of Engineering Science and Technology, 2012.
- [16] Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali, A. Ghorbani, "A Detailed Analysis of the KDD CUP 99 Data Set", proceeding of IEEE symposium on computational Intelligence in security and defense application, 2009.
- [17] KDD Cup 1999. Available on <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, October 2007.
- [18] R.O. Duda, P.E. Hart, and D. G. Stork, "Pattern Classification", Vol. 1, Wiley, 2002.
- [19] Chih-Chung Chang and Chih-Jen Lin, LIBSVM: a library for support vector machines. ACM Transactions on Intelligent Systems and Technology, 2:27:1--27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [20] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2009); the WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11, Issue 1.
- [21] Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A. Ghorbani A Detailed Analysis of the KDD CUP 99 Data Set proceeding of the 2009 IEEE symposium on computational Intelligence in security and defense application.