

Intrusion Detection Model Based on Ensemble Learning for U2R and R2L Attacks

Ployphan Sornsuwit

Department of Computer Science, Faculty of Science
King Mongkut's Institute of Technology
Bangkok, Thailand
ployphan@kpru.ac.th

Saichon Jaiyen

Department of Computer Science, Faculty of Science
King Mongkut's Institute of Technology
Bangkok, Thailand
kjsaicho@kmitl.ac.th

Abstract— Intrusion Detection System (IDS) is a tool for anomaly detection in network that can help to protect network security. At present, intrusion detection systems have been developed to prevent attacks with accuracy. In this paper, we concentrate on ensemble learning for detecting network intrusion data, which are difficult to detect. In addition, correlation-based algorithm is used for reducing some redundant features. Adaboost algorithm is adopted to create the ensemble of weak learners in order to create the model that can protect the security and improve the performance of classifiers. The U2R and R2L attacks in KDD Cup'99 intrusion detection dataset are used to train and test the ensemble classifiers. The experimental results show that reducing features can improve efficiency in attack detection of classifiers in many weak learners.

Keywords— Intrusion detection, Ensemble, Adaboost, KDD Cup'99, Feature Reduction

I. INTRODUCTION

Currently, intrusions have various patterns of attack behaviors that are more difficult to detect than in the past. For example, some attack patterns require long periods for analyzing packets, or some attacks have a few amount of traffic. Therefore, the efficiency of traditional methods may be poor to detect intrusions accurately. IDS are divided into 2 types which are signature-based detection and anomaly-based detection. Signature-based detection will match the unknown pattern with known pattern and then consider whether it is normal or abnormal. However, anomaly detection is to identify the behaviors that are deviated from normal patterns[1]. Both two types have difference advantages and drawbacks. Signature-based detection can give high accuracy because it can match predefined attack behavior in database, but it cannot detect novel attack. On the other hands, anomaly-based detection has low accuracy and high false alarm because anomaly-based method uses statistical methods to analyze packets, and it can detect novel attacks.

Anomaly detections using machine learning method have been investigated in a number of researches. Ensemble methods are adopted to detect anomaly patterns [2] and show the better performance when they use multiple classifiers [3-

5]. Some researches include preprocessing method in order to reduce redundant features. Consequently, they have only relevant features and produce the higher performance. However, some researches may be poor efficient to detect difficult attack types in datasets such as U2R and R2L types because both attack types have a few number of instants and more complicated behaviors.

In this paper, we present an algorithm that can overcome these problems, increase the accuracy, and decrease the false alarm rate of U2R and R2L attacks by using Correlation-based feature selection and multiple weak classifiers such as Naïve Bayes, Decision Tree, MLP, k-NN and SVM based on Adaboost algorithm. The rest of paper is organized as follows. In section 2, we present the related works. In Section 3, we present the proposed method. In section 4, the experimental results are described. In last section, conclusions and future work are presented.

II. RELATE WORKS

There are many research topics in intrusion detection system with several algorithms that improve accuracy and decrease false alarm rate. One of the topics to classify anomaly is to use ensemble methods to improve the performance of classifiers. Zhaza Merghani AbdElrahaman and Ajith Abraham [8] presented the comparisons of performance in many algorithms using Boosting and Bagging techniques with KDD Cup'99 datasets that contained 4 types of DoS Probe U2R and R2L. They compared the performance of these methods with Adaboost algorithm and Bagging technique. Then, they presented a new hybrid ensemble algorithm for detecting intrusion based on Error Collecting Output Code (ECOC). From their experimental results, they found that the proposed methods improved accuracy and false alarm rate.

Te-Shun Chou¹, Jeffrey Fan, Sharon Fan, and Kia Makki² [5] presented three layers of multiple classifiers for intrusion detection that was able to improve the overall accuracy. They applied three different patterns of learning including Naive Bayes, fuzzy k-NN, and back-propagation neural network for generating the decision boundary. The experiment used KDD cup'99 dataset with 30 features to test the model. The

experiment showed the better performance of combination methods.

Wei-Hu and Weiming Hu [9] presented NIDS using Adaboost methods with decision stumps weak learner. The proposed framework consists of four modules for anomaly detection including feature extraction, data labeling, weak classifier design, and strong classifier construction. The experiments used KDD Cup'99 dataset with 41 features. The experimental results showed that the proposed framework not only improved the false positive rate and detection rate but also reduced the computational complexity compared to genetic clustering, hierarchy SOM, and SVM.

Debojit Boro, Bernard Nongpoh, and Dhruva K. Bhattacharyya [10] presented the combination of several models for improving the performance of intrusion detection systems. They used four layers to combine the appropriated classifiers by using several weak learners. The experimental results showed that their proposed method was able to improved classification rate and false positive rate comparing to a single classifier.

P.Giffy Jeya, M Ravichandran, and C.S. Ravichandran [1] presented a new method for anomaly classification by using correlation based analysis to reduce the number of features and applying Fisher Linear Discriminant Analysis (FLDA) to detect intrusion data. KDD Cup'99 dataset was used in the experiments and the experimental results showed that the proposed method improved the accuracy of classification of U2R and R2L attacks comparing to other research.

Neelam Sharma and Sarubh Mukherjee [4] proposed a layer approach for improving the efficiency of attack detection rate by using domain knowledge and the sequential search to decrease feature sets and used Naïve Bayes classifier for classifying 4 classes of attack types. The experimental results showed that the proposed method was able to improve the recall in every attack type comparing to non-layer with feature selection method.

A. Ensemble Learning

Ensemble Learning is a method that combines multiple learners to solve the specific problems in order to improve the accuracy of classifiers. The traditional learning approaches are employed to construct each learner from training data and combine them to produce the final result. Ensemble learning is also called committee-based learning, or learning multiple classifier systems [11]. Figure 1 illustrates the architecture of the ensemble model.

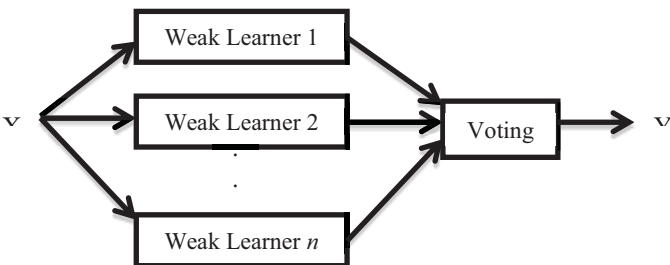


Figure1. The basic architecture of ensemble classifiers.

B. Adaboost Methods

Adaboost is an adaptive boosting algorithm for constructing a strong classifier as linear combination of weak classifiers. Each training sample uses a weight to determine the probability of being selected for a training set, and the combination is based on the weighted voting of weak classifiers [11]. The Adaboost algorithm is summarized as follows.

1. $D_x(x) = \frac{1}{m}$ (Initialize the distributions of data)
2. for $t=1, \dots, T$
3. $h_t = \mathcal{E}(D, D_t)$ (Train the base classifier h_t using D with distribution D_t)
4. $\epsilon_t = P_{x \sim D_t}(h_t(x) \neq f(x))$ (Evaluate the error of h_t)
5. if $\epsilon_t > 0.5$ then break.
6. $\alpha_t = \frac{1}{2} \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$; (Determine the weight of h_t)
7. $D_{t+1}(x) = D_t(x) \exp\left(\frac{-\alpha_t f(x) h_t(x)}{Z_t}\right)$ (Update distributions of data)

The output of the ensemble classifiers can be computed as:

$$\mathcal{H}(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

C. Naïve Bayes

Naïve Bayes is one of probability method based on Bayes' theorem [12] used to determine the appropriate class of unseen data, Naïve Bayes computes the posterior probability for each class C_i as follows:

$$P(C_i|X) = \frac{P(C_i) \prod_{k=1}^n P(X_k|C_i)}{P(X)}$$

C_i is the class label and X is an instance to be classified. After calculating the posterior probability for each class, it assigns the class label that has the highest probability to the unseen data.

D. Multilayer Perceptron (MLP)

MLP is one of artificial neural networks that perform linear mapping from input space to hidden space and from hidden space to output space. The network structure consists of an input layer, a hidden layer, and an output layer. The input layer consists of nodes for receiving the input. The output from previous layer is the input of the next layer. The layer between input layer and output layer is called hidden layer. The last layer is the output layer that is used to predict the class label of data. MLP neural network uses Back propagation learning algorithm to train the network which is described as follows [13]:

1. Initialize the weights and the learning rate η .
2. Present an input x_i into the network and calculate each output, y_k , by:

$$y_k = \sum_{j=1}^p w_{jk} x_j + \theta_k$$

where w_{jk} are weights and θ_k is a bias.

3. For each output node k , compute the error gradient, δ_k , by:

$$\delta_k = (y_{target} - y_k) y_k (1 - y_k)$$
4. For each hidden node j , compute the error gradient, δ_j , by:

$$\delta_j = o_j (1 - o_j) \sum_k w_{jk} \delta_k$$
5. Adjust all weights by:

$$w_{ij} = w_{ij} + \Delta w_{ij}$$
 where $\Delta w_{ij} = \eta \delta_j o_i$
6. Go to step 2.

F. Decision Tree

Decision tree is a classification model that uses a tree-like structure to perform a decision. Decision trees are commonly used in operation research and intrusion detection because it gives the better performance compared to other algorithms [14]. The leaning algorithm of decision tree is described as follows:

1. Select the attribute with the highest information gain
2. Let P_i be the probability of an arbitrary tuple in D belonging to class C_i estimated by $|C_i, D| / |D|$
3. Expected information (entropy) needed to classify a tuple can be calculated as:

$$Info(D) = - \sum_{i=1}^m P_i \log_2 P_i$$

4. Information needed to classify D can be computed as:

$$Info_A(D) = \sum_{i=1}^v \frac{D_i}{D} * Info(D_i)$$

5. Information gained by branching on attribute A can be computed as:

$$Gain_A = Info(D) - Info_A(D)$$

F. Support Vector Machine (SVM)

SVM is one of machine learning methods used to solve the classification problems based on an optimal hyperplane in a high-dimensional space. SVM is applied for classification, regression, and other tasks [15]. For a binary classification, all training data are classified based on the following constraints.

For all x that is a member of class +1, they satisfy the following constraints:

$$w^T x + b \geq +1$$

For all x that is a member of class -1, they satisfy the following constraints:

$$w^T x + b \leq -1$$

The purpose of SVM is to find the optimal Hyperplane defined by $w^T x + b = 0$ that maximizes the margin of the two conditions above. After finding the optimal hyperplane, the decision function is defined as:

$$f(x) = \text{sign}(w^T x + b)$$

G. k -Nearest Neighbor (k -NN)

k -NN method uses k records in the training set that are similar (neighboring) to a new record in order to classify this record into its class [16]. The main issue is how to measure the similarity between records. The most popular measurement of similarity is the Euclidean distance between two records (x_1, x_2, \dots, x_p) and (y_1, y_2, \dots, y_p) defined by the following equation.

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}$$

H. KDD Cup'99 dataset

KDD'99 dataset consists of approximately 4,900,000 records which contain 41 features [17]. Each record is labeled as either normal or attack. There are 4 attack types in this data set described as follows:

- Denial of Service Attack (DoS): an attack in which the attacker makes some computing or memory resource too busy or too full to handle legitimate requests, or denies legitimate users to access a machine, e.g. syn flood.
- User to Root Attack (U2R): a class of exploit in which the attacker starts out with access to a normal user account on the system (perhaps gained by sniffing passwords, a dictionary attack, or social engineering). It exploits some vulnerabilities to gain root access to the system, e.g. various "buffer overflow" attacks.
- Remote to Local Attack (R2L): occurrence when an attacker with hasan's ability sends packets to a machine over a network but does not have the account on that machine. The attacker exploits some vulnerability to gain the local access as a user of that machine, e.g. guessing password.
- Probing Attack: an attempt to gather information about a network of computers for circumventing its security controls, e.g. port scanning.

In this research, we focus on U2R and R2L types for training and testing. The training data consist of various subtypes and amounts as shown in Table 1 and the testing data are shown in Table 2. In testing dataset, some subtypes are not included in training data such as sendmail, snmpgetattack, snmpguess, worm, xlock, and xsnoop in U2R type and ps, sqlattack, and xterm in R2L type. It is hard to correct testing data that

include amount of new subtypes much more than subtype in training datasets.

Table 1. U2R and R2L in Training Data

Type	Amounts
R2L	
warezclient	1,020
guess_passwd	53
warezmaster	20
imap	12
ftp_write	8
multihop	7
phf	4
spy	2
U2R	
Buffer_overflow	30
Rootkit	10
Loadmodule	9
Perl	3
Total	1,178

Table 2. U2R and R2L Testing Data

Type	Amounts
R2L	
ftp_write	3
guess_password	4,367
imap	1
multihop	18
named	17
phf	2
sendmail	17
snmpgetattack	7,741
snmpguess	2,406
warezmaster	1,602
worm	2
xlock	9
xsnoop	4
httptunnel	158
U2R	
buffer_overflow	22
loadmodule	2
perl	2
rootkit	13
ps	16
sqlattack	2
xterm	13
Total	16,417

III. PROPOSE METHODS

In this paper, Adaboost methods are applied to protect network security. This is because Adaboost provides many advantages for classification of intrusions [18] such as:

- Implementation is simple, and it is applied to many problems of pattern recognitions.
- It is less sensitive to overfitting than other learning algorithms.
- Different feature types are difficult to find the relations, but the algorithm is implemented and improves performance of a weak classifier with a strong classifier.
- Simple and various weak classifiers are selectively used.

KDD Cup'99 Dataset from UCI Repository is used to train and test the ensemble classifiers and the results of ensemble classifiers are compared to single classifier and methods without reduce feature sets. Our implementation model consists of 4 stages showing in Figure 2.

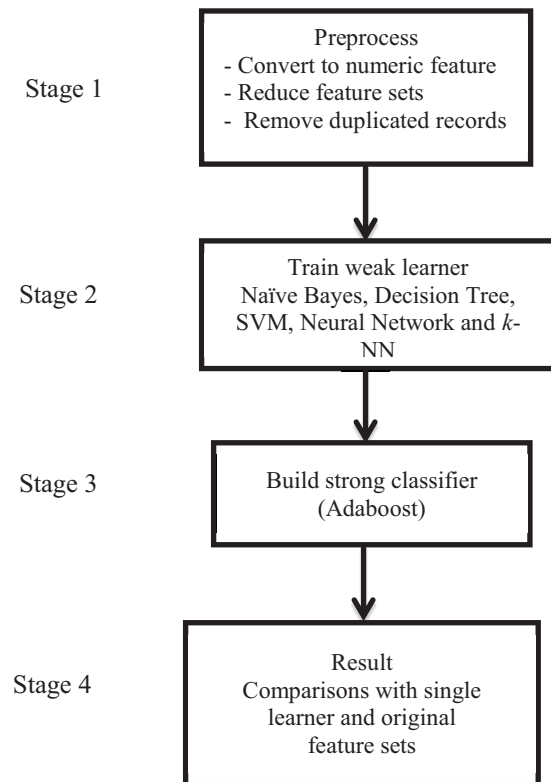


Figure 2. Implementation Model.

IV. EXPERIMENTAL RESULT

The experiments are conducted on WEKA mining tools, and Matlab 2014a. In stage 1, we convert every symbolic feature into numeric features in 3 fields including protocol type, service, and flag. Next, Correlation-based approach is

used for reducing redundant features. After reducing feature sets, 5 features are selected as follows:

1. `dst_bytes`: number of data bytes from destination to source
2. `num_compromised`: number of "compromised" conditions
3. `root_shell`: root shell is obtained; 0 otherwise if 1
4. `num_file_creations`: number of file creation operations
5. `is_host_login`: if the login belongs to the "hot" list; 0 1 otherwise

In preprocessing process, the duplicated data are removed in order to reduce the number of redundant data so the remaining data after removing are 1,051 records in training data and 3,128 records in testing data as shown in Table 3.

Table 3. Amount of training and testing data.

Attack Type	Number of original records		Number of duplicate records	
	Training	Testing	Training	Testing
R2L	1126	16347	999	3058
U2R	52	70	52	70
Total	1,178	16,417	1,051	3,128

In Stage 2, Naïve Bayes, Decision Tree, MLP, SVM, and k -NN are adopted as weak learners. Each learner is trained using preprocessed data from Stage 1. We set k as 4 for k -NN method and set the number of hidden neurons as 3 for MLP because these parameters provide the best result and apply Adaboost in Stage 3 by setting the number of iterations as 100.

In Stage 4, we show the comparisons of sensitivity and specificity in order to find the best classifier. In this experiment, all classifiers are trained and tested using the data set with all features and the data set with 5 selected features. Table 4 shows the performance of all classifiers tested on the dataset with all features. Table 5 shows the performance of all classifiers tested on the data set with 5 selected features.

From the experimental results in Table 4, the ensemble of MLPs can produce the highest sensitivity of 0.1753 and the ensemble of Naïve Bayes can achieve the highest specificity of 0.9856.

Table 4. The comparisons of sensitivity and specificity of single classifier and ensemble classifier with all features.

Weak Classifier	Single Classifier		Ensemble Classifier	
	Sensitivity	Specificity	Sensitivity	Specificity
Naïve Bayes	0.0443	0.9904	0.0370	0.9856
Decision Tree	0.0477	0.9850	0.0507	0.9850
MLP	0.1290	0.9880	0.1753	0.9855
SVM	0.0250	0.9791	0.0387	0.9834
k -NN	0.0326	0.9837	0.0344	0.9847

Table 5. The comparisons of sensitivity and specificity of single classifier and ensemble classifier with 5 selected features.

Weak Classifier	Single Classifier		Ensemble Classifier	
	Sensitivity	Specificity	Sensitivity	Specificity
Naïve Bayes	0.7272	0.9876	0.7600	0.9905
Decision Tree	0.0582	0.9870	0.0650	0.9890
MLP	0.7000	0.9841	0.7600	0.9896
SVM	0.7500	0.9861	0.7567	0.9864
k -NN	0.6727	0.9892	0.6851	0.9892

From the results in Table 5, the ensemble methods with feature selection can improve the sensitivity and specificity of all classifiers. The ensemble of Naïve Bayes, and MLP can achieve the highest sensitivity of 0.7600 and the ensemble of Naïve Bayes can provide the highest specificity of 0.9905. However, Decision Tree has the least sensitivity and specificity.

From experimental results, it can be concluded that the ensemble method with feature selection can improve the performance of all classifiers. Consequently, it can protect the security of U2R and R2L attacks. In addition, the ensemble classifiers based on Adaboost gives the higher performance than single classifiers.

V. CONCLUSION

In this paper, we adopt Adaboost algorithm to create the ensemble of Decision Tree, Naïve Bayes, SVM, and MLP classifiers for detecting U2R and R2L attacks which are difficult to detect. In addition, the correlation-based method is used to reduce redundant features. The performance of the ensemble classifiers is evaluated by using KDD CUP'99 datasets in UCI Repository datasets and the results are compared to single classifiers. The experiments are conducted on the data set with all features and some selected features. The experiments show that ensemble classifiers based on Adaboost can improve the performance of all classifiers. In addition, the ensemble of Naïve Bayes and MLP produce the highest sensitivity and the ensemble of Naïve Bayes has the highest specificity when they are tested on the data with some selected features. However, Decision Tree results the least performance because both attack types contain small data and have long period time pattern. U2R and R2L attacks are more similar pattern of normal users and the number of data are very small.

REFERENCES

- [1] P. G. Jeya, M. Ravichandran, and C. S. Ravichandran, "Efficient Classifier for R 2 L and U 2 R Attacks," *International Journal of Computer Applications*, vol. 45, no. 21, 2012.
- [2] P. Mrutyunjaya, and P. Manas Ranjan, "Ensemble of classifiers for detecting network intrusion," in *Proceedings of the International Conference on Advances in Computing, Communication and Control*, Mumbai, India, 2009.
- [3] P. Natesan, and P. Balasubramanie, "Multi stage filter using enhanced adaboost for network intrusion detection."
- [4] S. Neelam, and M. Saurabh, "Layered approach for intrusion detection using naïve Bayes classifier," in *Proceedings of the International Conference on Advances in Computing, Communications and Informatics*, Chennai, India, 2012.
- [5] C. Te-Shun, J. Fan, S. Fan, and M. Kia, "Ensemble of machine learning algorithms for intrusion detection." pp. 3976-3980.
- [6] Z. Safaa, E.-A. Mohammed, and K. Fakhri, "Features selection approaches for intrusion detection systems based on evolution algorithms," in *Proceedings of the 7th International Conference on Ubiquitous Information Management and Communication*, Kota Kinabalu, Malaysia, 2013.
- [7] C. Shi, W. Jinqiao, L. Yang, X. Changsheng, and L. Hanqing, "Fast feature selection and training for AdaBoost-based concept detection with large scale datasets," in *Proceedings of the international conference on Multimedia*, Firenze, Italy, 2010.
- [8] S. M. Abdelrahman, and A. Abraham, "Intrusion detection using error correcting output code based ensemble." pp. 181-186.
- [9] H. Wei, and H. Weiming, "Network-based intrusion detection using Adaboost algorithm." pp. 712-717.
- [10] B. Debojit, N. Bernard, and K. B. Dhruva, "Anomaly based intrusion detection using meta ensemble classifier," in *Proceedings of the Fifth International Conference on Security of Information and Networks*, Jaipur, India, 2012.
- [11] Z. Zhi-Hua, *Ensemble Methods: Foundations and Algorithms*: Chapman & Hall/CRC, 2012.
- [12] V. J. Finn, *Introduction to Bayesian Networks*: Springer-Verlag New York, Inc., 1996.
- [13] H. Simon, *Neural Networks: A Comprehensive Foundation*: Prentice Hall PTR, 1998.
- [14] J. R. Quinlan, *C4.5: programs for machine learning*: Morgan Kaufmann Publishers Inc., 1993.
- [15] R. G. Brereton, and G. R. Lloyd, "Support Vector Machines for classification and regression," *Analyst*, vol. 135, no. 2, pp. 230-267, 2010.
- [16] S. Galit, R. P. Nitin, and C. B. Peter, *Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner*: Wiley Publishing, 2010.
- [17] U. K. D. i. D. Archive.
- [18] H. Weiming, H. Wei, and S. Maybank, "AdaBoost-Based Algorithm for Network Intrusion Detection," *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, vol. 38, no. 2, pp. 577-583, 2008.