

A Survey and Taxonomy of Lightweight Intrusion Detection Systems

Sang Min Lee*
Korea Aerospace University
Goyang, South Korea
minuri33@kau.ac.kr

Dong Seong Kim
University of Canterbury
Christchurch, New Zealand
dongseong.kim@canterbury.ac.nz

Jong Sou Park
Korea Aerospace University
Goyang, South Korea
jspark@kau.ac.kr

Abstract

Internet and computer networks are exposed to an ever increasing number of security threats that can damage computer systems and communication channels. Firewalls are used to defend systems but still they are not enough to provide full protection to the systems. Then, the concern with Intrusion Detection Systems (IDSs) has been growing for network security over the past years. Due to the increasing of networks' speed and the amount of network traffic, it is essential that IDSs need to be lightweight to cope with it. Therefore, two representative methodologies have been applied to make IDSs lightweight, feature selection and parameter optimization. In this paper, we introduce concepts and algorithms of them and survey existing approaches which have used them. In particular, we review the previous approaches according to three broad categories: spam, Denial-of-Service (DoS) and Distributed Denial-of-Servive (DDoS) attacks detection since they are the most threatening intrusions these days. Finally, we conclude the survey by identifying trends and open challenges of lightweight IDSs research and development. Our hope is that this paper sheds some light on a fruitful direction of future research for lightweight IDSs.

Key Words: Intrusion detection system, parameters optimization, feature selection

1 Introduction

As the amount of information on networks has been increased tremendously, society has been relied on computer networks and then, network security is getting more essential. The confidentiality, integrity and availability of computer systems need to be defended against a number of threats. Especially, Denial-of-Service (DoS) and Distributed Denial-of-Service (DDoS) attacks which have badly damaged to systems should be paid attention. Also, spam is no more bothersome mail but threat since it recently includes virus attachments and spyware agents which compromise the recipients' system. Many security methods for protecting network systems, such as firewalls, access control, or encryption methods, have failed to protect computer systems and networks from sophisticated attacks and malwares. Among them, the Intrusion Detection Systems (IDSs) turn out to be the proper solution to this problem and become a critical security component. They detect automatically various kinds of intrusions or unwanted traffic against computer systems by monitoring and analyzing the behavior of users, computer systems or networks. They play a vital role in overall security fields. As new attacks appear and amount of audit data (or spam using bulk mailing tools) increases, IDSs should counteract them. In addition to this, as network becomes faster, the IDSs also should catch up the intrusions to cope with the increased network throughput and speed without degradation of detection rates. IDSs may utilize additional hardware such as network processor, System on Chip (SoC) and Field-Programmable Gate Array (FPGA) [37]. Additional hardware can increase packet capture speed and decrease processing time but it needs more costs and may not enhance detection rates of IDSs.

Journal of Internet Services and Information Security (JISIS), volume: 2, number: 1/2, pp. 119-131

*Corresponding author: Computer Engineering Department, Korea Aerospace University, 200-1, Hwajeon-dong, Dukyang-gu, Goyang-city, Gyeonggi-do, South Korea

Therefore, IDSs themselves should be lightweight while guaranteeing high detection rates. There are two main methods to realize such lightweight IDSs: (i) parameters optimization of intrusion detection model and (ii) feature selection. Firstly, parameters optimization is to find out optimal parameters of intrusion detection models based on various kinds of classification algorithms. Secondly, feature selection is used to find out only important features or feature set out of all the features of audit data. It can eliminate irrelevant features to avoid processing overheads.

In this paper, three lightweight intrusion detection models to detect spam, DoS, and DDoS attacks are introduced and analyzed. We focus on these threats since they are the most threatening intrusions these days. Parameters optimization and/or feature selection are adopted to make detection models lightweight. More detailed introductions for each detection model are described in the next Sections.

The rest of the paper is organized as follows. In Section 2, IDSs are briefly described. In Sections 3, parameters optimization and feature selection are presented in detail. Section 4 introduces lightweight spam detection models. Then, Section 5 states lightweight DoS attacks detection models are presented. We describe lightweight DDoS attacks detection models in Section 6. Finally, this survey is concluded with a presentation about the potential research arrears in the near future.

2 Intrusion Detection Systems

Intrusion is generally defined as violating confidentiality, Integrity, and Availability of computer or computer network system. Intrusion Detection Systems (IDSs) detect automatically computer intrusions to protect computers and computer networks safely from malicious uses or attacks. IDSs attempt to detect that an intruder breaks into computer system or legitimate user misuses system resources in real-time. According to intrusion detection model, IDSs can be divided into two major types, misuse detection and anomaly detection.

- Misuse detection (Signature detection) : It uses specifically known patterns of unauthorized behavior to detect intrusions. These specific patterns are called signatures. This method is quite strong to detect unknown intrusions. However, it has low degree of accuracy in detecting unknown intrusions since it relies on signatures extracted by human experts [29].
- Anomaly detection: It establishes a baseline of normal usage patterns. If it finds something that widely deviates from the baseline, the deviation is flagged as a possible intrusion. Although it is powerful to identify new types of intrusion as deviations from normal usage, a potential drawback is the high false alarm rate, i.e. previously unseen (yet legitimate) system behaviors may also be recognized as anomalies, and hence flagged as potential intrusions [29].

According to source of audit data, IDS can be divided into two main types, Host-based Intrusion Detection System (HIDS) and Network-based Intrusion Detection System (NIDS).

- HIDS : It is deployed in an individual host machine and can monitor audit data of a single host. HIDS uses a number of system characteristics to detect intrusions, such as file system, network events, system calls. A critical decision which is choosing the appropriate system characteristics involves a number of tradeoffs including the content of the data that is monitored, the volume of data that is captured, and the extent to which the IDS may modify the operating system of the host machine [5].
- NIDS : It monitors the packets sent and received by hosts. It reads raw packets off a network, usually after putting the network interface into promiscuous mode. The network interface in promiscuous mode will receive an entire traffic on the local network segment. Since the packets are not

actually addressed to the host that the NIDS resides on, the system is also impervious to an entire class of attacks such as the "ping-of-death" attack that can disable a host without ever triggering a HIDS [5].

3 Parameters Optimization and Feature Selection

There are two representative methods to build lightweight IDSs: (i) parameters optimization of data-mining and machine learning algorithms and (ii) feature selection of audit data.

- Parameters optimization : Data mining is the analysis of (often large) observational datasets to find patterns or models that are both understandable and useful to the data owner [16]. Data mining can efficiently generate classifiers, called intrusion detection models in IDSs, to detect attacks, especially for the vast amount of audit data. Parameters optimization is to find out optimal parameters of a detection model to build IDSs based on various kinds of data mining (classification) algorithms. For example, previous approaches have applied it to various kinds of classification and clustering algorithms including Support Vector Machine (SVM) [21][32], Random Forests (RF) [4][20][27][43] and modified Bayesian Additive Regression Trees (CBART) so on. Kim *et al.* [21] and Park *et al.* [32] regulated parameters of kernel function in SVM. Kim *et al.* [20] and Lee *et al.* [27] optimized two parameters of RF, *ntry* and *mtry*. Also, Abu-Nimeh *et al.* [1] found optimal number of trees and power parameter.
- Feature selection : IDSs process huge amount of audit data which contains a lot of features. However, all features are not essential to classify network audit data. In other words, some of these features are irrelevant or redundant because irrelevant and redundant features not only increase computational cost, such as time and overheads but also decrease the detection rates. Feature selection is used to find out only important features or feature set out of all the features of audit data. The selection of important features of audit data in IDSs is a significant issue. It can make IDSs lightweight and improve the classification performance. Two representative feature selection methods are wrapper [22][31] and filter [10][15]. Wrapper approach uses classification algorithms. It evaluates feature importance with respect to detection rates of detection model. The feature selection algorithm is wrapped inside the classification algorithm. Filter method does not utilize any classification algorithm to filter out the irrelevant and redundant features. It uses the primary characteristics of the training data to evaluate the relevance of the features or feature set by some independent measures such as distance measure, correlation measures, and consistency measures [2][11]. Furthermore, some studies [20][32] have proposed hybrid approaches which combine wrapper and filter methods.

4 Lightweight Spam Detection Models

The communication via an electronic mail (e-mail) is the most popular and useful service in the Internet since it is free, fast and easy to send. Concern about the proliferation of unsolicited bulk e-mail, commonly referred to as "spam", has been steadily increasing [9]. Before bulk mailing tools (bulk-mailers) have not been appeared, people receive a small amount of spam and it was not a significant problem. However, as the quantities of spam have been increased tremendously because of bulk-mailers, the recipients not only waste precious working time to delete spam mails but also become increasingly annoyed. In addition, Internet Service Providers (ISPs) have been deluged with complaints since spam generate considerable amount of network traffic and spam places a considerable burden on the system, i.e. it can

quickly fill up file server storage space. Moreover, virus attachments, spyware agents and phishing which are added in spam have become the most serious security threats to individuals and businesses recently.

As a result of this growing problem, spam detection could be considered to filter or delete spam mails automatically. A number of spam detection models using machine learning techniques have been proposed [1][6][23][28][39][41][45][46]. Existing detection models search for particular keyword patterns in e-mails by using machine learning algorithms. There are large amounts of emails and they contain a lot of keyword patterns which denote spam mails. They may burden spam detection system, so it should reduce the resources for processing to catch up with the huge amounts of e-mails. To reduce the amount of consuming resources with guaranteeing high detection rates, parameters optimization (e.g., threshold function value, the number of hidden layers in Artificial Neural Networks (ANN) and parameters of kernel function in SVM and so on) and feature selection (which figures out what features of mail are more important and need to be selected to detect spam mails) methods can be adopted.

4.1 Literature Review

Abu-Nimeh *et al.* [1] proposed a CBART method which is modified Bayesian Additive Regression Trees (BART) to make it applicable to classification. BART is primarily designed to predict quantitative (continuous) outcomes from observations via regression. They modified the current BART model and applied it to spam detection. They conducted parameters optimization which they used different numbers of trees ranging from 30 to 500 and also applied different power parameters for the tree prior, to specify the depth of the tree, ranging from 0.1 to 2.5. They found the optimal number of trees and power parameter. In their experimental results, CBART outperformed all the other classifiers and achieved the minimum average error rates. However, they did not apply feature selection.

Xie [41] used Support Vector Machine (SVM) in spam detection. They found two optimal parameters, cost and gamma. They used a good method for selecting proper values of them, which is called “grid search”, i.e. to search for the values of certain parameters over supplied parameter ranges. Although they performed parameters optimization, their detection rates were too low. They also did not perform feature selection as [1].

On the other hand, many researchers applied feature selection in spam detection. Multivariate Decision Tree (DT) based method for supervised linear feature extraction was presented by Bursteinas and Long [6]. Their approach is based on the wrapper model with randomized approach to the generation of the feature subsets. They carried out experiments with all features and selected features. Unfortunately, the result with selected features was worse than the result with all features.

Liang *et al.* [28] performed the feature selection based on a distance discriminant and feature ranking algorithm. Then, they used several detection algorithms such as DT, SVM, Nave Bayes (NB) and K Nearest Neighbor (KNN). The result of DT was the best but it was still very poor.

Thota *et al.* [39] proposed feature selection based by controlling False Discovery Rate (FDR). After feature selection using the proposed algorithm to control FDR, the classification performance of Logistic Regression (LR) classifier was improved. FDR controls the expected proportion of incorrectly rejected null hypotheses (errors). It could improve the efficiency of the regression classifier but the performance was not better than them of SVM and DT.

Zhao and Zhu [45] focused on an extended version of the rough set model called Variable Precision Rough Set (VPRS). They conducted feature selection according to the forward selection method. Zhu [46] dealt with feature selection using rough set theory, and then detected spam using SVM. However, both of them did not carry out enough experiments and compare with other approaches. Also, the detection rates were not outstanding.

They [6][28][39][45][46] conducted feature selection but not parameters optimization. Lee *et al.* [23] adopted both parameters optimization and feature selection simultaneously. For parameters optimization

of spam detection, two parameters ($mtry$ and $ntree$) of RF were regulated to increase spam detection rates. For feature selection, they provided variable importance of individual feature to select the important features on a scale between 0 and 1. This variable importance represents how each feature is significant for spam detection so that their approach could select relevant features and remove irrelevant ones. Then, the optimal number of selected features was figured out by using two methodologies: (a) only one parameters optimization during overall feature selection and (b) parameters optimization in every feature elimination phase.

Furthermore, Lee *et al.* [25] applied cost-sensitive measures that assign blocking legitimate messages (false positives) a higher cost than letting spam messages pass the spam detector (false negatives) to evaluate because false positives is more severe than false negatives.

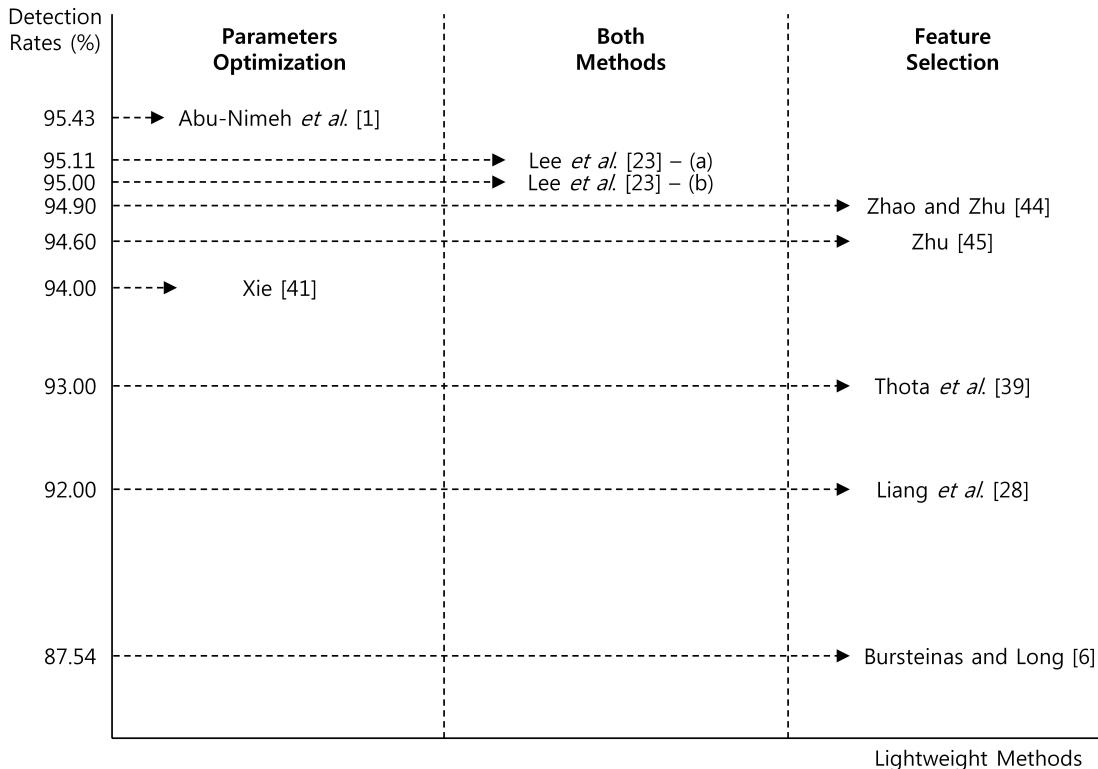


Figure 1: The experimental results of the previous spam detection approaches

4.2 Experiments and Discussions

In order to evaluate the detection rates of previous approaches, experiments were performed using same Spambase dataset. The Spambase dataset is an email message collection containing 4601 messages, being 1813 (39%) marked as spam, was created by Hopkins *et al.* [38]. The collection comes in pre-processed (not raw) form, and its instances have been represented as 58-dimensional vectors. The first 48 features are words extracted from the original messages, without stop list nor stemming, and selected as the most unbalanced words for the spam class. The next 6 features are the percentage of occurrences of the special characters “;”, “(”, “[”, “!”, “\$” and “#”. The following 3 features represent different measures of occurrences of capital letters in the text of the messages. Finally, the last feature is the class label which indicates spam or non-spam. Some researchers consider the Spambase dataset obsolete, as

it does not represent the state of practical spam mails, however, others consider it is a good test bed for evaluating learning techniques [44].

In Figure 1, the experimental results of the previous spam detection approaches are presented. The performance evaluation of various machine learning algorithms for detecting spam by applying parameters optimization, feature selection or both methods is focused. Abu-Nimeh *et al.* [1] and Xie [41] performed parameters optimization, [6][28][39][45][46] performed feature selection, and only Lee *et al.* [23] conducted both methods for building lightweight spam detection models. Abu-Nimeh *et al.* [1] shows the highest detection rates (95.43%) but they did not conduct feature selection. It means that they use all features for detecting spam mails and consume more computational resources. Furthermore, the detection rates of Lee *et al.* [23] were higher (95.52%) when they used 56 features. Even though the final experimental results of Lee *et al.* [23] shows a little degradation on detection rates, it is marginally small and less processing resources are consumed because they used much less features. The tradeoff between the detection rates and the processing resources should be considered more in the future works.

5 Lightweight DoS Attacks Detection Models

With the growth of the bandwidth and interconnectivity of computer networks, DoS attacks became a major intrusion to the Internet infrastructure integrity and one of the most devastating attacks possible to throw across the network. Unfortunately, automated tools make these attacks increasingly easier to execute and then they are becoming more sophisticated. A growing number of DoS attacks impose a significant threat on the availability of network services since DoS attacks attempt to overwhelm victim machines and it causes legitimate users to prevent from accessing their computing resources. Some DoS attacks target the bandwidth capabilities of computer systems while others target the machines' computational state. To foil the DoS attacks, a number of countermeasures are needed and detecting is the first step of countermeasures.

Many approaches based on very well-known detection (classification) algorithms, such as SVM, RF, and so on, have been proposed and focused on the achievable accuracy (detection rates). Since a simple DoS attack is normally that large amounts of traffic are generated and sent to a target machine, DoS attacks detection system should keep up with the large amounts of traffic. Thus, DoS attacks detection should be lightweight. Unlike previous spam detection approaches, there are many approaches which adopted parameters optimization and feature selection together.

5.1 Literature Review

Kim *et al.* [21] proposed feature selection method based on Genetic Algorithm (GA). GA built new chromosomes and searches the optimal detection model based on the fitness values obtained from the result of SVM classification. A chromosome was decoded into a set of features and parameters for a kernel function to be used by SVM classifier. The SVM was used to estimate the performance of a detection model represented by a chromosome. However, this approach costs much time to select features and results in a slow selecting process. Furthermore, it used detection rates as the criterion of evaluating IDSs. The detection rates are not a sufficient evaluation criterion for IDSs and for evaluation criterion of IDSs [8].

Park *et al.* [32] proposed Correlation-Based Hybrid Feature Selection (CBHFS) approach. GA was used to generate subsets of features from given feature set. CBHFS took full feature set as input and returned the optimal subset of feature after being evaluated by Correlation-Based Feature Selection (CFS) [14] and SVM. Each chromosome represented a feature vector. Merit of each chromosome was calculated by CFS. The chromosome having highest Merit represented the best feature subset in population.

This subset was then evaluated by SVM classification algorithm. Then, this procedure iterated with a new population of chromosomes, which is generated through performing genetic operations. The algorithm stopped if better subset was not found in next generation or when maximum number of generation was reached. All these methods yield good improvements but they are also fairly complex and computationally expensive as Kim *et al.*'s approach [21]. In other words, these two hybrid approaches [21][32] sometimes showed a little degradation on detection rates with more computations rather than the nave filter methods, did not provide the variable importance of features and were complicated to implement.

Chen *et al.* [8] conducted wrapper-based feature selection algorithm aiming at modeling lightweight IDSs. They used Modified Random Mutation Hill Climbing (MRMHC) as search strategy to specify a candidate subset for evaluation. Then, they adopted SVM as wrapper approach to obtain the optimum feature subset. Their MRMHC method can be enhanced in terms of speed compared to Kim *et al.* [21] and Park *et al.* [32]. They also used MRMHC to obtain the optimal parameters for kernels in SVM but this approach is still complex to implement.

A neurotree model for a classification engine and wrapper based feature selection algorithm for minimizing the computational complexity of the classifier were employed by Sivatha Sindhu *et al.* [36]. Then, they found optimal weight values to reduce the total error. However, the detection rates were not competent.

Zaman and Karray [42] proposed a lightweight IDS based on feature selection and IDS classification scheme. They applied Fuzzy Enhancing Support Vector Decision Function (Fuzzy ESVDF) to select appropriate feature set for the detection process. The Fuzzy ESVDF is iterative algorithm based on Support Vector Decision Function (SVDF) and Forward Selection (FS) with a fuzzy inferencing model. It is simple and efficient, and it does not require parameters optimization. Then, they used SVM and Neural Network (NM) for classifiers. Besides, the IDS classification scheme was divided into four types depending on the TCP/IP network model: Application Application layer IDS (AIDS), Transport layer IDS (TIDS), Network layer IDS (NIDS) and Link layer IDS (LIDS). Although the detection rates were high, the result was limited to TIDS. Therefore, they should integrate four different IDS types and evaluate the integrated IDS.

Zhang *et al.* [43] and Kim *et al.* [20] performed feature selection and parameter optimizations based on RF for lightweight DoS attacks detection. Feature importance ranking was performed according to the result of variable importance values. Then, irrelevant features are eliminated and only important features are selected. Zhang *et al.* [43] only cut off 3 features and optimized only *mtry* value. On the other hand, Kim *et al.* [20] eliminated much more irrelevant features and optimized both of *mtry* and *ntree*, two parameters of RF.

Moreover, one of the main problems of existing approaches is that IDSs provide only binary detection results: intrusion (attack) or normal. That is a main cause of high false rates and inaccurate detection rates in IDSs. If some attack or normal data are belonging to boundary, they may be classified wrong. To cope with it, Lee *et al.* [27] proposed Quantitative Intrusion Intensity Assessment (QIIA). It provides intrusion (or normal) quantitative intensity value. It is capable of representing how an instance of audit data is proximal to intrusion (DoS attacks) or normal in a numerical value such as "0.95" proximity to intrusion. It can be interpreted as the instance has a probability of 0.95 to be classified as an intrusion. This approach is very novel and refreshing paradigm. It can overcome the drawback of current binary detection and classify intrusions in more detail. For example, DoS attacks can be classified as Smurf, Neptune, Teardrop, etc.

5.2 Experiments and Discussions

Experiments were carried out on KDD 1999 dataset [17]. The dataset contains 41 features of each instance and 24 different types of attacks that fall into four main categories: DoS (Denial of Service),

R2L (unauthorized access from a remote machine), U2R (unauthorized access to root privileges) and Probing attacks. However, DoS attacks and normal instances were only used in these experiments since DoS attacks detection is only concerned in this paper.

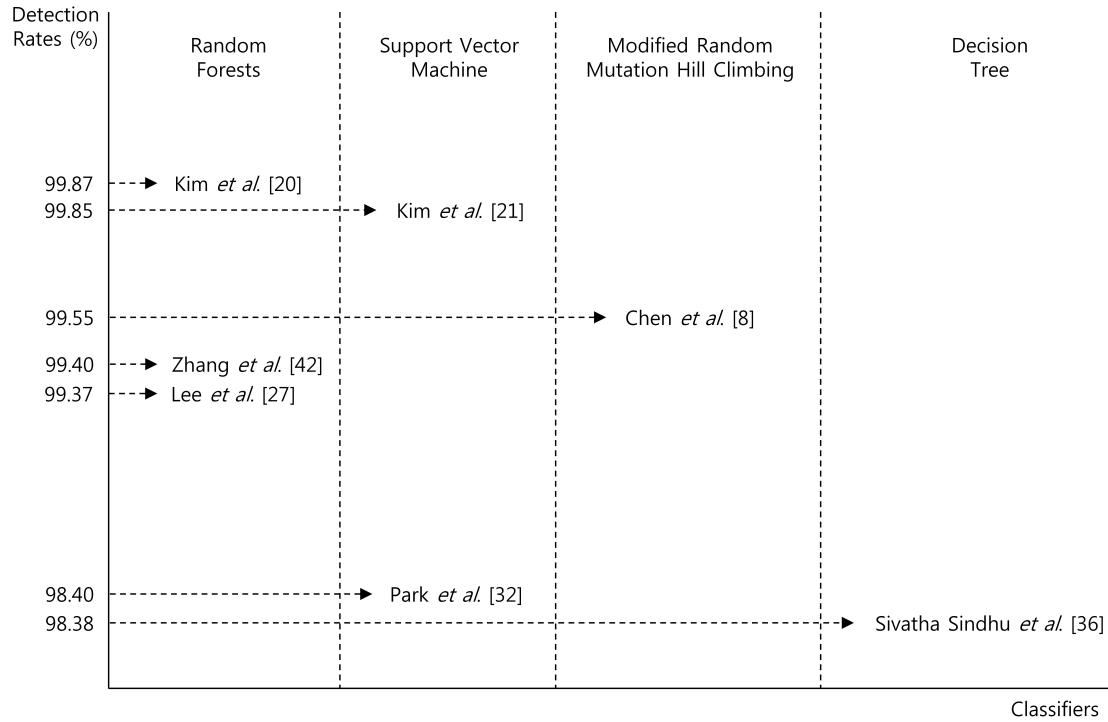


Figure 2: The experimental results of the previous DoS attacks detection approaches

The experimental results of the previous DoS attacks detection approaches are presented in Figure 2. They used same or different machine learning algorithms for classifiers. All approaches [8][20][21][27][36][43] conducted feature selection and parameters optimization simultaneously. Zaman and Karray [42] achieved 99.84% for detection rates but it was only for TIDS. Therefore, the result cannot be compared with others, so it is not described in Figure 1. Kim *et al.* [20] outperforms all the other approaches and achieved the maximum detection rates. As a result, feature selection using variable importance which is provided by RF is superior and the performance of RF is also outstanding. Lee *et al.* [27] is able to quantify the unknown audit data and classify attacks in more detail even though the degradation of detection rates is lower than Kim *et al.*'s [20]. Also, the difference is negligible. For this reason, Lee *et al.*'s QIIA approach [27] is a quite interesting research topic for the future works.

6 Lightweight DDoS Attacks Detection Models

DDoS attacks have emerged as one of the most significant threats among others [3]. DDoS attacks are large-scale and coordinated attacks targeting on the availability of services at a victim system or network resources. The intensity of DDoS attacks has become stronger through the development of network infrastructure and recent communication technology. DDoS attacks are normally launched by creating an extremely huge volume of attack traffic and they rapidly exhaust resources of target systems, such as network bandwidth and computing power. It is hard to detect and respond to DDoS attacks due to large and complex network environments and various types of DDoS attacks with different characteristics.

The sophisticated evolution of DDoS attacks techniques, the enhanced scale of Botnet and IP spoofing techniques also make detection difficult. To thwart DDoS attacks, they should be detected as soon as the attacks are launched. Thus, DDoS attacks detection system should guarantee both short detection delay and high detection rates with low false positives. Moreover, the DDoS attacks detection system should be lightweight to be able to deal with an enormous volume of real-time network traffics.

6.1 Literature Review

There have been researched lots of DDoS attacks detection approaches. They have been tried to detect DDoS attacks proactively.

Cabrera *et al.* [7] utilized network management systems for detecting DDoS attacks. Their approach depends only on information from Management Information Base (MIB) traffic variables intimating attack precursors. They showed how the relevant MIB traffic variables can be extracted automatically. There is a critical drawback that their approach is limited on one network management system domain. If attacker and victim are located in different domains, the correlations of variables cannot be detected. They should scale the approach to multi-domain environments.

Cluster analysis was performed for proactive detection of DDoS attacks in Kim *et al.*'s approach [24]. They selected detection parameters by observing the characteristics of DDoS attacks, and then entropy concept was adopted to analyze the traffic by using cluster analysis method. The approach can detect precursors of the DDoS attacks at early phases and be easily implemented since it uses only normalized distance. However, detection delay and detection rates were not provided in their experimental results, so it cannot be figured out whether the approach has reasonable success. Moreover, it is weak for new types of packet because it is customized for known attack.

Some previous approaches on anomaly detection rely on monitoring IP (internet protocol) attributes of incoming packets. Since these approaches use only few features, feature selection method is not needed. Only parameters optimization is needed to make DDoS attacks detection models lightweight.

Feinstein *et al.* [13] proposed a statistical detection model to identify DDoS attacks by computing entropy and frequency-sorted distributions of selected packet attributes. The entropy can be computed on consecutive packets. When a network is a normal state, the entropy values fall in a narrow range. On the other hands, the entropy values exceed the range when the network is under attacks. They implemented an entropy model as a plug-in for Snort [34] and performed experiments to validate it in various network traces. However, they did not provide an apparent method for optimizing the size of a sliding window. Furthermore, subnet spoofed attacks are not able to be detected by their approach since they did not concern about them.

Peng *et al.* [33] proposed a simple but robust detection scheme by monitoring the increase of new IP addresses. The approach monitors arrival rates of new source IP addresses and detects changes of them using Cumulative Sum (CUSUM) algorithm [40]. The CUSUM reduces the false positive and has good performance for a non-parametric model. Unlike Feinstein *et al.*'s approach [13], DDoS attacks including subnet spoofed IP addresses [12][30] can be detected by their approach. Besides, they achieved high detection accuracy with low computational overheads. But their experimental results showed that the detection delay was between 10 and 127.3 seconds which is not satisfactory in terms of the detection delay for a real-time DDoS attacks detection.

Kim *et al.* [18] collected a baseline profile on various attribute combinations. They analyzed actual traffic traces from two Internet traffic archives and verified the traffic stability in several aspects. Their research provides how to check whether a particular site has meaningful traffic stability, measure the stability within a site and decide the number of required traffic profiles. However, the combined attributes increased computational overheads so it is not suitable to detect real-time DDoS attacks.

Kim *et al.* [19] proposed a traffic matrix to detect DDoS attacks quickly and accurately. They

constructed a traffic matrix using source IP addresses of inbound traffic packets. Then, variance of the traffic matrix is computed and Weighted Moving Average (WMA) is adopted to decrease error rates. The WMA value of variance is compared with a certain threshold value to detect DDoS attacks. However, this approach could not find out to tune up parameters of the traffic matrix and time based window size, leading to computational overheads when DDoS attacks did not occur. Also, the proposed hash function creates many hash collisions.

To solve these problems, Lee *et al.* [26] proposed a revised traffic matrix. The traffic matrix is built up with packet based window size to reduce the computational overheads and a reformed hash function to reduce hash collisions. It made the model effectively in terms of processing overheads and detection delay. Thus, a revised approach could be used to detect DDoS attacks at the early stage in real-time. In addition, GA was adopted for optimization of parameters used in the traffic matrix. To increase detection rates, three parameters in their detection model were regulated: (i) size of traffic matrix, (ii) packet based window size and (iii) threshold value of variance from packet information.

All previous approaches have carried out experiments but their experimental environment and dataset were totally different. In this paper, therefore, their experimental results are not described since they cannot be compared with each other.

7 Concluding Remarks and Future Discussions

This survey provides a comprehensive overview of various lightweight IDSs. We discussed parameters optimization and feature selection which can be used to make IDSs lightweight. We classified the lightweight intrusion detection models according to the most threatening intrusions: spam, DoS and DDoS attacks. We made a thorough comparison on performance and delineated their similarities and differences.

In spam and DoS attacks detection, the previous approaches used datasets which have a lot of features. Many approaches adopted feature selection methods to implement lightweight IDSs. Among the feature selection methods, variable importance provided by RF is accurate and outperforms. Since it represents the importance of features by numerical values, irrelevant features can be eliminated very accurately. Also, parameters optimization methods are used to make IDSs lightweight. These methods depend on the data mining algorithms which are used for classification because they regulate the parameters of the algorithms. However, only one approach which adopted RF for a classifier was used the variable importance for feature selection. Thus, it needs to apply the variable importance to other classification algorithms and compare the results.

Furthermore, a novel QIIA approach was proposed to provide intrusion (or normal) quantitative intensity value. This approach used proximity values by RF. QIIA can be conducted by SVM since the proximity value can be replaced to the distance between hyperplane and support vectors.

In DDoS attacks detection, the previous approaches detected DDoS attacks using datasets that have only few features. So, the feature selection methods were not needed and parameters optimization methods were only conducted. Moreover, they used different dataset and even they created their own experimental dataset, so it was not possible to compare the results. Therefore, experiments with same dataset and experimental environments should be carried out to compare with each other.

Our hope is that this paper sheds some light on a fruitful direction of future research for lightweight IDSs.

References

- [1] S. Abu-Nimeh, D. Nappa, X. Wang, and S. Nair. Bayesian Additive Regression Trees-Based Spam Detection for Enhanced Email Privacy. In *Proc. of the 3rd Int. Conf. on Availability, Reliability and Security (ARES 2008), Barcelona, Spain*, pages 1044–1051, IEEE, March 2008.
- [2] H. Almuallim and T. G. Dietterich. Learning Boolean Concepts in the Presence of Many Irrelevant Features. *Artificial Intelligence*, 69(1-2):279–305, September 1994.
- [3] Arbor Networks. World Wide Infrastructure Security Report. 2008.
- [4] L. Breiman. Random Forests. *Machine Learning*, 45:5–32, October 2001.
- [5] D. J. Brown, B. Suckow, and T. Wang. A Survey of Intrusion Detection Systems. *Department of computer Science, University of California, USA*, 2001.
- [6] B. Bursteinas and J. A. Long. Transforming Supervised Classifiers for Feature Extraction. In *Proc. of the 12th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'00), Vancouver, British Columbia, Canada*, pages 274–280, IEEE, November 2000.
- [7] B. D. Cabrera, L. Lewis, X. Qin, W. Lee, R. K. Prasant, B. Ravichandran, and R. K. Mehra. Proactive Detection of Distributed Denial of Service Attacks Using MIB Traffic Variables-A Feasibility Study. In *Proc. of the 7th International Symposium on Integrated Network Management, Seattle, Washington, USA*, pages 1–14, IEEE, May 2001.
- [8] Y. Chen, W.-F. Li, and X.-Q Cheng. Toward Building Lightweight Intrusion Detection System through Modified RMHC and SVM. In *Proc. of the 15th IEEE Int. Conf. on Networks (ICON'07), Adelaide, Australia*, pages 83–88, IEEE, November 2007.
- [9] L. F. Cranor and B. A. LaMacchia. SPAM!. *Communications of the ACM*, 41(8):74–83, August 1998.
- [10] M. Dash, K. Choi, P. Scheuermann, and H. Liu. Feature Selection for Clustering - A Filter Solution. In *Proc. of the 2nd Int. Conf. on Data Mining (ICDM'02), Maebashi City, Japan*, pages 115–122, IEEE, December 2002.
- [11] M. Dash, H. Liu, and H. Motoda. Consistency Based Feature Selection. In *Proc. of the 4th Pacific Asia Conference on Knowledge Discovery and Data Mining (PAKDD'00), Kyoto, Japan, LNCS*, volume 1805, pages 98–109, Springer-Verlag, April 2000.
- [12] C. Douligeris and A. Mitrokotsa. DDoS Attacks and Defense Mechanisms: Classification and State-of-the-Art. *Computer Networks*, 44(5):643–666, April 2004.
- [13] L. Feinstein, D. Schnackenberg, R. Balupari, and D. Kindred. Statistical Approaches to DDoS Attack Detection and Response. In *Proc. of the 3rd DARPA Information Survivability Conf. and Exposition (DISCEX'03), Washington, District of Columbia, USA*, pages 303–314, IEEE, April 2003.
- [14] M. A. Hall. Correlation-Based Feature Selection for Discrete and Numeric Class Machine Learning. In *Proc. of the 17th International Conference on Machine Learning (ICML'00), Stanford, California, USA*, pages 359–366, Morgan Kaufmann, June 2000.
- [15] M. A. Hall and L. A. Smith. Feature Subset Selection: A Correlation Based Filter Approach. In *Proc. of the 4th Int. Conference on Neural Information Processing and Intelligent Information Systems (ICONIP'97), Dunedin, New Zealand*, pages 855–858, Elsevier, November 1997.
- [16] D. J. Hand, H. Mannila, and P. Smyth. Principles of Data Mining. The MIT Press, 2001.
- [17] KDD Cup 1999 Data, Available at <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.
- [18] Y. Kim, J.-Y. Jo, and K. Suh. Baseline Profile Stability for Network Anomaly Detection. In *Proc. of the 3rd International Conference on Information Technology: New Generations (ITNG'06), Las Vegas, Nevada, USA*, pages 720–725, IEEE, April 2006.
- [19] T. Kim, D. Kim, S. Lee, and J. Park. Detecting DDoS Attacks Using Dispersible Traffic Matrix and Weighted Moving Average. In *Proc. of the 3rd International Conference and Workshops on Advances in Information Security and Assurance (ISA'09), Seoul, Korea, LNCS*, volume 5567, pages 290–300, Springer-Verlag, June 2009.
- [20] D. Kim, S. Lee, and J. Park. Toward Lightweight Intrusion Detection System through Simultaneous Intrinsic Model Identification. In *Proc. of the 4th Int. Symp. on Parallel and Distributed Processing and Applications (ISPA'06), Sorrento, Italy, LNCS*, Vol. 4331, pages 981–989, Springer-Verlag, December 2006.

- [21] D. Kim, H.-N. Nguyen, S.-Y. Ohn, and J. Park. Fusions of GA and SVM for Anomaly Detection in Intrusion Detection System. In *Proc. of the 2nd International Symposium on Neural Networks (ISNN'05), Chongqing, China, LNCS*, volume 3498, pages 415–420, Springer-Verlag, May 2005.
- [22] R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1-2):273–324, December 1997.
- [23] S. Lee, D. Kim, J. Kim, and J. Park. Spam Detection Using Feature Selection and Parameters Optimization. In *Proc. of the 4th International Conference on Complex, Intelligent and Software Intensive Systems (CISIS'10), Krakow, Poland*, pages 883–888, IEEE, February 2010.
- [24] K. Lee, J. Kim, K. Kwon, Y. Han, and S. Kim. DDoS Attack Detection Method Using Cluster Analysis. *Expert Systems with Applications*, 34(3):1659–1665, April 2008.
- [25] S. Lee, D. Kim, and J. Park. Cost-Sensitive Spam Detection Using Parameters Optimization and Feature Selection. *Journal of Universal Computer Science*, 17(6):944–960, 2011.
- [26] S. Lee, D. Kim, and J. Park. Detection of DDoS Attacks Using Optimized Traffic Matrix. *Computers and Mathematics with Application*, 63(2):501–510, January 2012.
- [27] S. Lee, D. Kim, Y. Yoon, and J. Park. Quantitative Intrusion Intensity Assessment using Important Feature Selection and Proximity Metrics. In *Proc. of the 15th IEEE Pacific Rim Int. Symp. on Dependable Computing (PRDC'09), Shanghai, China*, pages 127–134, IEEE, November 2009.
- [28] J. Liang, S. Yang, and A. Winstanley. Invariant Optimal Feature Selection: A Distance Discriminant and Feature Ranking based Solution. *Pattern Recognition*, 41(1):1429–1439, January 2008.
- [29] V. Marinova. A Short Survey of Intrusion Detection Systems. *Problems of Engineering Cybernetics and Robotics*, 58:23–30, 2007.
- [30] J. Mirkovic and P. Reiher. A Taxonomy of DDoS Attacks and DDoS Defense Mechanisms. *ACM SIGCOMM Computer Comm. Rev.*, 34(2):39–53, April 2004.
- [31] S.-M. Noelia. A New Wrapper Method for Feature Subset Selection. In *Proc. of the European Symp. on Artificial Neural Networks (ESANN'05), Bruges, Belgium*, pages 515–520, April 2005.
- [32] J. Park, K. M. Shazzad, and D. Kim. Toward Modeling Lightweight Intrusion Detection System through Correlation-Based Hybrid Feature Selection. In *Proc. of the 1st SKLOIS Conf. on Information Security and Cryptology (CISC 2005), Beijing, China, LNCS*, volume 3822, pages 279–289, Springer-Verlag, December 2005.
- [33] T. Peng, C. Leckie, and R. Kotagiri. Proactively Detecting Distributed Denial of Service Attacks Using Source IP Address Monitoring. In *Proc. of the 3rd Int. IFIP-TC6 Networking Conf. (Networking'04), Athens, Greece, LNCS*, volume 3042, pages 771–782, Springer-Verlag, May 2004.
- [34] M. Roesch. Snort - Lightweight Intrusion Detection for Networks. In *Proc. of the 13th USENIX Conf. on Systems Administration (LISA'99), Seattle, Washington, USA*, pages 229–238, November 1999.
- [35] A. Shevtkar, K. Anantharam, and N. Ansari. Low Rate TCP Denial-of-Service Attack Detection at Edge Routers. *IEEE Communications Letters*, 9(4):pp. 363–365, April 2005.
- [36] S. Sivatha Sindhu, S. Geetha, and A. Kannan. Decision Tree Based Light Weight Intrusion Detection Using a Wrapper Approach. *Expert Systems with Applications*, 39:129–141, 2012.
- [37] H. Song and J. W. Lockwood. Efficient Packet Classification for Network Intrusion Detection using FPGA. In *Proc. of the ACM/SIGDA 13th Int. Symposium on Field-Programmable Gate Arrays, New York, New York, USA*, pages 238–245, ACM Press, January 2005.
- [38] Spambase Dataset, Available at <ftp://ftp.ics.uci.edu/pub/machine-learningdatabases/spambase>.
- [39] H. Thota, R. N. Miriyala, S. P. Akula, K. M. Rao, C. S. Vellanki, A. A. Rao, and S. Gedela. Performance Comparative in Classification Algorithms Using Real Datasets. *Journal of Computer Science and Systems Biology*, 2(1):97–100, January 2009.
- [40] H. Wang, D. Zhang, and K. G. Shin. Detecting SYN Flooding Attacks. In *Proc. of the 2002 IEEE INFOCOM, New York, USA*, pages 1530–1539, IEEE, June 2002.
- [41] Y. Xie. An Introduction to Support Vector Machine and Implementation in R. May 2007. available at http://yihui.name/cv/images/SVM_Report_Yihui.pdf

- [42] S. Zaman and F. Karray. Lightweight IDS based on Feature Selection and IDS Classification Scheme. In *Proc. of 12th International Conference on Computational Science and Engineering (CSE'09), Vancouver, Canada*, pages 365–370, IEEE, August 2009.
- [43] J. Zhang and M. Zulkernine. Network Intrusion Detection Using Random Forests. In *Proc. of 3rd Annual Conf. on Privacy, Security and Trust (PST'05), St. Andrews, New Brunswick, Canada*, pages 53–61, IEEE, October 2005.
- [44] C. Zhao. Towards Better Accuracy for Spam Predictions. *Technical Report, University of Toronto, Canada*, December 2004.
- [45] W. Zhao and Y. Zhu. Classifying Email Using Variable Precision Rough Set Approach. In *Proc. of the 1st International Conference on Rough Sets and Knowledge Technology (RSKT'06), Chongqing, China, LNCS*, volume 4062, pages 766–771, Springer-Verlag, July 2006.
- [46] Z. Zhu. An Email Classification Model Based on Rough Set and Support Vector Machine. In *Proc. of the 5th Int. Conf. on Fuzzy Systems and Knowledge Discovery (FSKD'08), Shandong, China*, pages 236–240, IEEE, October 2008.



Sang Min Lee received the B.S. degree in Information and Telecommunication Engineering from Korea Aerospace University, Republic of Korea in 2005. And he received M.S. and Ph.D. degree in Computer Engineering from Korea Aerospace University, Republic of Korea in 2007, 2012, respectively. He is currently a part-time lecturer in Computer Engineering Department at Korea Aerospace University since March 2012. His research interests are in network security, especially intrusion detection systems, and survivability of cloud computing.



Dong Seong Kim received the B.S. degrees in Electronic Engineering from Korea Aerospace University, Republic of Korea in 2001. And he received M.S. and Ph.D. degree in Computer Engineering from Korea Aerospace University, Republic of Korea in 2003, 2008, respectively. And he was a visiting researcher in University of Maryland at College Park, USA in 2007. He was a postdoctoral researcher in Duke University from June 2008 to July 2011. He is currently a lecturer in Computer Science and Software Engineering Department at the University of Canterbury, New Zealand since August 2011. His research interests are in dependable and secure systems and networks. In particular, intrusion detection systems, security for wireless ad hoc and sensor networks, virtualization, and cloud computing systems.



Jong Sou Park received the M.S. degree in Electrical and Computer Engineering from North Carolina State University in 1986. And he received his Ph.D in Computer Engineering from The Pennsylvania State University in 1994. From 1994 - 1996, he worked as an assistant Professor at The Pennsylvania State University in Computer Engineering Department and he was president of the KSEA Central PA, Chapter. He is currently a full professor in Computer Engineering Department, Korea Aerospace University. His main research interests are information security, embedded system and hardware design. He is a member of IEEE and IEICE and he is an executive board member of the Korea Institute of Information Security and Cryptology and Korea Information Assurance Society