

A Neural Network Based System for Intrusion Detection and Classification of Attacks

Mehdi MORADI and Mohammad ZULKERNINE

Abstract-- With the rapid expansion of computer networks during the past decade, security has become a crucial issue for computer systems. Different soft-computing based methods have been proposed in recent years for the development of intrusion detection systems. This paper presents a neural network approach to intrusion detection. A Multi Layer Perceptron (MLP) is used for intrusion detection based on an off-line analysis approach. While most of the previous studies have focused on classification of records in one of the two general classes - normal and attack, this research aims to solve a multi class problem in which the type of attack is also detected by the neural network.

Different neural network structures are analyzed to find the optimal neural network with regards to the number of hidden layers. An early stopping validation method is also applied in the training phase to increase the generalization capability of the neural network. The results show that the designed system is capable of classifying records with about 91% accuracy with two hidden layers of neurons in the neural network and 87% accuracy with one hidden layer.

Index Terms—Artificial Neural Networks, Intrusion Detection, Multilayer Perceptron, Training Strategies.

I. INTRODUCTION

THE rapid development and expansion of World Wide Web and local network systems have changed the computing world in the last decade. However, this outstanding achievement has an Achilles' heel: The highly connected computing world has also equipped the intruders and hackers with new facilities for their destructive purposes. The costs of temporary or permanent damages caused by unauthorized access of the intruders to computer systems have urged different organizations to increasingly implement various systems to monitor data flow in their networks [14]. These systems are generally referred to as Intrusion Detection Systems (IDSs).

There are two main approaches to the design of IDSs. In a misuse detection based IDS, intrusions are detected by looking for activities that correspond to known signatures of intrusions or vulnerabilities. On the other hand, an anomaly detection based IDS detects intrusions by searching for abnormal network traffic. The abnormal

traffic pattern can be defined either as the violation of accepted thresholds for frequency of events in a connection or as a user's violation of the legitimate profile developed for his/her normal behavior.

One of the most commonly used approaches in expert-system based intrusion detection systems is rule-based analysis using Denning's [1] profile model. Rule-based analysis relies on sets of predefined rules that are provided by an administrator or created by the system. Unfortunately, expert systems require frequent updates to remain current. This design approach usually results in an inflexible detection system that is unable to detect an attack if the sequence of events is even slightly different from the predefined profile. The problem may lie in the fact that the intruder is an intelligent and flexible agent while the rule-based IDSs obey fixed rules. This problem can be tackled by the application of soft computing techniques in IDSs.

Soft computing is a general term for describing a set of optimization and processing techniques that are tolerant of imprecision and uncertainty. The principal constituents of soft computing techniques are Fuzzy Logic (FL), Artificial Neural Networks (ANNs), Probabilistic Reasoning (PR), and Genetic Algorithms (GAs) [15]. The idea behind the application of soft computing techniques and particularly ANNs in implementing IDSs is to include an intelligent agent in the system that is capable of disclosing the latent patterns in abnormal and normal connection audit records, and to generalize the patterns to new (and slightly different) connection records of the same class.

In the present study, an off-line intrusion detection system is implemented using Multi Layer Perceptron (MLP) artificial neural network. While in many previous studies [2], [3], [10] the implemented system is a neural network with the capability of detecting normal or attack connections, in the present study a more general problem is considered in which the attack type is also detected. This feature enables the system to suggest proper actions against possible attacks. The promising results of the present study show the potential applicability of ANNs for developing practical IDSs.

Different structures of MLP are examined to find a minimal architecture that is reasonably capable of classification of network connection records. The results show that even an MLP with a single layer of hidden neurons can generate satisfactory classification results. Because the generalization capability of the IDS is

This work was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC).

M. Moradi is with the School of Computing, Queen's University, Kingston, Ontario, Canada K7L 3N6 (e-mail: moradi@cs.queensu.ca).

M. Zulkernine is with the School of Computing, Queen's University, Kingston, Ontario, Canada K7L 3N6 (e-mail: mzulker@cs.queensu.ca).

critically important, the training procedure of the neural networks is carried out using a validation method that increases the generalization capability of the final neural network.

Paper Organization: Section I has introduced the basic ideas in intrusion detection and the motivations for this study. Section II reviews some basic ideas in neural network theory and presents an overview of some of the previous studies that have applied neural networks in intrusion detection. Section III deals with the dataset, attack types, and the features used for classifying network connection records in this study. Section IV describes the implementation procedure and training-validation method. Section V presents the experimental results and Section VI concludes the paper with a discussion of the results and possibilities for future work.

II. ARTIFICIAL NEURAL NETWORKS (ANNs) IN INTRUSION DETECTION REVIEW STAGE

The ability of soft computing techniques for dealing with uncertain and partially true data makes them attractive to be applied in intrusion detection. Some studies have used soft computing techniques other than ANNs in intrusion detection. For example, genetic algorithms have been used along with decision trees to automatically generate rules for classifying network connections [13]. However, ANNs are the most commonly used soft computing technique in IDSs [2], [4], [6], [10], [11].

An ANN is an information processing system that is inspired by the way biological nervous systems, such as the brain, process information. It is composed of a large number of highly interconnected processing elements (neurons) working with each other to solve specific problems. Each processing element (neuron) is basically a summing element followed by an activation function. The output of each neuron (after applying the weight parameter associated with the connection) is fed as the input to all of the neurons in the next layer. The learning process is essentially an optimization process in which the parameters of the best set of connection coefficients (weights) for solving a problem are found and includes the following basic steps [8]:

- Present the neural network with a number of inputs (vectors each representing a pattern)
- Check how closely the actual output generated for a specific input matches the desired output.
- Change the neural network parameters (weights) to better approximate the outputs.

Some IDS designers exploit ANN as a pattern recognition technique. Pattern recognition can be implemented by using a feed-forward neural network that has been trained accordingly. During training, the neural network parameters are optimized to associate outputs (each output represents a class of computer network

connections, like normal and attack) with corresponding input patterns (every input pattern is represented by a feature vector extracted from the characteristics of the network connection record). When the neural network is used, it identifies the input pattern and tries to output the corresponding class. When a connection record that has no output associated with it is given as an input, the neural network gives the output that corresponds to a taught input pattern that is least different from the given pattern [6].

The most commonly reported application of neural networks in IDSs is to train the neural net on a sequence of information units, each of which may be an audit record or a sequence of commands. The input to the net consists of the current command and the past w commands (w is the size of window of commands under examination). Once the net is trained on a set of representative command sequences of a user, it constitutes (learns) the profile of the user and when put in action, it can discover the variance of the user from its profile [4], [6]. Usually recurrent neural networks are used for this purpose.

Ryan et al. [3] described an off-line anomaly detection system (NNID) which utilized a back-propagation MLP neural network. The MLP was trained to identify users' profile and at the end of each log session, the MLP evaluated the users' commands for possible intrusions (off-line). The authors described their research in a small computer network with 10 users. Each feature vector described the connections of a single user during a whole day. 100 most important commands are used to describe a user's behavior. They used a 3 layer MLP (2 hidden layers). The MLP identified the user correctly in 22 cases out of 24.

Cannady [2] used a three layer neural network for offline classification of connection records in normal and misuse classes. The system designed in this study was intended to work as a standalone system (not as a preliminary classifier whose result may be used in a rule-based system). The feature vector used in [2] was composed of nine features all describing the current connection and the commands used in it. A dataset of 10,000 connection records including 1,000 simulated attacks was used. The training set included 30% of the data. The final result is a two class classifier that succeeded in classification of normal and attack records in 89-91% of the cases.

In yet another study [10], the authors used three and four layer neural networks and reported results of about 99.25% correct classification for their two class (normal and attack) problem.

Cunningham and Lippmann [11] used ANNs in misuse detection. They used an MLP to detect Unix-host attacks by searching for attack specific keywords in the network traffic. Different groups used self-organizing maps (SOM) for intrusion detection [5].

In most of the previous studies [2], [3], [10], the

implemented systems were neural networks with two possible outputs: normal or anomaly. In these studies, some types of attacks and a set of normal records were included in the dataset; however, the output of the neural network was 1 or 0 for normal or attack conditions (the attack type was not determined by the neural network). The present study is aimed to solve a multi class problem in which not only the attack records are distinguished from normal ones, but also the attack type is identified.

III. EVALUATION DATASET: ATTACK TYPES AND FEATURES

The 1999 version of MIT Lincoln Laboratory - DARPA (Defense Advanced Research Projects Agency) intrusion detection evaluation data was used in this research [16]. The sample version of the dataset included more than 450,000 connection records. A subset of the data that contained the desired attack types and a reasonable number of normal events were selected manually. The final dataset used in this study included 20,055 records.

A. Attack Types

There are at least four different known categories of computer attacks including denial of service attacks, user to root attacks, remote to user attacks and probing attacks [9]. Two different attack types were included in the dataset used for this study: *SYN Flood (Neptune)* and *Satan*. These two attack types were selected from two different attack categories (denial of service and probing) to check for the ability of the intrusion detection system to identify attacks from different categories. Availability of enough data records was the other factor in choosing these two specific types. Furthermore, there are studies that have used the same attack types [2]. Therefore, evaluation of the results by comparing them to previous studies was possible. In the following paragraphs, a description of the attack types is provided.

SYN Flood (Neptune) is a denial of service attack to which every TCP/IP implementation is vulnerable (to some degree). For distinguishing a Neptune attack network traffic is monitored for a number of simultaneous SYN packets destined for a particular machine. The host sending these packets is usually unreachable [9].

Satan is a probing intrusion which automatically scans a network of computers to gather information or find known

vulnerabilities. The network probes are quite useful for attackers planning a future attack [9].

Table 1 shows detailed information about the number of records from normal and two attack types included in training, validation, and testing sets. There were 9,830 records of normal connections, 7,051 records of Neptune attack, and 3,174 records of Satan attack in the dataset.

B. Features: Selection, Numerical Representation, and Normalization

In DARPA dataset each event (connection) is described with 41 features. 22 of these features describe the connection itself and 19 of them describe the properties of connections to the same host in last two seconds. In many attack scenarios, the signature of the attack record is identified through examination of some features in a sequence of records. Therefore, the IDS should analyze the service types used by the same user in previous connections and for this purpose these 19 features describing past events in the computer network were included in the feature vector.

A complete description of all 41 features is available [10], [16]. Instead of describing all the features, here we divide them into three groups and provide descriptions and examples for each group.

Group 1 includes features describing the *commands* used in the connection (instead of the commands themselves). These features describe the aspects of the commands that have a key role in defining the attack scenarios. Examples of this group are number of file creations, number of operations on access control files, number of root accesses, etc..

Group 2 includes features describing the *connection specifications*. This group includes a set of features that present the technical aspects of the connection. Examples of this group include: protocol type, flags, duration, service types, number of data bytes from source to destination, etc..

Group 3 includes features describing the *connections to the same host in last 2 seconds*. Examples of this group are: number of connections having the same destination host and using the same service, % of connections to the current host that have a rejection error, % of different services on the current host, etc..

During inspection of the data it turned out that the values of six features (land, urgent, num_failed_logins, num_shells, is_host_login num_outbound_cmds) were constantly zero over all data records (see [10] for descriptions). Clearly these features could not have any effect on classification and only made it more complicated and time consuming. They were excluded from the data vector. Hence the data vector was a 35 dimensional vector.

Different possible values for selected features were extracted and a numerical value was attributed to each of them. For example, for the protocol type the possible

TABLE 1.
DISTRIBUTION OF DATA VECTORS IN DIFFERENT SUBSETS FOR
TRAINING, VALIDATION, AND TESTING SETS
(THERE WERE TOTALLY 20,055 VECTORS IN THE DATASET).

Record Types	Training SET	Validation Set	Test Set
Normal	5,922	300	3,608
Neptune	4,430	300	2,321
Satan	1,807	300	1,067

numerical values were: tcp=0, udp=1, icmp=2. This numerical representation was necessary because the feature vector fed to the input of the neural network has to be numerical.

The ranges of the features were different and this made them incomparable. Some of the features had binary values where some others had a continuous numerical range (such as duration of connection). As a result, the features were *normalized* by mapping all the different values for each feature to [0, 1] range.

IV. IMPLEMENTATION: TRAINING AND VALIDATION METHOD

The present study was aimed to solve a multi class problem. Here, a three class case is described which can be extended to cases with more attack types. An output layer with three neurons (output states) was used: [1 0 0] for normal conditions, [0 1 0] for Neptune attack and [0 0 1] for the Satan attack. The desired output vectors used in training, validation, and testing phases were simply as mentioned above. In practice, sometimes the output of the neural network showed other patterns like [1 1 0] which were considered irrelevant. It is straightforward to show that there are 6 possible irrelevant cases.

In this paper, a three layer¹ neural network means a neural network with two hidden layers (the input layer is not counted because it acts just like a buffer and no processing takes place in it; however, the output layer is counted). The universal approximation theorem states that an MLP (with one or more hidden layers) can approximate any function with arbitrary precision and of course the price is an increase in the number of neurons in the hidden layer [8]. The question is if anything is gained by using more than one hidden layer. One answer is that using more than one layer may lead to more efficient approximation or to achieving the same accuracy with fewer neurons in the neural network.

The performance of a 2 layer neural network is seldom reported in the previous studies as described in Section II. One of the objectives of the present study is to evaluate the possibility of achieving the same results with this less complicated neural network structure. Using a less complicated neural network is more computationally efficient. Also it would decrease the training time.

MATLABTM Neural Network Toolbox [12] was used for the implementation of the MLP networks. Using this tool one can define specifications like number of layers, number of neurons in each layer, activation functions of neurons in different layers, and number of training epochs. Then the training feature vectors and the corresponding desired outputs can be fed to the neural network to begin training.

¹ There are different traditions for naming neural networks based on the number of hidden neuron layers [8].

All the implemented neural networks had 35 input neurons (equal to the dimension of the feature vector) and three output neurons (equal to the number of classes). Number of the hidden layers and neurons in each were parameters used for the optimization of the architecture of the neural network. Error back-propagation algorithm was used for training.

A. The Over-fitting Problem

One problem that can occur during neural network training is over-fitting. In an over fitted ANN, the error (number of incorrectly classified patterns) on the training set is driven to a very small value, however, when new data is presented, the error is large. In these cases, the ANN has memorized the training examples; however, it has not learnt to generalize the solution to new situations.

One possible solution for the over-fitting problem is to find the suitable number of training epochs by trial and error. In this study, the training time was too long (25 hours in the first experiment). Therefore, it was not reasonable to find the optimal number of epochs by trial and error. A more reasonable method for improving generalization is called early stopping. In this technique, the available data is divided into three subsets. The first subset is the training set, which is used for training and updating the ANN parameters. The second subset is the validation set. The error on the validation set is monitored during the training process. The validation error will normally decrease during the initial phase of training similar to the training set error. However, when the ANN begins to over-fit the data, the error on the validation set will typically begin to rise. When the validation error increases for a specified number of iterations, the training is stopped, and the weights that produced the minimum error on the validation set are retrieved [12]. In the present study, this training-validation strategy was used in order to maximize the generalization capability of the ANN.

V. EXPERIMENTAL RESULTS

The first implemented intrusion detector was a three layer MLP (two hidden layers with 35 neurons in each). This structure is referred to as: {35 35 35 3}. At this stage,

TABLE 2.
CORRECT CLASSIFICATION RATES IN THREE DIFFERENT TRAINING-VALIDATION-TESTING SESSIONS FOR THE {35 35 35 3} MLP NEURAL NETWORK. THE EARLY STOPPING VALIDATION METHOD IS APPLIED; THEREFORE, THE NUMBER OF TRAINING EPOCHS IS NOT THE SAME IN DIFFERENT SESSIONS.

Training Session	CORRECT CLASSIFICATION ON TRAINING SET	Correct Classification on Test Set
1	98.2	89.2
2	98.1	90.9
3	96.9	90.3
Average	97.46	90.13

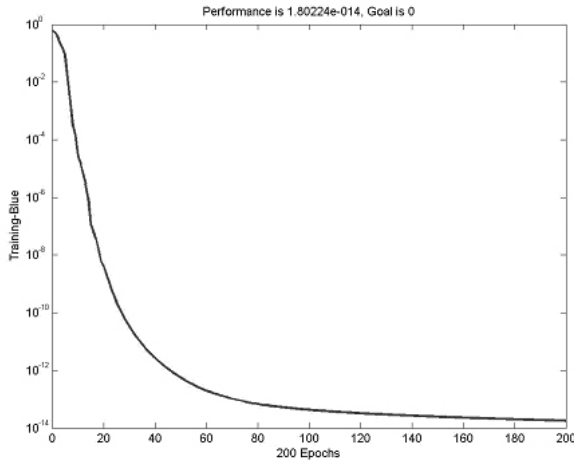


Fig. 1. The mean square error of the back-propagation training procedure versus training epochs for a 3 layer neural network {35 35 35 3}. The decrease in the error was completely satisfactory. The network was over-fitted.

early stopping validation was not applied and the training was performed for 200 times. The training process took more than 25 hours. Figure 1 shows the mean square error of the back propagation training process versus the progress of training epochs. The error clearly decreased to an outstanding level (comparable to zero). Therefore, it was expected to have good classification results. The final correct classification rate on training set confirmed this theory: it was very close to 100%. However, when unseen data (test set) was fed to the neural network, the result was undesirable. The correct classification rate was less than 80%.

A. Application of Early Stopping Validation Method

The initial result was a clear indication of over-fitting of the neural network (a description of this problem is presented in Section IV.A). As explained, the reasonable solution was to define a validation data set and monitor the classification error on this data set while the neural network was being trained.

The validation set used in this study consisted of 900 data records (300 of each class). The same neural network {35 35 35 3} was trained this time by applying early stopping validation method. Figure 2 shows the error of the training process versus progress of training epochs for one training session. The error on the training set (darker curve) was decreasing after epoch number 45; however, the training process was stopped because the error on the validation set was constant for ten epochs.

As expected, the correct classification rate on the training set declined slightly (98% compared to 100% in the first experiment). Instead, when unseen data (test set) was fed to the neural network the result was considerably better than the first experiment in which the early stopping method was not applied. The correct classification rate was more than 90% showing an 11% increase (from 80% in the first

experiment). There was another advantage associated with application of early stopping method: the training time was decreased because the number of training epochs was restricted by early stopping. The training-validation time in this implementation was less than 5 hours which is an improvement over 25 hours training time in the first experiment.

Because of the stochastic nature of the neural networks, it is usually common to report results of multiple training-testing procedures. Table 2 illustrates the results of three training-validation-testing sessions of {35 35 35 3} MLP used in this study. The correct classification results are reported separately for training and test data sets.

B. Two Layer Neural Network

As described in Section IV, one of the objectives of this study was to evaluate the possibility of application of a two layer neural network (one hidden layer) in the classification of normal and attack records. The back-propagation training algorithm becomes more complicated and consumes more memory and time when a hidden layer is added to the neural network. Furthermore, the resulting neural network is more complicated and less memory efficient. Therefore, it is always desirable to solve the problem with a simpler classifier.

The best two layer neural network used in this study was {35 45 35}. The early stopping validation method was applied. The best result was attained in a training session that was stopped on 48th epoch. The result was 93.1% correct classification on training and 87% on the testing set. The result of multiple training sessions also led to an average of 86% correct classification on unseen data.

Although the classification efficiency of the best two layer neural network was less than the best three layer neural network, the difference was just 4%.

C. Discussion

There were three categories of incorrect outputs: false positive, false negative, and irrelevant neural network output. The irrelevant outputs were those that did not represent any of the output classes in the data set (normal, Neptune attack, Satan attack). While in a two state neural network implemented with one output neuron there is no irrelevant output state, in a three output neural network, there are 6 irrelevant states. An analysis showed that in the three layer neural network with 90.9% correct classification, more than half of the incorrect results were from the category of irrelevant results. The number of incorrect classifications of this category can be decreased by classifying each irrelevant pattern in the class corresponding to the output neuron that has the highest value of activation function.

Are the results presented in the previous section satisfactory? To answer this question, they should be compared to the results of similar studies. In a previous

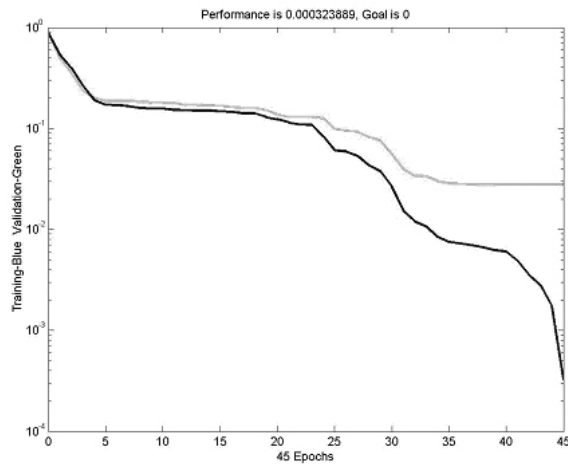


Fig. 2. The training process error when the early stopping validation method is applied (for the same network as in figure 1). The darker curve shows the error on the training set and the brighter curve presents the error on validation set.

study [10], a result of more than 99% correct classification on this dataset using the neural network structure {41-40-40-1} was reported. However, a two class problem was solved in which the records were classified either in normal or in attack classes. In another similar study with different dataset [2], the success rate was comparable to the results of the present study (89-99%) and again a two class problem was implemented.

VI. CONCLUSION AND FUTURE WORK

An approach for a neural network based intrusion detection system, intended to classify the normal and attack patterns and the type of the attack, has been presented in this paper. We applied the early stopping validation method which increased the generalization capability of the neural network and at the same time decreased the training time. It should be mentioned that the long training time of the neural network was mostly due to the huge number of training vectors of computation facilities. However, when the neural network parameters were determined by training, classification of a single record was done in a negligible time. Therefore, the neural network based IDS can operate as an *online* classifier for the attack types that it has been trained for. The only factor that makes the neural network off-line is the time used for gathering information necessary to compute the features.

A two layer neural network was also successfully used for the classification of connection records. Although the classification results were slightly better in the three layer network, application of a less complicated neural network was more computationally and memory wise efficient.

From the practical point of view, the experimental results imply that there is more to do in the field of artificial neural network based intrusion detection systems. The implemented system solved a three class problem. However, its further development to several classes is

straightforward. As a possible future development to the present study, one can include more attack scenarios in the dataset. Practical IDSs should include several attack types. In order to avoid unreasonable complexity in the neural network, an initial classification of the connection records to normal and general categories of attacks can be the first step. The records in each category of intrusions can then be further classified to the attack types.

REFERENCES

- [1] D. E. Denning, "An intrusion detection model," *IEEE Transactions on Software Engineering*, vol. 13, no. 2, pp. 222–232, 1987.
- [2] James Cannady, "Artificial neural networks for misuse detection," Proceedings of the 1998 National Information Systems Security Conference (NISSC'98), Arlington, VA, 1998.
- [3] J. Ryan, M. Lin, and R. Miikkulainen, "Intrusion Detection with Neural Networks," *AI Approaches to Fraud Detection and Risk Management: Papers from the 1997 AAAI Workshop*, Providence, RI, pp. 72-79, 1997.
- [4] K. Fox, R. Henning, J. Reed, and R. Simonian, "A neural network approach towards intrusion detection," Proceedings of 13th National Computer Security Conference, Baltimore, MD, pp. 125-134, 1990.
- [5] P. Lichodziejewski, A.N. Zincir Heywood, and M. I. Heywood, "Host-based intrusion detection using self-organizing maps," *Proceedings of the 2002 IEEE World Congress on Computational Intelligence*, Honolulu, HI, pp. 1714-1719, 2002.
- [6] H. Debar, M. Becker, and D. Siboni, "A neural network component for an intrusion detection system," Proceedings of 1992 IEEE Computer Society Symposium on Research in Security and Privacy, Oakland, California, pp. 240 – 250, 1992.
- [7] David Poole, Alan Makworth, and Randi Goebel, *Computational Intelligence*, New York: Oxford University Press, 1998.
- [8] Sergios Theodorios and Konstantinos Koutroumbas, *Pattern Recognition*, Cambridge: Academic Press, 1999.
- [9] Kristopher Kendall, "A Database of Computer Attacks for the Evaluation of Intrusion Detection Systems," *Masters Thesis, MIT*, 1999.
- [10] Srinivas Mulkamala, "Intrusion detection using neural networks and support vector machine," *Proceedings of the 2002 IEEE International Honolulu, HI*, 2002.
- [11] R. Cunningham and R. Lippmann, "Improving intrusion detection performance using keyword selection and neural networks," *Proceedings of the International Symposium on Recent Advances in Intrusion Detection*, Purdue, IN, 1999.
- [12] MATLAB online support:
www.mathworks.com/access/helpdesk/help/techdoc/matlab.shtml.
- [13] C. Sinclair, L. Pierce, and S. Matzner, "An application of machine learning to network intrusion detection," Proceedings of 15th Annual Computer Security Applications Conference (ACSAC '99), Phoenix, AZ, pp. 371-377, 1999.
- [14] R. A. Kemmerer and G. Vigna, "Intrusion detection: a brief history and overview," *Computer*, vol. 35, no. 4, pp. 27–30, 2002.
- [15] Piero P. Bonissone, "Soft computing: the convergence of emerging reasoning technologies," *Soft Computing Journal*, vol.1, no. 1, pp. 6-18, Springer-Verlag 1997.
- [16] MIT Lincoln Laboratory, <http://www.ll.mit.edu>.