
Mutating network scans for the assessment of supervised classifier ensembles

JAVIER SEDANO and SILVIA GONZÁLEZ, *Instituto Tecnológico de Castilla y León, C/López Bravo 70, Pol. Ind. Villalonquejar, 09001, Spain.*
E-mail: javier.sedano@itcl.es; silvia.gonzalez@itcl.es

ÁLVARO HERRERO and BRUNO BARUQUE, *Department of Civil Engineering, University of Burgos, University of Burgos, C/ Francisco de Vitoria s/n, 09006 Burgos, Spain.*
E-mail: ahcosio@ubu.es; bbaruque@ubu.es

EMILIO CORCHADO, *Departamento de Informática y Automática, Universidad de Salamanca, Plaza de la Merced, s/n, 37008 Salamanca, Spain; Department of Computer Science, VSB-Technical University of Ostrava, Ostrava, Czech Republic.*
E-mail: escorchado@usal.es

Abstract

As it is well known, some Intrusion Detection Systems (IDSs) suffer from high rates of false positives and negatives. A mutation technique is proposed in this study to test and evaluate the performance of a full range of classifier ensembles for Network Intrusion Detection when trying to recognize new attacks. The novel technique applies mutant operators that randomly modify the features of the captured network packets to generate situations that could not otherwise be provided to IDSs while learning. A comprehensive comparison of supervised classifiers and their ensembles is performed to assess their generalization capability. It is based on the idea of confronting brand new network attacks obtained by means of the mutation technique. Finally, an example application of the proposed testing model is specially applied to the identification of network scans and related mutations.

Keywords: Network intrusion detection, IDS performance, classifier ensembles, machine learning, zero-day attacks, mutation.

1 Introduction and motivation

The ever-changing nature of attack strategies and technologies is certainly one of the most harmful issues of attacks and intrusions. It greatly increases the difficulty of protecting computer networks and systems. As is well known, Intrusion Detection Systems (IDSs) figure prominently among the various kinds of tools proposed to date to confront such problems [3, 58]

In the context of computer networks, an IDS can roughly be defined as a tool that is designed to analyze the computer data network, in order to detect suspicious patterns that may be related to a network or system attack, assuming that some other security (prevention) tools may have failed. The main purpose of an IDS is to alert the security personnel to any suspicious and possibly intrusive event taking place in the system that is under analysis. Intrusion detection (ID) is therefore a field that focuses on the identification of attempted or ongoing attacks on a computer system (Host IDS) or network (Network IDS).

2 Network ID

A standard characterization of IDSs, based on their detection method, or model of intrusions, defines the following paradigms:

- Anomaly-based IDSs (also known as behaviour-based IDSs): the IDS detects intrusions by looking for activity that differs from the previously defined ‘normal’ behaviour of users and/or systems. In keeping with this idea, the observed activity is compared against ‘predefined’ profiles of expected normal usage. It is assumed that all intrusive activities are necessarily anomalous. However, in real-life environments, instead of their being identical, the set of intrusive activities only intersects the set of anomalous activities in certain cases. As a consequence [57], some anomalous activities that are not intrusive are flagged as intrusive (i.e. false positives) and intrusive activities that are not anomalous are not flagged up (i.e. false negatives). Anomaly-based IDSs can support detection of novel (zero-day) attack strategies, but may suffer from a relatively high rate of false positives [49].
- Misuse-based IDSs (also known as knowledge-based IDSs): intrusions are detected by checking activity that corresponds to known intrusion techniques (signatures) or system vulnerabilities. Misuse-based IDSs are therefore commonly known as signature-based IDSs. They detect intrusions by exploiting the available knowledge on specific attacks and vulnerabilities. As opposed to anomaly detection, misuse detection assumes that each intrusive activity can be represented by a unique pattern or signature [46]. The main problem with this approach is that intrusions with signatures that have not been archived by the system cannot be detected. As a consequence, a misuse-based IDS will never detect a zero-day attack [46]. The completeness of such IDSs requires regular updating of their knowledge of attacks.
- Specification-based ID: programme behavioural specifications reflect system policies that are used as a basis to detect attacks [6].

Some authors claim that one of the main challenges for researchers seeking to implement and to validate a new ID method is to assess and compare its performance with other available approaches [26]. Benchmarking an IDS is a common way of establishing its effectiveness, although this task is neither easy nor clear-cut [44, 2].

Some interesting works have proposed IDS testing in real (or as realistic as possible) environments [37] and in experimental settings [23, 43]. These two approaches have been proven to have some advantages and disadvantages, as pointed in [4], in which a combination of them is proposed.

The most widely accepted approach to IDS benchmarking is based on a publicly available data set. In this approach, the performance of a new IDS can be compared against those of previous IDSs by analyzing the same data set. Two of the most well-known data sets available for such purposes are DARPA [33] and the associated KDD Cup 1999 [32]. DARPA initiated its ID Evaluation Program in 1998 [21]. This programme generated background traffic interlaced with malicious activity so that IDSs and algorithms can be tested and compared. Certain problems concerning the inner structure and features of these data sets have been identified [36, 48]. They also have some other up-to-date problems, such as the outdated attacks and the data rate, which is not comparable with that in real current networks [9]. Despite these limitations, this is still the *de facto* standard corpus for IDS testing so far, as there are important constraints (privacy and confidentiality mainly) that prevent a new network traffic data set from being made publicly available.

Apart from testing data sets, some methodologies [4, 43], frameworks [23, 24] and tools [34, 42] for assessing IDSs have also been proposed up to now. Additionally, comparative studies on commercial IDSs have been presented, mainly in network journals. Despite the fact that increasing effort has been devoted over recent years to IDSs in general, and to IDSs assessment in particular, there is

no comprehensive and scientifically rigorous methodology to test the effectiveness of these systems [30]. Thus, it remains an open issue and a significant challenge [26].

In view of the above-mentioned issues, a fair comparison between different methods for ID is difficult to achieve. In consequence, a mutation-based technique [20, 18] is proposed in this study to check the generalization capability of different classifiers and classifier ensembles [10, 52]. This novel testing technique, which specializes in IDSs that rely on numerical packet features, is based on measuring the evaluated IDS results when they encounter unknown anomalous situations, as the identification of zero-day (previously unseen) attacks is possibly the most challenging and important ID problem [18, 20, 44]. In this case, the method is used to assess a comprehensive set of classifiers and classifier ensembles in the detection of mutated network scans. The ability to detect such scans can help to identify wider and potentially more dangerous threats to a network. This testing model has the main advantage of providing the classifiers with brand new attacks—network scans in this case.

Three types of scans (or sweeps) were defined [56]: network scans, port scans and the hybrid block scans [56]. Unlike other attacks, scans must use a real source IP address, because the results of the scan (open ports or responding IP addresses) must be returned to the attacker [45].

A port scan (or sweep) may be defined as series of messages sent to different port numbers to gain information on its activity status. These messages can be sent by an external agent attempting to access a host to find out more about the network services that this host is providing. So, a scan is an attempt to count the services running on a machine (or a set of machines) by probing each port for a response. A network scan is one of the most common techniques used to identify services that might then be accessed without permission [1]. This work focuses on the identification of network scans, in which the same port is the target for a number of computers in an IP address range.

The remaining sections of this article are structured as follows: Section 2 introduces the mutation testing technique. The applied classifiers and classifier ensembles are then described in Section 3 and the experimental results in Section 4. Finally, the conclusions of this study along with future work are discussed in Section 5.

2 A mutation testing technique for IDSs

Misuse IDSs rely on models of known attacks. The effectiveness of these IDSs depends on the ‘goodness’ of their models. This is to say, if a model of an attack does not cover all the possible modifications, the performance of the IDS will be greatly impaired. As previously stated, this is one of the main drawbacks of the misuse-based ID approach.

This study focuses on verifying the performance of many machine-learning models (classifiers and classifier ensembles) when applied to ID and when confronted with unknown (not previously seen) anomalous situations that can lead to a network attack. To generate such brand new situations, a mutation testing model is applied. This mutation testing model was previously proposed for a visualization-based IDS [18, 20] and mainly consists on mutating attack traffic. In general, a mutation can be defined as a random change. In keeping with this idea, the testing model modifies different features of the numerical information extracted from the packet headers.

The modifications created by this model may involve changes in such aspects as: attack length (amount of time that each attack lasts), packet density (number of packets per time unit), attack density (number of attacks per time unit) and time intervals between attacks. The mutations can also concern both source and destination ports, varying between the different three ranges of TCP/UDP port numbers: well known (from 0 to 1,023), registered (from 1,024 to 49,151) and dynamic and/or private (from 49,152 to 65,535).

4 Network ID

Time is a further issue of great importance when considering intrusions since the chances of detecting an attack increase in relation to their duration. There are therefore two main strategies to defeat IDSs:

- to make a significant reduction in the time used to perform a scan and
- to spread the packets out over time, which is to say, to reduce the number of packets sent per time unit, so they are more likely to slip past unnoticed.

In this study, the mutations are applied to data that relate to network scans. Some of the possible mutations may be meaningless, such as a sweep of very few hosts in the case of a network scan. Taking that into account, changes can be made to attack packets addressing the following issues:

- number of scans in the attack (i.e. number of addressed port numbers);
- destination port numbers at which scans are aimed;
- time intervals when scans are performed; and
- number of packets (density) forming the scans (number of scanned hosts).

3 Classifiers and ensembles for ID

As previously explained, one of the most interesting features of an ideal IDSs would be automatic detection of traffic that constitutes an attack against the network, as opposed to normal traffic. Automated learning techniques are algorithms specifically designed to take decisions with regard to new data that are presented on the basis of information extracted from previous data.

Current study is centred on the use of supervised algorithms [11] to detect new anomalies. Although some previous researches have proposed classifier ensembles for ID [13, 38, 39]; in this case, ensembles of supervised learning algorithms are presented to assess their generalization capability. The intention is to ascertain whether, once a certain type of attack is discovered, such information can be incorporated in the IDS, in order to detect new attacks in the future—with a certain similarity to those detected previously.

Usually, most automated learning algorithms suffer from common problems: like; over-fitting of the data used for training—and therefore, poor generalization capabilities, getting stuck on local minima in their learning function and high computational complexity when dealing with complex data [41]. One of the most widespread and useful techniques to avoid such problems is the ensemble learning scheme [54]. The main idea behind this kind of meta-algorithm is to train several, less complex classifiers and combine their results in a single and usually more complete classifier, which will generate even better results [47].

In this study, several different algorithms have been considered both for the base classifiers and for the ensemble construction, in order to obtain a significant wide array of possible algorithms to compare their performance results on mutated data sets. Among the base classifiers, all of which belong to the family of instance-based statistical classification algorithms, an array of techniques has been studied with different approaches to that problem; the instance-based k -Nearest Neighbours (IBk), classification trees (such as the Simple Classification and Regression Decision Tree (CART) and the REP-Tree) and Artificial Neural Networks such as the Radial Basis Function Network [10, 11].

Among the ensemble meta-algorithms that use the above-mentioned simple algorithms, this study tests basic algorithms, such as the MultiClass Classifier: used to adapt binary classifiers to multi-class problems; Bagging; Adaptive Boosting (AdaBoost); and Random Forest [41]. Their results are compared with more modern boosting algorithms, such as the LogitBoost [25] or the StackingC [60].

All of these techniques have been previously used to good effect, in order to solve practical problems. The Multi-class classifier has been used in the field of image recognition [50] or bio-informatics [28]; Bagging and AdaBoost are among the best known algorithms that have been used in a wide array of fields, such as food quality assessment [17], financial risk assessment [59], robotics [5], image [12] or signal processing [7], weather conditions forecasting [8] or operational costs optimization [51], to name but a few.

As explained in previous studies [53], ensemble meta-algorithms work well, among other reasons, because their use makes it relatively easy to increase the number of hypotheses used to classify new data. This increase is achieved by partitioning the data used to train the different components of the ensemble, allowing each of them to concentrate on a specific data subspace within the overall space of the problem. In the case of ID, this aspect is very interesting as an intrusion can be regarded as the kind of data that the classifier is unable to recognize. This could be because the data have never been presented to the classifier before or because the size of the attack data is very small when compared to normal data, which the classifier fails to notice. By partitioning the data set to construct different—and hopefully sufficiently diverse—classifiers, the probability of all composing classifiers failing to notice data decreases; also decreasing the probability of the classifier ensemble not noticing an attack.

As results prove, ensemble learning adds an important value to the analysis, as almost all variants consistently improve on the results obtained by their single classifier equivalent.

4 Experiments and results

This section describes the data sets used for evaluating the proposed testing method and how they were generated before detailing their results.

4.1 Data sets

Real-life data sets have previously been applied to perform ID [18, 20]. These data sets were generated ‘made-to-measure’ in a small-size (28 hosts) network where ‘normal’ traffic was known in advance. The captured packets relate to 63 different protocols. Further details about the network in which the traffic was captured are unavailable due to the security policy. This data set generation methodology has been described above in detail [18].

A set of features extracted from the packet headers characterized packets travelling along the monitored network. Once codified, the following five features contribute to building up the input vectors of the machine-learning models ($x \in \mathfrak{R}^5$) are:

- Timestamp: the time when the packet was sent;
- Source port: the port number of the device that sent the packet;
- Destination port: the port number of the target host, i.e. the host to which the packet is sent;
- Protocol ID: an integer number that identifies the protocol over TCP of the packet; and
- Size: the packet size (in Bytes).

It has been proven that these low-dimensional data sets allow for the detection of some anomalous situations, mainly those related to SNMP (Simple Network Management Protocol) [18, 29]. Different data sets were generated, containing different examples of SNMP anomalous situations. SNMP [14, 22] was chosen because it was ranked as one of the top five most vulnerable services by CISCO [16]. Especially the first two versions [14, 15] of this protocol that still are widely used

at present. SNMP attacks were also listed by the SANS Institute as one of the top 10 most critical Internet security threats [31, 40]. In addition to the SNMP packets, the data sets contain traffic related to some other protocols, considered as ‘normal’ traffic. As this network was isolated and protected from external attacks, ‘normal’ traffic was known in advance and has been empirically tested.

The three main anomalous situations related to the SNMP are distributed throughout the different data sets in this study, namely: network scans, SNMP community searches and MIB (Management Information Base) information transfers.

A hacker generates an SNMP community search by sending SNMP queries to the same port number of different hosts, and by using different strategies (brute force, dictionary, etc.) to guess the SNMP community string. Once the community string has been obtained, all the information stored in the MIB is available to the intruder. The unencrypted community string can be seen as the SNMP password for versions 1 and 2. In fact, it is the only SNMP authentication mechanism. As reported in [31], *The default community string used by the vast majority of SNMP devices is ‘public’, with a few clever network equipment vendors changing the string to ‘private’ for more sensitive information.* This facilitates SNMP intrusions as intruders can guess the community string with little effort.

MIB information transfer involves a transfer of some (or all the) information contained in the SNMP MIB, generally through the *Get* command or similar primitives such as *GetBulk* [35, 55]. This kind of transfer is potentially quite a dangerous situation because anybody possessing some free tools, basic SNMP knowledge and the community string (in SNMP versions 1 and 2), will be able to access all sorts of interesting and sometimes useful information. As specified by the Internet Activities Board, the SNMP is used to access MIB objects. Thus, protecting a network from malicious MIB information transfer is crucial.

By performing the SNMP community searches between the port/network scans and the MIB transfers, all anomalous situations constitute an SNMP attack from the outset. After attempting to establish whether and where (hosts and port numbers) SNMP is running, the community string is guessed, in order to access the information contained in the MIB. When the community string has been found, the MIB information is read.

The experimental setting of the present study comprises two different data sets:

- Data set 1: contains examples of the three anomalous situations (scans, community searches and MIB transfers). All packets have been labelled according to the following six classes:
 - (1) normal traffic;
 - (2) scans to port number 161 (SNMP default port number);
 - (3) scans to port number 162 (SNMP default port number);
 - (4) scans to port number 3,750 (alternative port number);
 - (5) MIB information transfer [19]; and
 - (6) community search.
- Data set 2: this data set was generated by applying the mutation testing technique (see Section 2) to Data set 1. As a result, it only contains examples of one of the anomalous situations, namely network scans. The network scans aimed at 162 and 3,750 were mutated, whereas the other two anomalous situations and the scan to port number 161 were not present in Data set 2. As a consequence of the mutations, the scan to port number 161 was removed, the scan to port number 162 was redirected to port number 1,434 and the remaining one is now aimed at port number 65,788. Moreover, the mutations also modified the number of packets in the network scans. Therefore, the scans of this data set are labelled as Classes 3 (scan to port number 1,434) and 4 (scan to port number 65,788). This data set has been labelled in the same way as data set 1 (Table 1).

TABLE 1. Grouping of classes to the training and the validation data set

Situations in the network	Number of classes and labels in each class			
	6 classes	4 classes	3 classes	2 classes
Normal traffic	1	1	1	1
Scans to port number 161 (SNMP default port number)	2	2	2	2
Scans to port number 162 (SNMP default port number)	3			
Mutation to port number 1,434	4			
Scans to port number 3,750. Mutation to port number 65,788	5	3	3	
MIB information transfer	6	4		
Community search				

4.2 Experiments

A combination of classifier ensembles was used to identify the new anomalous situations. In the experiment, a 10-fold cross-validation schema was selected. The final classification rate presented is obtained using the two previously described data sets. The 10-fold cross-validation is performed for each classifier/ensemble using only Data set 1 (original data set), and the best performing classifier is selected. The final classification accuracy for each model is then calculated over Data set 2 (the one which contains the mutated network scans). The number of samples used for training and validation was 5,866 and 64, respectively. WEKA software [27] was used to train and classify the data sets with ensembles and classifiers.

Each ensemble uses 10 base classifiers of the same type. The experimentation comprised a total of 1,534 tests from the combination of 25 ensembles: including ‘Bagging’, ‘Boosting’, ‘Adaboost’, ‘Random Subspaces’, ‘Decorate’, ‘Rotation Forest’, etc.; 59 classifiers: including ‘NaiveBayes’, ‘Ibk’, ‘LinearRegression’, ‘JRip’, ‘RBFNetwork’, ‘SMO’, etc.; adding the baseline classifiers tests (i.e. without ensembles). The best results are presented in the following section.

4.3 Results

The classifier ensembles try to differentiate between the six possible situations to identify these anomalous situations in the network, thereby labelling the training data sets in different ways. Thus, two different labels were assigned to differentiate attacks from normal traffic; three different classes were assigned to differentiate among MIB and community search, scans to all ports and normal traffic; and finally, four different classes were defined to obtain differences between normal traffic, community search, scans and MIB transfer. Table 1 presents the number of classes and the labels in each class, the situations in the network and the grouping of the classes.

Tables 2–5 show the training and the validation/classification performance of the applied models. The values are obtained for each combination ensemble–classifier. The tables present the worst indices of training and classification in the chosen sets of classifiers.

TABLE 2. Classification results for normal traffic and other situations in the network (two classes)

Ensembles	Classifiers	Correctly classified instances (%)
FilteredClassifier	Naive Bayes, NaiveBayesSimple, NaiveBayesUpdateable, IBK, IBI, PART AODE, RBF network, LMT, NBTree	Training (99.8636); classification (100)
Grading	SmlpleCart, Jrip, REPTree, PART, logistic, SMO, SPegasos, Ridor, ADTree, BFTree, FT, LADTree, LMT, NBTree, RandomForest	Training (99.284); classification (100)
MultiScheme	SmlpleCart, JRip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, FT, LADTree, LMT, NBTree	Training (99.2499); classification (100)
OrdinalClassClassifier	SmlpleCart, JRip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, FT, LADTree, LMT, NBTree	Training (99.2499); classification (100)
RandomCommittee	SmlpleCart, BFTree, RandomForest	Training (99.9318); classification (100)
RotationForest	SmlpleCart, JRip, REPTree, Logistic, SPegasos, OmrR, BFTree, FT, RandomForest	Training (99.3181); classification (100)
StackingC	SmlpleCart, JRip, REPTree, PART, Logistic, SMO, SPegasos, Ridor, ADTree, BFTree, FT, LADTree, LMT, NBTree, RandomForest	Training (99.2499); classification (100)
RacedIncrementalLogitBoost	LinearRegression, SMOReg, ConjunctiveRule, DecisionTable, M5Rules, DecisionStump, M5P	Training (98.8919); classification (100)
Bagging	SmlpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, Ridor, ADTree, BFTree, FT, LADTree, LMT, NBTree, RandomForest	Training (99.3011); classification (100)
AdaBoost	SmlpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, ConjunctiveRule, NNge, Ridor, ADTree, BFTree, DecisionStump, J48, NBTree, RandomForest	Training (99.3863); classification (100)
LogitBoost	REPTree, LinearRegression, Isotonic Regression, LinearRegression, ConjunctiveRule, DecisionTable, M5Rules, DecisionStump	Training (96.7951); classification (100)
MultiBoostAB	SmlpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, Ridor, ADTree, BFTree, FT, J48, LADTree, LMT, NBTree, RandomForest	Training (99.4033); classification (100)

RandomSubSpace	SimpleCart, Jrip, REPTree, PART, Ridor, ADTree, BFTree, FT, J48, LADTree, LMT, RandomForest, RandomTree	Training (99.9489); classification (100)
Dagging	PART, Logistic, SPegasos, ADTree, FT, J48, LADTree, LMT, RandomForest, RandomTree	Training (99.6591); classification (100)
Decorate	Logistic, SMO, SPegasos, DecisionTable, ADTree, FT, LADTree, LMT, RandomForest	Training (99.284); classification (100)
MultiClassClassifier	SimpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, FT, LADTree, LMT, NBTree	Training (99.284); classification (100)
ClassificationViaRegression	REPTree, LinearRegression, PaceRegression, DecisionTable, M5Rules, M5P	Training (96.6758); classification (100)
CVParameterSelection	SimpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, FT, LADTree, LMT, NBTree	Training (99.284); classification (100)
AttributeSelectedClassifier	IBK, SimpleCart, Jrip, REPTree, IB1, PART, VotedPerceptron, DTNB, DecisionTable, NNge, OneR, Ridor, ADTree, BFTree, FT, J48, LADTree, LMT, RandomForest, RandomTree	Training (82.3219); classification (100)
ThresholdSelector	SimpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, LADTree, LMT, NBTree	Training (82.5264); classification (100)
Vote	SimpleCart, Jrip, REPTree, PART, Logistic, SMO, SPegasos, ADTree, BFTree, FT, LADTree, LMT, NBTree	Training (99.284); classification (100)

TABLE 3. Classification results for normal traffic, scans and other situations in the network (three classes)

Ensembles	Classifiers	Correctly classified instances (%)
FilteredClassifier	Naive Bayes, NaiveBayesSimple, NaiveBayesUpdateable, IBK, IB1, RBF network, JRip, REPTree, PART, Id3, AODE, WAODE, Logistic, SMO, LBR, RBF network, LMT, NNge, OneR, Ridor, FT, J48, LADTree, RandomForest, RandomTree	Training (99.8636); classification (100)
Grading	SImpleCart, REPTree, PART, SMO, OneR, Ridor, BFTree, LADTree, RandomForest	Training (99.5397); classification (100)
MultiScheme	SImpleCart, REPTree, PART, SMO, OneR, Ridor, BFTree, LADTree	Training (99.5397); classification (100)
OrdinalClassClassifier	PART, LADTree, NBTree	Training (99.9148); classification (100)
RandomCommittee	SImpleCart, BFTree, RandomForest	Training (99.983); classification (100)
StackingC	SImpleCart, REPTree, PART, OneR, Ridor, BFTree, LADTree, RandomForest	Training (99.8636); classification (100)
RacedIncrementalLogitBoost	REPTree, ConjunctiveRule, DecisionStump, M5P	Training (99.8977); classification (100)
Bagging	SImpleCart, REPTree, PART, SMO, OneR, Ridor, BFTree, RandomForest	Training (99.5909); classification (100)
Adaboost	SImpleCart, PART, SMO, ConjunctiveRule, NNge, Ridor, BFTree, RandomForest, RandomTree	Training (99.5909); classification (100)
LogitBoost	ConjunctiveRule, M5Rules, DecisionStump	Training (99.9318); classification (100)
MultiboostAB	SImpleCart, REPTree, PART, Ridor, BFTree, RandomForest, RandomTree	Training (99.9148); classification (100)
RandomSubSpace	SImpleCart, REPTree, PART, DTNB, DecisionTable, OneR, Ridor, BFTree, FT, J48, LADTree, LMT, NBTree, RandomForest, RandomTree	Training (99.7272); classification (100)
Dagging	PART, DTNB, DecisionTable, BFTree, J48, LADTree, RandomForest, RandomTree	Training (99.6761); classification (100)
Decorate	SImpleCart, SMO, OneR, LADTree, RandomForest, RandomTree	Training (99.5397); classification (100)
MultiClassClassifier	PART, LADTree, NBTree	Training (99.9318); classification (100)
ClassificationViaRegression	DecisionTable, M5Rules, M5P	Training (99.9659); classification (100)
CVPParameterSelection	SImpleCart, REPTree, PART, SMO, OneR, Ridor, BFTree, LADTree	Training (99.5397); classification (100)
AttributesSelectedClassifier	IBK, SImpleCart, REPTree, IB1, PART, DTNB, DecisionTable, NNge, OneR, Ridor, BFTree, FT, J48, LADTree, LMT, RandomForest, RandomTree	Training (99.7443); classification (100)
Vote	SImpleCart, REPTree, SMO, OneR, Ridor, BFTree, LADTree	Training (99.5397); classification (100)

TABLE 4. Classification results for normal traffic, scans, MIB and community search (four classes)

Ensembles	Classifiers	Correctly classified instances (%)
FilteredClassifier	NaiveBayes, IBK, SimpleCart, JRip, REPTree, Ib1, PART, Id3, AODE, NaiveBayesSimple, NaiveBayesUpdateable, WAODE, SMO, NNge, OneR, Ridor, BFTree, FT, J48, LADTree, LMT, NBTree, RandomForest, RandomTree	Training (99.8636); classification (100)
Grading	SimpleCart, REPTree, SMO, OneR, BFTree, LADTree, NBTree, RandomForest	Training (99.5568); classification (100)
MultiScheme	SimpleCart, REPTree, SMO, OneR, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.5568); classification (100)
OrdinalClassClassifier	SimpleCart, REPTree, PART, ADTree, BFTree, LADTree, NBTree, RandomForest	Training (99.8807); classification (100)
RandomCommittee	SimpleCart, BFTree, RandomForest, RandomTree	Training (99.983); classification (100)
RotationForest	OneR, Ridor	Training (99.8977); classification (100)
StackingC	SimpleCart, REPTree, OneR, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.8636); classification (100)
RacedIncrementalLogitBoost	REPTree, M5Rules, DecisionStump, M5P	Training (99.8807); classification (100)
Bagging	SimpleCart, REPTree, SMO, NNge, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.5568); classification (100)
Adaboost	SimpleCart, REPTree, Jrip, J48, LADTree, NBTree, RandomForest, RandomTree	Training (99.9489); classification (100)
LogitBoost	DecisionTable, M5Rules, M5P	Training (99.9489); classification (100)
MultiboostAB	SimpleCart, REPTree, OneR, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.8636); classification (100)
RandomSubSpace	SimpleCart, IBk, REPTree, IB1, DecisionTable, NNge, OneR, Ridor, BFTree, FT, LADTree, LMT, RandomForest, PART, DecisionTable, BFTree, J48, LADTree, NBTree, RandomForest, RandomTree	Training (99.8125); classification (100)
Dagging	SMO, OneR, Ridor, LADTree, NBTree	Training (99.6931); classification (100)
Decorate	PART, LADTree, NBTree	Training (99.5568); classification (100)
MultiClassClassifier	DecisionTable, MRules, M5P	Training (98.7385); classification (100)
ClassificationViaRegression	SimpleCart, REPTree, SMO, OneR, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.9659); classification (100)
CVParameterSelection	IBK, SimpleCart, REPTree, IB1, Logistic, DecisionTable, NNge, OneR, Ridor, BFTree, FT, LADTree, LMT	Training (99.5568); classification (100)
AttributeSelectedClassifier	SimpleCart, REPTree, SMO, OneR, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (98.551); classification (100)
Vote		Training (99.5568); classification (100)

TABLE 5. Classification results for the whole classes (six classes)

Ensembles	Classifiers	Correctly classified instances
FilteredClassifier	NaiveBayesSimple, NaiveBayesUpdateable, LBR	Training (99.8636); classification (100)
Grading	SImpleCart, REPTree, BFTree, LADTree, NBTree, RandomForest	Training (99.9148); classification (100)
MultiScheme	SImpleCart, JRip, REPTree, BFTree, LADTree, NBTree, RandomTree	Training (99.8807); classification (100)
OrdinalClassifier	SImpleCart, PARTADTree, BFTree, NBTree, RandomForest	Training (99.8977); classification (100)
RandomCommittee	SImpleCart, BFTree, RandomForest	Training (99.9659); classification (100)
RotationForest	REPTree, J48	Training (99.983); classification (100)
StackingC	SImpleCart, REPTree, BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.9148); classification (100)
RacedIncrementalLogitBoost	REPTree, LWL, M5Rules, DecisionStump, M5P	Training (99.8466); classification (100)
Bagging	SImpleCart, JRip, REPTree, BFTree, LADTree, LMT, NBTree, RandomForest, RandomTree	Training (99.8977); classification (100)
Adaboost	Jrip, J48, LADTree, NBTree, RandomTree	Training (99.9148); classification (100)
LogitBoost	REPTree, SimpleLinearRegression, LWL, ConjunctiveRule, M5Rules, DecisionStump, M5P	Training (95.1074); classification (100)
MultiboostAB	Jrip, REPTree, BFTree, J48, LADTree, NBTree, RandomForest	Training (99.9148); classification (100)
RandomSubSpace	IBk, JRip, REPTree, IB1, NNge, Ridor, FT, LADTree, LMT, NBTree, RandomForest,	Training (99.9148); classification (100)
Dagging	BFTree, LADTree, NBTree, RandomForest, RandomTree	Training (99.5568); classification (100)
Decorate	JRip, BFTree, NBTree	Training (99.8636); classification (100)
MultiClassifier	SImpleCart, JRip, REPTree, ADTree, BFTree, LADTree, RandomForest,	Training (99.8977); classification (100)
CVPParameterSelection	SImpleCart, JRip, REPTree, BFTree, LADTree, NBTree, RandomTree	Training (99.8807); classification (100)
AttributeSelectedClassifier	IBk, SImpleCart, IB1, NNge, BFTree, FT, LADTree, LMT	Training (99.8125); classification (100)
Vote	SImpleCart, JRip, BFTree, LADTree, NBTree, RandomForest	Training (99.8636); classification (100)

The training percentage is calculated from k -folder strategy using Data set 1, while the classification percentage corresponds to the classification of anomalous situations of scans, Data set 2, with a classifier ensemble that has previously been trained with all values of the training data set.

From Tables 2–5, it can be concluded that classification models are able to carry out the classification of the data set and get the right classification for the whole classes. Besides, proving that the use of classifier–ensembles to identify the attacks improves the use of the baseline classifiers in one value between 5% and 50% according to the different classifiers and the class studied. Therefore, these models can classify different conditions of the attack network that are associated with the scans labelled as Classes 3 and 4. The characteristics and options of the chosen ensembles together with their tuned values are shown in Table 6.

TABLE 6. Selected options for ensembles in the experimental study

Ensembles	Selected options
FilteredClassifier	Name of the filter ‘Discretize’.
Adaboost	Number of boost iterations (10), seed for resampling (1), use resampling instead of reweighting (false), percentage of weight mass (100).
MultiboostAB	Number of boost iterations (10), number of subcommittees (3), seed for resampling (1), use resampling instead of re-weighting (false), percentage of weight mass (100).
RandomSubSpace	Number of iterations (10), size of each sub-Space (0.5), seed for resampling (1).
RotationForest	Maximum size of a group (3), minimum size of a group (3), number of iterations to be performed (10), number of groups (false), filter used (Principal Components), percentage of instances to be removed (50), seed for resampling (1).
Attribute selected classifier	Name of an attribute evaluator (CfsSubsetEval), name of a search method (BestFirst).
Bagging	Size of each bag (100), compute out of bag error (false), number of bagging iterations (10), seed for resampling (1).
Classification via regression LogitBoost	Threshold (-1.79), number of folds for internal cross-validation (0), number of iterations (10), number of runs for internal cross-validation (1), seed (1), shrinkage (1.0), resampling is used (false), weight pruning (100).
Grading	Debug (false), metaClassifier (ZeroR), numFolds (10), seed (1).
MultiScheme	Debug (false), numFolds (0), seed (1).
OrdinalClassClassifier	Debug (false).
RandomCommittee	Classifier (RandomTree), debug (false), numIterations (10), seed (1).
StackingC	Debug (false), metaClassifier (LinearRegression), numFolds (10), seed (1).

continued

TABLE 6. Continued

Ensembles	Selected options
RacedIncrementalLogitBoost	Debug (false), maxChunkSize (2000), minChunkSize (500), pruningType (Log likelihood pruning), seed (1), useResampling (false), validationChunkSize (1000).
Dagging	Debug (false), numFolds (10), seed (1), verbose (false).
Decorate	ArtificialSize (10), debug (false), desiredSize (10), numIterations (10), seed (1).
MultiClassClassifier	Debug (false), method (1-against-all), randomWidthFactor (2.0), seed (1), usePairwiseCoupling (false).
CVParameterSelection	CVParameters (0 java.lang.String), debug (false), numFolds (10), seed (1).
Vote	CombinationRule (Average of Probabilities), debug (false), seed (1).
ThresholdSelector	Debug (false), designatedClass ('yes', 'positive', '1'), evaluationMode (singletuned fold), manualThresholdValue (-1.0), measure (FMEASURE), numXValFolds (3), rangeCorrection (norange correction), seed (1).

5 Conclusions and future work

This study has applied a mutation testing model for ensembles of supervised classifiers performing ID on numerical traffic data sets. It aims to assess the generalization capability of the applied classifier ensembles when confronting potential zero-day attacks.

Experimental results show that the applied classifier ensembles deal properly with the analyzed data, containing mutated network scans. It can then be concluded that the applied models are able to properly detect new attacks related to previously unseen scans.

Future work will be based on the mutation of some other attack situations and the broadening of considered classifiers and ensembles.

Acknowledgements

This research has been partially supported through the projects of the Spanish Ministry of Science and Innovation (TIN2010-21272-C02-01 and CIT-020000-2009-12) (both funded by the European Regional Development Fund). The authors would also like to thank the vehicle interior manufacturer, Grupo Antolin Ingenieria S.A., within the framework of the MAGNO2008 - 1028.- CENIT. Project also funded by the MICINN, the Spanish Ministry of Science and Innovation (PID 560300-2009-11) and the Regional Government of Castile-Leon (CCTT/10/BU/0002). This work was also supported in the framework of the IT4Innovations Centre of Excellence project, reg. no. (CZ.1.05/1.1.00/02.0070) supported by the Operational Program 'Research and Development for Innovations' funded through the Structural Funds of the European Union and the state budget of the Czech Republic.

References

- [1] K. Abdullah, C. Lee, G. Conti, and J. A. Copeland. Visualizing network data for intrusion detection. In *Sixth Annual IEEE Information Assurance Workshop*, pp. 100–108, West Point, NY, 2005. IEEE Systems, Man and Cybernetics (SMC).
- [2] J. Allen, A. Christie, W. Fithen, J. Mchugh, J. Pickel, and E. Stoner. *State of the practice of intrusion detection technologies. Technical Report CMU/SEI-99-TR-028, ESC-99-028*, Carnegie Mellon University, Pittsburgh, PA 15213-3890, January 2000.
- [3] J. P. Anderson. Computer security threat monitoring and surveillance. *Technical report*, James P. Anderson Co., Fort Washington, Pennsylvania, 1980.
- [4] N. Athanasiades, R. Abler, J. Levine, H. Owen, and G. Riley. Intrusion detection testing and benchmarking methodologies. In *Proceedings of the First IEEE International Workshop on Information Assurance (IWIA'03), IEEE-IWIA '03*, pp. 63–72, Washington, DC, USA, 2003. IEEE Computer Society, 2003.
- [5] B. Bacca, J. Salvi, and X. Cufi. Appearance-based mapping and localization for mobile robots using a feature stability histogram. *Robotics and Autonomous Systems*, **59**, 840–857, 2011.
- [6] I. Balepin, S. Maltsev, J. Rowe, and K. N. Levitt. Using specification-based intrusion detection for automated response. In *RAID*, pp. 136–154, 2003.
- [7] P. Baraldi, G. Gola, E. Zio, D. Roverso, and M. Hoffmann. A randomized model ensemble approach for reconstructing signals from faulty sensors. *Expert Systems with Applications*, **38**, 9211–9224, 2011.
- [8] B. Baruque, E. Corchado, A. Mata, and J. M. Corchado. A forecasting solution to the oil spill problem based on a hybrid intelligent system. *Information Sciences*, **180**, 2029–2043, 2010.
- [9] M. Bermúdez-Edo, R. Salazar-Hernández, J. Díaz-Verdejo, and P. García-Teodoro. Proposals on assessment environments for anomaly-based network intrusion detection systems. In *Critical Information Infrastructures Security*, J. Lopez, ed., Vol. 4347 of *Lecture Notes in Computer Science*, pp. 210–221. Springer, 2006.
- [10] M. Berry and G. Linoff. *Data Mining Techniques*. Wiley Computer Publishing, 2004.
- [11] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2007.
- [12] R. Blagojevic, B. Plimmer, J. Grundy, and Y. Wang. Using data mining for digital ink recognition: dividing text and shapes in sketched diagrams. *Computers & Graphics-UK*, **35**, 976–991, 2011.
- [13] J. B. D. Cabrera, C. Gutiérrez, and R. K. Mehra. Ensemble methods for anomaly detection and distributed intrusion detection in mobile ad-hoc networks. *Information Fusion*, **9**, 96–119, 2008.
- [14] J. D. Case, M. Fedor, M. L. Schoffstall, and J. Davin. *Simple Network Management Protocol (SNMP)*, Technical report, RFC Editor, United States, 1990.
- [15] J. Case, K. McCloghrie, M. Rose, and S. Waldbusser. *Introduction to version 2 of the internet-standard network management framework*, Technical report, RFC Editor, United States, 1993.
- [16] Secure Consulting. *Vulnerability statistics report*. Cisco Secure Consulting, 2000.
- [17] E. Corchado, B. Baruque, and H. Yin. Boosting unsupervised competitive learning ensembles. In *ICANN*, pp. 339–348, 2007.
- [18] E. Corchado and Á. Herrero. Neural visualization of network traffic data for intrusion detection. *Applied Soft Computing*, **11**, 2042–2056, 2011.
- [19] E. Corchado, Á. Herrero, and J. M. Sáiz. Detecting compounded anomalous snmp situations using cooperative unsupervised pattern recognition. In *15th International Conference on Artificial Neural Networks (ICANN 2005)*, W. Duch, J. Kacprzyk, E. Oja, and S. Zadrozny, eds, Vol. 3697 of *LNCS*, pp. 905–910. Springer, 2005.

- [20] E. Corchado, Á. Herrero, and J. M. Sáiz. Testing CAB-IDS through mutations: on the identification of network scans. In *International Conference in Knowledge-Based and Intelligent Engineering; Information Systems (KES 2006)*, R. J. Howlett, L. C. Jain and B. Gabrys, eds, Vol. 4252 of *LNAI*, pp. 433–441. Springer, 2006.
- [21] DARPA. *Darpa intrusion detection evaluation*, 1999.
- [22] J. Davin, J. Galvin, and K. McCloghrie. *SNMP Administrative Model. RFC 1351 (Historic)*, July 1992.
- [23] H. Debar, M. Dacier, A. Wespi, and S. Lampart. An experimentation workbench for intrusion detection systems. *Research Report RZ 2998*. IBM Research Division Zurich Research Laboratory, 1998.
- [24] H. Debar and B. Morin. Evaluation of the diagnostic capabilities of commercial intrusion detection systems. In *Recent Advances in Intrusion Detection*. A. Wespi, G. Vigna, and L. Deri, eds, Vol. 2516, of *Lecture Notes in Computer Science*, pp. 177–198. Springer, 2002.
- [25] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *The Annals of Statistics*, **28**, 337–407, 2000.
- [26] P. Garcia-Teodoro, J. E. Dáz-Verdejo, G. Maciá-Fernández, and E. Vázquez. Anomaly-based network intrusion detection: techniques, systems and challenges. *Computers & Security*, **28**, 18–28, 2009.
- [27] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The WEKA data mining software: an update. *ACM SIGKDD Explorations Newsletter*, **11**, 10–18, 2009.
- [28] N. Hatami. Thinned ECOC decomposition for gene expression based cancer classification. In *Eighth International Conference on Intelligent Systems Design and Applications ISDA 2008*, J. S. Pan, A. Abraham, and C. C. Chang, eds, Vol. 1, pp. 216–221. *IEEE; IEEE SMC Soc; Natl Sci Council; Minist Educ, IEEE COMPUTER SOC. 8th International Conference on Intelligent Systems Design and Applications (ISDA 2008)*, Kaohsiung, Taiwan, November 26–28, 2008.
- [29] Á. Herrero, E. Corchado, P. Gastaldo, and R. Zunino. Neural projection techniques for the visual inspection of network traffic. *Neurocomputing*, **72**, 3649–3658, 2009.
- [30] V. Hu, P. Mell and R. Lippmann. An overview of issues in testing intrusion detection systems. *NIST Interagency Reports*. National Institute of Standards and Technology - Information Technology Laboratory, 2003.
- [31] SANS Institute. *The top 10 most critical internet security threats (2000-2001 archive)*. SANS Institute, 2001.
- [32] KDD. *Kdd cup 1999 dataset*, 1999.
- [33] R. Lippmann, J. W. Haines, D. J. Fried, J. Korba, and K. Das. The 1999 DARPA off-line intrusion detection evaluation. *Computer Networks*, **34**, 579–595, 2000.
- [34] R. P. Lippmann, J. W. Haines, L. M. Rossey and R. Cunningham. Extending the darpa off-line intrusion detection evaluations. *DARPA Information Survivability Conference and Exposition II*, pp. 35, Los Alamitos, CA, USA, 2001. IEEE Computer Society.
- [35] M. Malowidzki. Getbulk worth fixing. *The Simple Times*, **10**, 3–6, 2002.
- [36] J. McHugh. Testing intrusion detection systems: a critique of the 1998 and 1999 DARPA intrusion detection system evaluations as performed by Lincoln laboratory. *ACM Transactions on Information and System Security*, **3**, 262–294, 2000.
- [37] P. Mueller and G. Shipley. Cover story: dragon claws its way to the top. *Network Computing*, **12**, 45–67, 2001.

- [38] A. Munson and R. Caruana. Cluster ensembles for network anomaly detection. *Technical Report 2006-2047*, Cornell University, 2006.
- [39] S. Nawata, M. Uchida, Y. Gu, M. Tsuru, and Y. Oie. Unsupervised ensemble anomaly detection through time-periodical packet sampling. In *INFOCOM IEEE Conference on Computer Communications Workshops, 2010*, 2010.
- [40] S. Northcutt, M. Cooper, K. Fredericks, and M. Fearnow. *Intrusion Signatures and Analysis*. New Riders Publishing, Thousand Oaks, 2001.
- [41] R. Polikar. Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine*, **6**, 21–45, 2006.
- [42] N. Puketza, M. Chung, R. A. Olsson, and B. Mukherjee. A software platform for testing intrusion detection systems. *IEEE Software*, **14**, 43–51, 1997.
- [43] N. J. Puketza, K. Zhang, M. Chung, B. Mukherjee, and R. A. Olsson. A methodology for testing intrusion detection systems. *IEEE Transactions Software Engineering*, **22**, 719–729, 1996.
- [44] M. J. Ranum. *Experiences benchmarking intrusion detection systems. Technical report, NFR Security Technical Publications*, 2001.
- [45] P. Ren, Y. Gao, Z. Li, Y. Chen, and B. Watson. IDGraphs: intrusion detection and analysis using stream compositing. *IEEE Computer Graphics and Application*, **26**, 28–39, 2006.
- [46] J. M. Rizza. *Computer Network Security*. Springer, 2005.
- [47] D. Ruta and B. Gabrys. An overview of classifier fusion methods. *Computing and Information Systems*, **7**, 1–10, 2000.
- [48] M. Sabhnani and G. Serpen. Why machine learning algorithms fail in misuse detection on KDD intrusion detection data set. *Intelligent Data Analysis*, **8**, 403–415, 2004.
- [49] C. Schafer, P. Laskov, P. Dussel, and K. Rieck. Learning intrusion detection: supervised or unsupervised? In *ICIAP, 13th International Conference on Image Analysis and Processing*, F. Roli and S. Vitulano, eds, pp. 50–57. Springer, 2005.
- [50] B. Schuller, B. J. B. Schmitt, D. Arsic, S. Reiter, M. Lang, and G. Rigoll. Feature selection and stacking for robust discrimination of speech, monophonic singing, and polyphonic music. In *2005 IEEE International Conference on Multimedia and Expo (ICME), Vols 1 and 2*, pp. 840–843, 345. IEEE, IEEE. *IEEE International Conference on Multimedia and Expo (ICME)*, July 6–8, 2005.
- [51] J. Sedano, A. Berzosa, J. R. Villar, E. Corchado, and E. de la Cal. Optimising operational costs using soft computing techniques. *Integrated Computer-aided Engineering*, **18**, 313–325, 2011.
- [52] M. Sewell. *Ensemble Learning. Technical Report*, University College London - Department of Computer Science, 2011.
- [53] A. J. C. Sharkey and N. E. Sharkey. Diversity, selection and ensembles of artificial neural nets. In *Proceedings of Third International Conference on Neural Networks and their Applications. IUSPIM*, pp. 205–212, University of Aix-Marseille III, Marseilles, France, 1997.
- [54] A. J. C. Sharkey and N. E. Sharkey. Combining diverse neural nets. *Knowledge Engineering Review*, **12**, 231–247, 1997.
- [55] R. Sprenkels and J. P. Martin-Flatin. *Bulk Transfers of MIB Data. Communication Systems Division*. Swiss Federal Institute of Technology Lausanne, 1999.
- [56] S. Staniford, J. A. Hoagland, and J. M. McAlerney. Practical automated detection of stealthy portscans. *Journal of Computer Security*, **10**, 105–136, 2002.
- [57] A. Sundaram. An introduction to intrusion detection. *Crossroads*, **2**, 3–7, 1996.
- [58] C.-F. Tsai, Y.-F. Hsu, C.-Y. Lin, and W.-Y. Lin. Intrusion detection by machine learning: a review. *Expert Systems with Applications*, **36**, 11994–12000, 2009.

- [59] G. Wang and J. Ma. Study of corporate credit risk prediction based on integrating boosting and random subspace. *Expert Systems with Applications*, **38**, 13871–13878, 2011.
- [60] D. Wolpert. Stacked generalization. *Neural Networks*, **5**, 241–259, 1992.

Received 29 April 2010