

Intrusion Detection System with Data Mining Approach A Review

DR. Madjid Khalilian¹ and Norwati Mustapha²

¹ university putra malaysia

Received: 7 February 2011 Accepted: 27 February 2011 Published: 9 March 2011

Abstract

Despite of growing information technology widely, security has remained one challenging area for computers and networks. Recently many researchers have focused on intrusion detection system based on data mining techniques as an efficient strategy. The main problem in intrusion detection system is accuracy to detect new attacks therefore unsupervised methods should be applied. On the other hand, intrusion in system must be recognized in realtime, although, intrusion detection system is also helpful in off-line status for removing weaknesses of network's security. However, data mining techniques can lead us to discover hidden information from network's log data. In this survey, we try to clarify: first, the different problem definitions with regard to network intrusion detection generally; second, the specific difficulties encountered in this field of research; third, the varying assumptions, heuristics, and intuitions forming the basis of erent approaches; and how several prominent solutions tackle different problems.

Index terms— Data mining, Intrusion Detection, Clustering, Classification.

1 INTRODUCTION

owadays we have many applications with massive amount of data which causes limitation in data storage capacity and processing time. Furthermore, many applications must operate in realtime to achieve theirs objectives. As an important case for these kinds of application, Network Intrusion Detection System (NIDS) can be pointed. Generally we define NIDS as the detection of intrusions or intrusions attempts either manually or via software expert systems that operate on logs or other information available from the system or the network. An intrusion is a deliberate, unauthorized attempt to access or manipulate information or system and to render them unreliable or unusable. If a suspicious activity is from your internal network or system it will also be classified as intrusion. Some popular intrusion as follows:

? Denial of service (DoS): attempts to starve a host of resources needed to function correctly. ? Scan: reconnaissance on the network or a particular host. ? Worms and viruses: replicating on other hosts.

? Compromises: obtain privileged access to a host by known vulnerabilities.

About-Computer science and information technology University Putra Malaysia E-mails-khalilian@ieee.org,{norwati, nasir, ali}@fsktm.upm.edu.my Furthermore, we identify some important objectives for IDS as below:

? Detect wide variety of attacks.

? Detect intrusions in timely fashion.

? Present analysis in simple, easy-to-understand format. ? Minimize false positives, false negatives:

1. False positive: An event, incorrectly identified by the IDS as being an intrusion when none has occurred

2. False negative: An event that the IDS fails to identify as an intrusion when one has in fact occurred There are many solutions for intrusion detection that we categorize them into four main groups: anomaly detection, signature based misuse, host based and network based. Many researchers have applied data mining techniques,

44 which are powerful methods for extracting hidden information from huge datasets, for network intrusion detection
45 system. On the other part, traditional data mining is not suitable for this kind of applications so they should be
46 tuned and changed or designed with new algorithms. Besides of speed up and storage capacity, real-life concepts
47 tend to change over time e.g. new attacks should be recognized.

48 The growth of volume of existing data and insufficiency of data storage capacity lead us to the dynamic
49 processing data and extracting knowledge. The problem is that current IDS are tuned specifically to detect
50 known service level network attacks. At the same time, enough data exists or could be collected to allow network
51 administrators to detect these policy violations. Unfortunately, the data is so enormous, and the analysis process
52 so time-consuming, that the administrators don't have the resources to proactively analyze the data for policy
53 violations, especially in the presence of a high number of false positives that cause them to waste their limited
54 resources.

55 The nature solution is utilizing data mining techniques. However, data mining can be applied in offline status.
56 Most previous work focused on off-line environment while on-line system for detecting policy violations is needed.

57 In next section we addressed general problems in this domain, after that we discuss different solutions in four
58 groups with pros and cons, finally we will have the conclusion.

59 2 II.

60 3 Gaps

61 4 traffic based Methods

62 If we want to categorize intrusion detection methods, we will recognize two main aspects for grouping approaches,
63 which one group refers to type of attack includes host based and network based. Another group of approaches
64 refers to solutions techniques which are signature based and anomaly detection methods. In continue we review
65 these techniques with their pros and cons.

66 5 a) Host based methods

67 This method is based on data source category; consequently, its data comes from the records of various activities
68 of hosts, including system logs, audit operation system information, etc. the main architecture for this kind
69 of methods is similar to network based which is described in the next section. Ref [1] presents a host-based
70 combinatorial method based on k-Means clustering and ID3 decision tree learning algorithms for unsupervised
71 classification of anomalous and normal activities in computer network.

72 6 IV. approaches to solutions

73 A powerful survey can be found in [2] that it discusses data mining for cyber security applications. For example,
74 anomaly detection techniques could be used to detect unusual patterns and behaviors, Link analysis may be used
75 to trace the viruses to the perpetrators, Classification may be used to group various cyber attacks and then use
76 the profiles to detect an attack when it occurs, Prediction may be used to determine potential future attacks
77 depending in a way on information learnt about terrorists through email and phone conversations. This paper
78 also mentioned about real-time problem in IDS and other challenges include mining unstructured data types.

79 We divide approaches in two main groups: misuse detection which the main study is the classification
80 algorithms and anomaly detection which the main study is the pattern comparison(association rules and sequence
81 rules) and the cluster algorithms.

82 7 a) Signature-based methods

83 The research [3] compares accuracy, detection rate, false alarm rate and accuracy of other attacks under different
84 proportion of normal information. For comparison results of C4.5 and SVM, they demonstrate that C4.5 is
85 superior to SVM in accuracy and detection; in accuracy for Probe, Dos and U2R attacks, C4.5 is also better than
86 SVM; but in false alarm rate, SVM is better. Through test and comparison, the accuracy and detection rate
87 of C4.5 is higher than that of SVM, but false alarm rate of SVM is better. In sampling, the research supposes
88 that the distribution of attack data other than normal data is even, which cannot surely get optimal results, and
89 this should be improved and validated. Another weakness refers to C4.5 parameters that is not optimal, thus
90 the future work should optimize the parameters according to C4.5 parameters and different training dataset.
91 For huge datasets optimizing parameter in SVM takes too much time; however, it is not suitable, for intrusion
92 detection system requires realtimeliness. The future research should aim at the direction where the parameters
93 can be optimized rapidly.

94 With the concept of field in physics, [4] proposed a data field based method for discrimination of network
95 behaviors. Similar to electric charge or particle, each data point we concerned has its own influence region and
96 the influence is a function of position giving the force on each point placed at that position. Furthermore, the
97 positive potential and negative potential have been described, by which it can determine the test point's class.
98 This scheme is based on "supervised" learning, whereas unsupervised methods are preferred. Some advantages
99 and disadvantages are as follows:

100 ? Advantages o Specifying exact class of attacks. o Efficiency is high and complexity is low.
101 ? disadvantages o Many false positives: prone to generating alerts when there is no problem in fact. o Cannot
102 detect unknown intrusions.

103 8 b) Anomaly based methods

104 The basic idea of clustering analysis originates in the difference between intrusion and normal pattern;
105 consequently, we can put data sets into different categories and detect intrusion by distinguish normal and
106 abnormal behaviors. The common clustering algorithms in data mining include two main categories: hierarchical
107 and partitioning clustering algorithms. Clustering intrusion detection is detection for anomaly with no
108 supervision, and it detects intrusion by training the unmarked data.

109 Ref [5] considered the outlier factor of clusters for measuring the deviation degree of a cluster. A method
110 has been proposed to compute the cluster radius threshold. The data classification has been performed by an
111 improved nearest neighbor (INN) method. For the unsupervised intrusion detection, they applied a clustering
112 based method that its time complexity is linear with the size of dataset and the number of attributes.

113 Ref [6] outline a data mining framework for constructing intrusion detection models. To facilitate adaptability
114 and extensibility, they use of meta-learning as a means to construct a combined model that incorporate evidence
115 from multiple base models. They also extend the basic association rules and frequent episodes algorithms to
116 accommodate the special requirements in analyzing audit data. The main shortcoming is lack of devising
117 a mechanical procedure to translate automatically learned detection rules into modules for real-time IDS. [7]
118 proposed a new weighted support vector clustering algorithm; it can cluster large data set and high-dimensional
119 data effectively. They also introduced the new weighted SVC method to network intrusion detection. The
120 experiments with KDD Cup1999 data demonstrate that proposed method achieves highly detection rate with
121 low false alarm rate.

122 Ref [8] have presented a fast distributed outlier detection algorithm for mixed attribute datasets that deals
123 with sparse high-dimensional data. The algorithm called outlier detection for mixed attribute datasets identifies
124 outliers based on the categorical attributes first, and then focuses on subsets of data in the continuous space by
125 utilizing information about these subsets from the categorical attribute space.

126 Ref [9] improved the speed of intrusion detection system, keep the high detect date and the low false positive
127 rate using the Parallel Clustering Ensemble based on Evidence Accumulation algorithm, it overcomes the
128 disadvantages of conventional Parallel K-means algorithm. Through paralleling, the algorithm clusters more
129 speedily facing to mass data, and keep the advantages of the Evidence Accumulation which combines the results
130 of multiple clustering into a single data partition.

131 Ref [10] present a method for outlier detection that uses HPSO clustering based on swarm intelligence, which
132 is capable of providing clustering at different levels of compactness. Merging clusters and attribute evolution help
133 in learning about the correct cluster solution and outlier data. Experiments show that the approach is capable
134 of identifying true outliers as well as a good clustering configuration of data. Setting parameters automatically
135 is a challenging problem in this method. High dimension data is also problematic in this research.

136 Ref [11] is an anomaly detection algorithm based on hierarchical clustering, called ADBHC. ADBHC generates
137 clusters using density-based partitioning method which has less computational cost. It uses the improved
138 hierarchical clustering tree to carry out fast scalable and adaptive anomaly detection. The improved hierarchical
139 clustering tree supports updating profiles at any time. They extend the clustering algorithm and apply branch
140 and bound mechanism for filtering noise. ADBHC had lower false alarm rate and higher detection rate. The
141 superior performance of detection was mainly due to the high accuracy of normality profiles and the capability
142 of filtering noise. Various parameters had pernicious impacts on the adaptive captivity of ADBHC.

143 ? Advantages: Anomaly detection can detect novel attacks to increase the detection rate. Compared
144 to supervised approaches, unsupervised approach breaks the dependency on attack-free training datasets.
145 The performance of unsupervised anomaly detection approaches achieve higher detection rate over supervised
146 approach. Also, unsupervised approach have high false positive rate over supervised approach. Using
147 unsupervised anomaly detection techniques, however, the system can be trained with unlabeled data and is
148 capable of detecting previously unseen attacks [12].

149 ? Disadvantages: Obviously, not all typical behaviors are attacks or intrusion attempts. This represents one
150 drawback of intrusion detection methods based on clustering [13].

151 9 c) Hybrid methods

152 Through analyzing the advantages and disadvantages between anomaly detection and misuse detection, a mixed
153 intrusion detection system (IDS) model is designed. [14]First, data is examined by the misuse detection module,
154 then abnormal data detection is examined by anomaly detection module.

155 Ref [1]proposed combinatorial approach for unsupervised classification of anomalous and normal activities in
156 computer network. The proposed approach combines the two well-known machine learning methods: the k-Means
157 clustering and the ID3 decision tree learning approaches. The k-Means method was first applied to partition
158 the training instances into k disjoint clusters. The ID3 decision tree built on each cluster learns the subgroups

159 within the cluster and partitions the decision space into finer classification regions; thereby improving the overall
160 classification performance.

161 Ref [15] An incremental intrusion detecting model is proposed. This model integrates unsupervised Self
162 Organizing Map and supervised Radial Basis Function to complete incremental learning. Self Organizing Map
163 can get new type intrusion information and generate new nodes in Radial Basis Function. By this model, intrusion
164 of unknown type can be detected online.

165 Fuzzy clustering algorithm is an unsupervised anomaly detection technique without training; it does not need
166 to know the type of attack in Intrusion Detection data samples, so it can detect a variety of known and unknown
167 characteristics of network intrusion simultaneously. This article combined QPSO with the FCM algorithm, using
168 QPSO algorithm has better features to find the global optimal value, using particle swarm flying in the solution
169 space search best value to replace FCM iterative process to obtain a more suitable mix of clustering algorithm
170 [16].

171 In order to reduce or eliminate the noise impact on constructing the hyper plane of SVM, firstly it preprocesses
172 the data, after that the fuzzy membership function is introduced into SVM. The fuzzy membership function
173 acquires different values for each input data according to different effects on the classification result. Because
174 different network protocol has different attributes, that must affect the detection effect. This paper proposes
175 cooperative network intrusion detection Based on Fuzzy SVM. Three types of detecting agents are generated
176 according to TCP, UDP and ICMP protocol. How to improve the accuracy of UDP detection agent in existing
177 data set will be the major weakness [17].

178 V.

179 10 Conclusions

180 In this paper we have demonstrated some difficulties in Network Intrusion Detection Systems where its log files
181 are high scale and dimensions; consequently, new methods need to be developed for processing these huge data
182 sources. Furthermore concept drift is nature of data in IDS and should be managed by new methods. On the
183 other hand, efficiency in terms of accuracy is one of the most critical measurements which are mostly defined
184 by ratio of false positive and false negative alarms. Therefore, we need to design efficient algorithms whereas
185 scan data once and extract hidden patterns inside it. Evolving data, visiting data once, accuracy in intrusion
186 detections and space limitations are major issues in intrusion detection systems. However, there are two main
187 approaches for intrusion detection: first group employs signature-based methods to identify attacks and second one
188 refers to anomaly detection techniques but devising new framework with combining these two main approaches
189 can overcome most drawbacks.

190 VI. ¹

¹Intrusion Detection System with Data Mining Approach A Review ©2011 Global Journals Inc. (US)

.1 Acknowledgement

- 191 This work was supported by grant 03-04-10-875FR from the Basic Research Program of the University Putra
192 Malaysia.
193
- 194 [Jiang et al. ()] ‘A clustering-based method for unsupervised intrusion detections’. S Y Jiang , X Song , H Wang
195 , J J Han , Q H Li . *Pattern Recognition Letters* 2006. 27 p. .
- 196 [Teng et al. ()] ‘A Cooperative Network Intrusion detection Based on Fuzzy SVMs’. S Teng , H Du , N Wu , W
197 Zhang , J Su . *Journal of Networks* 2010. 5 p. 475.
- 198 [Lee et al. ()] *A data mining framework for building intrusion detection models*, W Lee , S J Stolfo , K W Mok
199 . 2002.
- 200 [Koufakou and Georgiopoulos (2011)] *A fast outlier detection strategy for distributed high-Global Journal of*
201 *Computer Science and Technology Volume XI Issue V Version*, M Koufakou , Georgiopoulos . April 2011.
- 202 [Zhang et al. ()] *A Mixed Unsupervised Clustering-Based Intrusion Detection Model*, G Zhang , S Zhang , Sun .
203 2009.
- 204 [Yasami and Mozaffari ()] ‘A novel unsupervised classification approach for network anomaly detection by k-
205 Means clustering and ID3 decision tree learning methods’. Y Yasami , S P Mozaffari . *The Journal of*
206 *Supercomputing* 2010. 53 p. .
- 207 [Gao et al. ()] *A Parallel Clustering Ensemble Algorithm for Intrusion Detection System*, H Gao , D Zhu , X
208 Wang . 2010.
- 209 [Alam et al. ()] *A swarm intelligence based clustering approach for outlier detection*, S Alam , G Dobbie , P
210 Riddle , M A Naeem . 2010.
- 211 [Sun and Wang ()] *A Weighted Support Vector Clustering Algorithm and its Application in Network Intrusion*
212 *Detection*, S Sun , Y Z Wang . 2009.
- 213 [Liang et al. ()] *An Adaptive Anomaly Detection Based on Hierarchical Clustering*, H Liang , R Wei-Wu , R Fei
214 . 2009.
- 215 [Gogoi et al. ()] ‘Anomaly Detection Analysis of Intrusion Data using Supervised & Unsupervised Approach’. P
216 Gogoi , B Borah , D K Bhattacharyya . *Journal of Convergence Information Technology* 2010. 5.
- 217 [Thuraisingham et al. ()] *Data mining for security applications*, B Thuraisingham , L Khan , M M Masud , K
218 W Hamlen . 2009.
- 219 [Wu and Yen ()] ‘Data mining-based intrusion detectors’. S Y Wu , E Yen . *Expert Systems with Applications*
220 2009. 36 p. .
- 221 [Tian and Liu ()] *Incremental intrusion detecting method based on SOM/RBF*, L Y Tian , W P Liu . 2010.
- 222 [Singh et al. ()] ‘Mining Common Outliers for Intrusion Detection’. G Singh , F Masegla , C Fiot , A Marascu
223 , P Poncelet . *Advances in Knowledge Discovery and Management*, 2010. p. .
- 224 [Wang et al. ()] *Network intrusion detection based on hybrid Fuzzy Cmean clustering*, H Wang , Y Zhang , D Li
225 . 2010.
- 226 [US) dimensional data sets with mixed attributes Data Mining and Knowledge Discovery ()] ‘US) dimensional
227 data sets with mixed attributes’. *Data Mining and Knowledge Discovery* 2010. 20 p. . (©2011 Global Journals
228 Inc.)
- 229 [Xie and Bai ()] *Using Data Field to Analyze Network Intrusions*, F Xie , S Bai . 2006. p. . (Information Security
230 Practice and Experience)