



HAL
open science

Data Mining for Intrusion Detection: from Outliers to True Intrusions

Goverdhan Singh, Florent Masegla, Céline Fiot, Alice Marascu, Pascal Poncelet

► **To cite this version:**

Goverdhan Singh, Florent Masegla, Céline Fiot, Alice Marascu, Pascal Poncelet. Data Mining for Intrusion Detection: from Outliers to True Intrusions. The 13th Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD-09), Apr 2009, BANGKOK, Thailand. pp.891-898. inria-00359206v1

HAL Id: inria-00359206

<https://inria.hal.science/inria-00359206v1>

Submitted on 6 Feb 2009 (v1), last revised 28 Oct 2009 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Data Mining for Intrusion Detection: from Outliers to True Intrusions

Goverdhan Singh¹, Florent Massegli¹, Céline Fiot¹, Alice Marascu¹, and
Pascal Poncelet²

¹ INRIA Sophia Antipolis, 2004 route des lucioles - BP 93, FR-06902 Sophia Antipolis
Email: First.Last@sophia.inria.fr

² LIRMM UMR CNRS 5506, 161 Rue Ada, 34392 Montpellier Cedex 5, France
Email: poncelet@lirmm.fr

Abstract. Data mining for intrusion detection can be divided into several sub-topics, among which unsupervised clustering has controversial properties. Unsupervised clustering for intrusion detection aims to i) group behaviors together depending on their similarity and ii) detect groups containing only one (or very few) behaviour. Such isolated behaviours are then considered as deviating from a model of normality and are therefore considered as malicious. Obviously, all atypical behaviours are not attacks or intrusion attempts. Hence, this is the limits of unsupervised clustering for intrusion detection. In this paper, we consider to add a new feature to such isolated behaviours before they can be considered as malicious. This feature is based on their possible repetition from one information system to another. We propose a new outlier mining principle and validate it through a set of experiments.

1 Introduction

Intrusion detection is a very important topic of network security that has received much attention [13, 20, 12, 17] since potential cyber threats are making the organizations vulnerable. *Intrusion Detection Systems (IDS)* are intended to protect information systems against intrusions and attacks and are traditionally based on signatures of known attacks [19, 2]. Therefore, new kinds of attacks regularly have to be added to the signature list. The main drawback is that in case of an emerging attack, based on the recent discovery of a new security hole for instance, the IDS will ignore it since this new attack has not yet been listed in the base of signatures.

Protecting a system against new attacks, while keeping an automatic and adaptive framework is an important topic in this domain. One answer to that problem could rely on data mining. Data mining tools have been used to provide IDS with more adaptive detection of cyber threats [4, 21]. Among those data mining approaches, anomaly detection tries to deduce intrusions from atypical records [12, 6]. The overall principle is generally to build clusters, or classes, of usage and find outliers (*i.e.* events that do not belong to any class or group identifying normal usage). Actually, outlier detection aims to find records that deviate significantly from a well-defined notion of normality. It has a wide range of applications, such as fraud detection for credit card [1], health care, cyber security [4] or safety of critical systems [8]. However, the main drawback of detecting intrusions by means of anomaly (outliers) detection is the high rate of false alarms since an alarm can be triggered because of a new kind of usages that has never been

seen before (and is thus considered as abnormal). Considering the large amount of new usage patterns emerging in the Information Systems, even a weak percent of false positive will give a very large amount of spurious alarms that would be overwhelming for the analyst. Therefore, the goal of this paper is to propose an intrusion detection algorithm that is based on the analysis of usage data coming from multiple partners in order to reduce the number of false alarms. Our main idea is that a new usage is likely to be related to the context of the information system on which it occurs (so it should only occur on this system). On the other hand, when a new security hole has been found on a system, the hackers will want to use it in as many information systems as possible. Thus a new anomaly that occurs on two (or more) information systems is probably not a new kind of usage, but rather an intrusion attempt. Let us consider A_x , an anomaly detected in the usage of web site S_1 corresponding to a php request on the staff directory for a new employee: John Doe, who works in room 204, floor 2, in the R&D department. The request will have the following form: `staff.php?FName=John\&LName=Doe\&room=204\&floor=2\&Dpt=RD`. This new request, due to the recent recruitment of John Due in this department, should not be considered as an attack. On the other hand, let us consider A_y , an anomaly that corresponds to a true intrusion. A_y will be based on a security hole of the system (for instance a php vulnerability) and might, for instance, look like: `staff.php?path=./etc/passwd%00`. One can see in this request that the parameters are not related to the data accessed by the php script, but rather to a security hole that has been discovered on the *staff* script. If two or more firms use the same script (say, a directory requesting script bought to the same software company) then the usage of this security hole will certainly be repeated from one system to another and the request having parameter `./etc/passwd%00` will be the same for all the victims. In this paper, we propose to provide the end-user with a method that takes only one parameter: n , the number of desired alarms. Then, based on the analysis of the usage data coming from the different partners, our algorithm will detect n common outliers they share. Such common outliers are likely to be true attacks and will trigger an alarm. In a real-world application of this technique, privacy preserving will be a major issue in order to protect partners' data. In this paper we focus on clustering and outlier detection techniques in a distributed environment. However, privacy issues in our framework are presented in another paper currently being submitted.

The paper is organized as follows. In Section 2 we present the motivation of this approach and our general framework and Section 3 gives an overview of existing works in this domain. Section 4 presents COD, our method for detecting outliers and triggering true alarms. Eventually, our method is tested through a set of experiments in Section 5 and Section 6 gives the conclusion.

2 Motivation and General Principle

In this section, we present the motivation of our work, based on the main drawbacks of existing anomaly-based methods for intrusion detection and we propose the general principle of COD, our algorithm for comparing the anomalies on different systems.

2.1 Motivation

Anomaly-based IDS [6] can be divided into two categories; semi-supervised and unsupervised. The semi-supervised methods build a model of “normal” behaviours on the system. Every behaviour that is not considered as normal is an anomaly and should trigger an alarm. The unsupervised methods do not use any labelled data. Usually, based on a clustering algorithm, they try to detect outliers and consider them as anomalies. Obviously, anomaly-based IDS will suffer from a very high number of false alarms since a new kind of behaviour will be considered as an anomaly (and an attack).

In this paper, we propose to improve the results of unsupervised IDS by means of a collaborative framework involving different network-based systems. Section 3 gives an overview of the existing IDS based on the principles presented above and the existing collaborative IDS. However, to the best of our knowledge, our proposal is the first unsupervised IDS using the common anomalies of multiple partners in order to detect the true intrusion attempts. The main idea of our proposal is that multiple partners do not share the same data, but they share the same systems (the Web server can be Apache or IIS, the data server can run Oracle, the scripts accessing the data can be written with PHP or CGI, etc). When a security hole has been found for one system (say a php scripts with specific parameters leading to privileged access to the hard drive), then this weakness will be the same for all the partners using the same technology. Our goal is to reduce the rate of false alarms based on this observation, as explained in section 2.2

2.2 General Principle

In this paper we present COD (Common Outlier Detection) a framework and algorithm intended to detect the outliers shared by at least two partners in a collaborative IDS. Outliers are usually small clusters. Some methods used to find them are presented in section 3. Our goal is to use such outlier lists from different systems (based on a similar clustering, involving the same similarity measure). If an outlier occurs for at least two systems, then it is considered as an attack. COD is indeed based on the assumption that an intrusion attempt trying to find a weakness of a script will look similar for all the victims of this attack. We propose to detect intrusion attempts among the records of a Web server, such as an Apache access log file. Such a file keeps record, for each access on the Web site, of the IP, date, requested URL and referrer (among other informations). For clarity of presentation we present our framework on the collaboration of two Web sites, S_1 and S_2 and we consider the requests that have been received by the scripts of each site (cgi, php, sql, etc). Our goal is to perform a clustering on the usage patterns of each site and find the common outliers. However, that would not be enough to meet the second constraint of our objective: to require only one parameter, n , the number of alarms to return. Our similarity measure (presented in section 4.1) will allow normal usage patterns to be grouped together rather than grouped with intrusion patterns. On the other hand, our similarity measure also has to ensure distinguishing an intrusion pattern from normal usage patterns and from other intrusion patterns (since different intrusion patterns will be based on a different security hole and will have very different characteristics). Our algorithm performs successive clustering steps for each site. At each step we check the potentially matching outliers between

both sites. The clustering algorithm is agglomerative and depends on the maximum dissimilarity (MD) that has to be respected between two objects. Let us consider that n , the desired number of alarms, is set to 1 and the usage patterns are distributed as illustrated in figure 1. Let us also consider that for these sites the cluster labeled D at step 1 is the only one that corresponds to an intrusion attempt. For step one, MD is initialized with a very low value, so the clusters will be as tight and small as possible. Then we check correspondences between outliers of S_1 and S_2 . Let us consider the clustering results on S_1 and S_2 at step one in figure 1. There are four matching of outliers between both sites (A , B , C and D). That would lead to 4 alarms (among which only one is true) which is more than desired by the user and . We thus have to increase the clustering tolerance (*i.e.* increase MD) so that bigger clusters can be built. After a few steps, we will find the clusters of step n in figure 1. The only common outlier is A , which corresponds to the intrusion attempt. Furthermore, this will trigger one alarm, as desired by the user, and there is no need to continue increasing MD until step m .

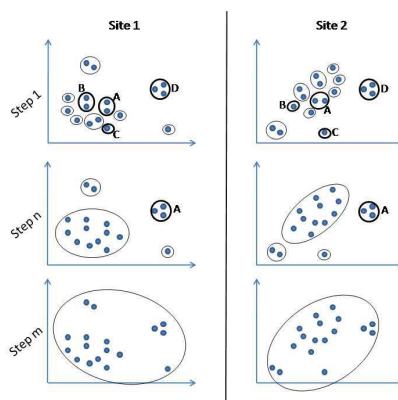


Fig. 1. Detection of common outliers in the usage patterns of two Web sites

This work is intended to explore the solutions for monitoring a network in real time. Then, the potential alarms will be triggered at each step of the monitoring (for instance with a frequency of one hour). Depending on the number of true or false alarms, the user might want to adjust n for the next step, until no (or very few) false alarm is returned. Our assumption is that common outliers, sorted by similarity from one site to another, will give the intrusions at the beginning of the list. Our challenge, in this paper, is to reply to important questions underlying our method ; what is the similarity between two usage patterns? How to separate clusters in order to give the list of outliers? How to detect common outliers? Our main algorithm, corresponding to the framework presented in this section, is given in section 4.1. Our similarity measure and our clustering algorithm are given in section 4.2. As explained in section 4.3 our outlier detection method is parameterless, thanks to a wavelet transform on the cluster distribution. In contrast to most previous methods [24, 18, 9] it does not require a percent of cluster and it does not depend on a top-n param-

eter given by the user. The correspondance between outliers of S_1 and S_2 also has to be parameterless. As explained in section 4.4 it will automatically find clusters that are close enough to trigger an alarm.

3 Related Works

Outlier detection has been extensively studied these past years, since it has a wide range of applications, such as fraud detection for credit card [1], cyber security [4] or safety of critical systems [8]. Those fields of application rely on methods to find patterns which deviate significantly from a well-defined notion of normality. The concept of outlyingness has been studied by statistics [15, 11] where statistical approaches construct probability distribution models under which outliers are objects of low probability [3, 14]. Within the context of intrusion detection, data dimensionality is high. Therefore, to improve overall performances and accuracy, it has become necessary to develop data mining algorithms using the whole data distribution as well as most of data features [10]. In this paper, we focus on clustering-based outlier detection algorithms [10, 7, 6, 16]. Such techniques rely on the assumption that normal points belong to large clusters while outliers either do not belong to any cluster [10] or form very small and tight clusters [7]. In other words, outlier detection consists in identifying among data those that are far from being significant clusters. However all the anomaly-based intrusion detection techniques suffer of the number of false alarms they trigger. On the contrary, misuse techniques (*i.e.* approaches that detect elements similar to well-known malicious usage) will precisely detect attacks but they will miss every intrusion that differs from these already known attack signatures. Therefore some works proposed collaborative frameworks in order to improve performance and both true and false alarm rates [20, 22]. These approaches rely on propagating in a distributed IDS IP blacklist after individual misuse or anomaly detection. Also this communication can lead to more accurate results, it does not allow the system to uncover totally unknown attacks or to avoid high false alarm rates. For these reasons we propose in this paper an anomaly detection approach that uses collaboration between systems in order to discriminate attacks from emerging or novel usage behaviours, thus leading to a reduced number of false alarms.

4 COD: Common Outlier Detection

The principle of COD is to perform successive clustering steps on usage patterns of different partners sites, until the number of common outliers meets the number of alarms desired by the user. We present in this section an algorithm designed for two information systems. Extending this work to more than two systems would require a central node coordinating the comparisons and triggering the alarms, or a peer-to-peer communication protocol. This is not the goal of this paper. Our objects are the parameters given to script files in the requests received on a Web site. In other words, the access log file is filtered and we only keep lines corresponding to requests with parameters to a script. For each such line, we separate the parameters and for each parameter we create an object. Let us consider, for

instance, the following request: `staff.php?FName=John&LName=Doe`. The corresponding objects are $o_1 = \text{John}$ and $o_2 = \text{Doe}$. Once the objects are obtained from the usage data of multiple Web sites, COD is applied and gives their common outliers.

4.1 Main Algorithm

As explained in section 2.2, COD algorithm will process the usage patterns of both sites step by step. For each step, a clustering result is provided and analyzed for intrusion detection. First, MD is set to obtain very tight and numerous clusters (very short similarity is allowed between two objects in a cluster). Then, MD is relaxed by an amount of 0.05 step after step in order to increase the size of resulting clusters, decrease their number and lower the number of alarms. When the number of alarms desired by the user is reached, then COD ends.

Algorithm Cod

Input: U_1 and U_2 the usage patterns of sites S_1 and S_2
and n the number of alarms.

Output: I the set of clusters corresponding
to malicious patterns.

1. Build M , the distance matrix between each pattern ;
2. $\forall p \in M, Neighbours_p \leftarrow$ sorted list of neighbours for p (the first usage pattern in the list of p is the closest to p).
3. $DensityList \leftarrow$ sorted list of patterns by density ;
4. $MD \leftarrow 0$;
5. $MD \leftarrow MD + 0.05$;
6. $C_1 \leftarrow Clustering(U_1, MD)$;
 $C_2 \leftarrow Clustering(U_2, MD)$;
7. $O_1 \leftarrow Outliers(C_1)$; $O_2 \leftarrow Outliers(C_2)$;
8. $I \leftarrow CommonOutliers(O_1, O_2, MD)$;
9. If $|I| \leq n$ then return I ;
10. If $MD = 1$ then return I ; // No common outlier
11. Else return to step 5 ;

End algorithm Cod

4.2 Clustering

COD Clustering algorithm is based on an agglomerative principle. The goal is to increase the volume of clusters by adding candidate objects, until the Maximum Dissimilarity (MD) is broken (*i.e.* there is one object o_i in the cluster such that the similarity between o_i and the candidate object o_c is greater than MD).

Similarity between objects. We consider each object as a sequence of characters. Our similarity is then based on the longest common subsequence (LCS), as described in definition 1.

Definition 1. Let s_1 and s_2 be two sequences. Let $LCS(s_1, s_2)$ be the length of the longest common subsequences between s_1 and s_2 . The dissimilarity $d(s_1, s_2)$ between s_1 and s_2 is defined as follows: $d(s_1, s_2) = 1 - \frac{2 \times LCS(s_1, s_2)}{|s_1| + |s_2|}$

Example 1. Let us consider two parameters p_1 =intrusion and p_2 =induction. The LCS between p_1 and p_2 is L =inuion. L has length 6 and the similarity between p_1 and p_2 is $d = 1 - \frac{2 \times L}{|p_1| + |p_2|} = 33.33\%$. Which also means a similarity of 77.77% between both parameters.

Centre of clusters. When an object is inserted into a cluster we maintain the centre of this cluster, since it will be used in the CommonOutliers algorithm. The centre of a cluster C is the LCS between all the objects in C . When object o_i is added to C , its center C_c is updated. The new value of C_c is the LCS between the current value of C_c and o_i .

Algorithm Clustering

Input: U , the usage patterns
and MD , the Maximum Dissimilarity.

Output: C , the set of as large clusters as possible,
respecting MD .

1. $i \leftarrow 0$; $C \leftarrow \emptyset$;
2. $p \leftarrow$ next unclassified pattern in *DensityList* ;
3. $i++$; $c_i \leftarrow p$;
4. $C \leftarrow C + c_i$;
5. $q \leftarrow$ next unclassified pattern in *Neighbours_p* ;
6. $\forall o \in c_i$
 If $d(o, q) > MD$ then return to step 2 ;
7. add q to c_i ;
8. $C_c \leftarrow LCS(C_c, q)$; // C_c is the center of C
9. return to step 5 ;
10. If unclassified patterns remain then return to step 2 ;
11. return C ;

End algorithm Clustering

4.3 Wavelet-based Outlier Detection

Most previous work in outlier detection require a parameter [24, 18, 9], such as a percent of small clusters that should be considered as outliers, or the top- n outliers. Their key idea is generally to sort the clusters by size and/or tightness. We consider that our clusters will be as tight as possible, according to our clustering algorithm and we want to extract outliers by sorting the cluster by size. The problem is to separate “big” and “small” clusters. Our solution is based on an analysis of cluster distribution, once they are sorted by size. The usual distribution of clusters is illustrated by Figure 2 (screenshot made with our real data). We propose to use a wavelet transform to cut down the distribution. In figure 2, the y axis stands for the size of the clusters, whereas their index in the sorted list is represented on x , and the two plateaux allow separating small and big clusters. With a prior knowledge on the number of plateaux (we want two plateaux, the first one standing for small clusters, or outliers, and the second one standing for big clusters) we can cut the distribution in a very effective manner. Actually, each cluster mapped to the first plateau will be considered as an outlier.

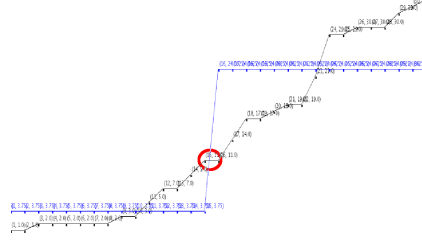


Fig. 2. Detection of outliers by means of Haar Wavelets

The wavelet transform is a tool that cuts up data or functions or operators into different frequency components, and then studies each component with a resolution matched to its scale [5]. In other words, wavelet theory represents series of values by breaking them down into many interrelated component pieces; when the pieces are scaled and translated wavelets, this breaking down process is termed wavelet decomposition or wavelet transform. Wavelet reconstructions or inverse wavelet transforms involve putting the wavelet pieces back together to retrieve the original object [23]. Mathematically, the continuous wavelet transform is defined by:

$$T^{wav} f(a, b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{+\infty} f(x) \psi^* \left(\frac{x-b}{a} \right) dx$$

where z^* denotes the complex conjugate of z , $\psi^*(x)$ is the analyzing wavelet, $a (> 0)$ is the scale parameter and b is the translation parameter. This transform is a linear transformation and it is covariant under translations and dilations. This expression can be equally interpreted as a signal projection on a function family analyzing $\psi_{a,b}$ constructed from a mother function in accordance with the following equation: $\psi_{a,b}(t) = \frac{1}{\sqrt{a}} \psi \left(\frac{t-b}{a} \right)$. Wavelets are a family of basis functions that are localized in time and frequency and are obtained by translations and dilations from a single function $\psi(t)$, called the mother wavelet. For some very special choices of a , b , and ψ , $\psi_{a,b}$ is an orthonormal basis for $L^2(\mathbb{R})$. Any signal can be decomposed by projecting it on the corresponding wavelet basis function. To understand the mechanism of wavelet transform, we must understand the multiresolution analysis (MRA). A multiresolution analysis of the space $L^2(\mathbb{R})$ consists of a sequence of nested subspaces such as:

$$\begin{aligned} \dots \subset V_2 \subset V_1 \subset V_0 \subset V_{-1} \dots \subset V_{j+1} \subset V_j \dots \\ \bigcup_{j \in \mathbb{Z}} V_j = L^2(\mathbb{R}) \\ \bigcap_{j \in \mathbb{Z}} V_j = \{0\} \end{aligned}$$

$$\forall j \in \mathbb{Z} \text{ if } f(x) \in V_j \iff f(2^{-1}x) \in V_{j+1} \text{ (or } f(2^j x) \in V_0)$$

$$\forall k \in \mathbb{Z} \text{ if } f(x) \in V_0 \iff f(x - k) \in V_0$$

There is a function $\varphi(x) \in L^2(\mathbb{R})$, called scaling function, which by dilation and translation generates an orthonormal basis of V_j . Basis functions are constructed according to the following relation : $\varphi_{j,n}(x) = 2^{-\frac{j}{2}} \varphi(2^{-j}x - n)$, $n \in \mathbb{Z}$, and the basis is orthonormal if $\int_{-\infty}^{+\infty} \varphi(x) \varphi^*(x+n) dx = \delta(n)$, $n \in \mathbb{Z}$. For each V_j , its orthogonal complement W_j in V_{j-1} can be defined as follows: $V_{j-1} = V_j \oplus W_j$ and $L^2(\mathbb{R}) = \bigoplus_{j \in \mathbb{Z}} W_j$. As W_j is orthogonal to V_{j-1} , then W_{j-1} is orthogonal to W_j , so $\forall j, k \neq j$ then $W_j \perp W_k$.

There is a function $\psi(x) \in \mathbb{R}$, called wavelet, which by dilations and translations generates an orthonormal basis of W_j , and so of $L^2(\mathbb{R})$. The basis functions are constructed as follows: $\psi_{j,n}(x) = 2^{-\frac{j}{2}}\psi(2^{-j}x - n)$, $n \in \mathbb{Z}$. Therefore, $L^2(\mathbb{R})$ is decomposed into an infinite sequence of wavelet spaces, i.e. $L^2(\mathbb{R}) = \bigoplus_{j \in \mathbb{Z}} W_j$. To summarize the wavelet decomposition: given a f_n function in V_n , f_n is decomposed into two parts, one part in V_{n-1} and the other in W_{n-1} . At next step, the part in V_{n-1} continues to be decomposed into two parts, one part in V_{n-2} and the other in W_{n-2} and so on. A direct application of multiresolution analysis is the fast discrete wavelet transform algorithm. The idea is to iteratively smooth data and keep the details all along the way. More formal proofs about wavelets can be found in [5]. We use the Haar wavelets to illustrate our outlier detection method. Let us consider the following series of values: [1, 1, 1, 2, 7, 10, 11, 12]. Its Haar wavelet transform is illustrated by the following table:

Level	Approximations	Coefficients
8	1, 1, 1, 2, 7, 10, 11, 12	
4	1, 1.5, 8.5, 11.5	0, -0.5, -1.5, -0.5
2	1.25, 10	-0.25, -1.5
1	5.625	-4.375

Then, we keep only the most two significant coefficients and we make the others zero. In our series of coefficients ([5.625, -4, 375, -0.25, -1.5, 0, -0.5, -1.5, -0.5]) the most two significant ones are 5.625 and -4, 375, meaning that the series becomes [5.625, -4, 375, 0, 0, 0, 0, 0, 0]. In the following step, the inverse operation is calculated and we obtain an approximation of the original data [1.25, 1.25, 1.25, 1.25, 10.0, 10.0, 10.0, 10.0]. This gives us two plateaus and allows cutting the series after index 4 in order to separate big and small values.

4.4 Comparing Outliers

Since we want our global algorithm to require only one parameter (the number of alarms), we want to avoid introducing a similarity degree for comparing two lists of outliers. For this comparison, CommonOutliers algorithm will use the centre of outliers. For each pair of outliers, CommonOutliers calculates the similarity between centers of these outliers. If this similarity is below the current MD (C.f. Subsection 4.2), then we consider those outliers as similar and add them to the alarm list.

Algorithm CommonOutliers

Input: O_1 and O_2 , two lists of outliers
and MD , the maximum dissimilarity.

Output: A , the list of alarms (common outliers).

1. $A \leftarrow \emptyset$
2. $\forall i \in O_1$ do
3. $\forall j \in O_2$ do
4. $centre_i \leftarrow centre(i)$;
5. $centre_j \leftarrow centre(j)$;
6. If $d(centre_i, centre_j) < MD$
 Then $A \leftarrow A + i \cup j$;

7. done ;
8. done ;
9. Return A ;

End algorithm CommonOutliers

5 Experiments

The goal of this section is to analyze our results (*i.e.* the number of outliers and true intrusions and the kind of intrusions we have detected). Our datasets come from two different research organizations; (*anonymized for submission*). We have analyzed their Web access log files from March 1 to March 31. The first log file represents 1.8 Gb of rough data. In this file, the total number of objects (parameters given to scripts) is 30,454. The second log file represents 1.2 Gb of rough data and the total number of objects is 72,381. COD has been written in Java and C++ on a PC (2.33GHz i686) running Linux with 4Gb of main memory. Parameters that are automatically generated by the scripts have been removed from the datasets since they cannot correspond to attacks (for instance “`publications.php?Category=Books`”). This can be done by listing all the possible generation of parameters in the scripts of a Web site.

5.1 Detection of common outliers

As described in Section 2.2, COD proceeds by steps and slowly increases the value of MD , which stands for a tolerance value when grouping objects during the clustering process. In our experiments, MD has been increased by steps of 0.05 from 0.05 to 0.5. For each step, we report our measures in table 1. The meaning of each measure is as follows. O_1 (resp. O_2) is the number of outlying objects in site 1 (resp. site 2). $\%_1$ (resp $\%_2$) is the fraction of outlying objects on the number of objects in site 1 (resp. site 2). For instance, when MD is set to 0.3, for site 1 we have 5,607 outlying objects, which represents 18.4% of the total number of objects (*i.e.* 30,454) in site 1. COD is the number of common outliers between both sites and $\%_{FA}$ is the percentage of false alarms within the common outliers. For instance, when MD is set to 0.05, we find 101 alarms among which 5 are false (which represents 4.9%). One first observation is that outliers cannot be directly used to trigger alarms. Obviously, a number as high as 5,607 alarms to check, even for one month, is not realistic. On the other hand, the results of COD show its ability to separate malicious behaviour from normal usage. Our false alarms correspond to normal requests that are common to both sites but rarely occur. For instance, on the references interrogation script of *anonym.lab1*, a user might request papers of “John Doe” and the request will be `publications.php?FName=John\&LName=Doe`. If another user requests papers of “John Rare” on the Web site of *anonym.lab2*, the request will be `biblio.php?FName=John\&LName=Rare` and the parameter “John” will be given as a common outlier and trigger an alarm. As we can see, $\%_{FA}$ is very low (usually we have at most 5 false alarms in our experiments for both Web sites) compared to the thousands of outliers that have been filtered by COD. Another lesson from these experiments is that a low MD implies very small clusters and numerous outliers. These outliers are shared between both sites, among which some are false alarms due to rare but common normal

usage. When MD increases, the clustering process gets more agglomerative and alarms are grouped together. Then one alarm can cover several ones of the same kind (*e.g.* the case of easter eggs explained further). At the same time, the number of outliers corresponding to normal usage decreases (since they are also grouped together). Eventually, a too large value of MD implies building clusters that do not really make sense. In this case, outliers will get larger, and the matching criteria will get too tolerant, leading to a large number of matching outliers capturing normal usage. In a streaming environment involving the real data of these experiments, one could decide to keep 70 as the number of desired alarms and watch the ratio of false alarms. If this ratio decreases, then the end-user should consider increasing the number of desired alarms.

	0.05	0.1	0.15	0.2	0.25	0.3	0.35	0.4	0.45	0.5
O_1	13197	10860	8839	7714	6547	5607	5184	4410	3945	3532
$\%_{O_1}$	43.3%	35.6%	29%	25.3%	21.5%	18.4%	17%	14.4%	12.9%	11.6%
O_2	35983	27519	24032	20948	18152	14664	12738	11680	10179	8734
$\%_{O_2}$	49.6%	37.9%	33.1%	28.9%	25%	20.2%	17.5%	16.1%	14%	12.1%
COD	101	78	74	70	67	71	71	85	89	90
$\%_{FA}$	4.9%	5.12%	4%	2.85%	1.5%	2.8%	2.8%	10.6%	11.2%	16.6%

Table 1. Results on real data

5.2 A sample of our results

None of the attacks found in our experiments have been successful on the considered Web sites. However, our security services and our own investigations allow us to confirm the intrusion attempts that have been discovered by our method:

- **Code Injection:** a recent kind of attack aims to inject code in PHP scripts by giving a URL in the parameters. Here is a sample of such URLs detected by COD:
 - <http://myweddingphotos.by.ru/images?>
 - <http://levispotparty.eclub.lv/images?>
 - <http://0xg3458.hub.io/pb.php?>

Depending on the PHP settings on the victim’s Web server, the injected code allows modifying the site. These URLs are directly, automatically and massively given as parameters to scripts through batches of instructions.

- **Passwords:** another kind of (naive and basic) attack aims to retrieve the password file. This results in outliers containing parameters like `../etc/password` with a varying number of `../` at the beginning of the parameter. This is probably the most frequent attempt. It is generally not dangerous but shows the effectiveness of our method.
- **Easter Eggs:** this is not really an intrusion but if one adds the code `?=PHPE9568F36-D428-11d2-A769-00AA001ACF42` to the end of any URL that is a PHP page, he will see a (funny) picture on most servers. Also on April 1st (April Fool’s Day), the picture will replace the PHP logo on any `phpinfo()` page. This code (as well as two other ones, grouped into the same outlier) has been detected as a common outlier by COD.

6 Conclusion

In this paper, we have proposed i) an unsupervised clustering scheme for isolating atypical behaviours, ii) a parameterless outlier detection method based on wavelets and iii) a new feature for characterizing intrusions. This new feature is based on the repetition of an intrusion attempt from one system to another. Actually, our experiments show that atypical behaviours cannot be directly used to trigger alarms since most of them correspond to normal requests. On the other hand, this very large number of outliers can be effectively filtered (reducing the amount of atypical behaviours up to 0.21%) in order to find true intrusion attempts (or attacks) if we consider more than one site. Eventually, our method guarantees a very low ratio of false alarms, thus making unsupervised clustering for intrusion detection effective, realistic and feasible.

References

1. E. Aleskerov, B. Freisleben, and B. Rao. Cardwatch: A neural network based database mining system for credit card fraud detection. In *IEEE CIFE*, 1997.
2. D. Barbara, N. Wu, and S. Jajodia. Detecting novel network intrusions using bayes estimators. In *1st SIAM Conference on Data Mining*, 2001.
3. N. Billor, A. S. Hadi, and P. F. Velleman. BACON: blocked adaptive computationally efficient outlier nominators. *Computational Statistics and Data Analysis*, 34, 2000.
4. E. Bloedorn, A. D. Christiansen, W. Hill, C. Skorupka, and L. M. Talbot. Data mining for network intrusion detection: How to get started. Technical report, MITRE, 2001.
5. Ingrid Daubechies. *Ten lectures on wavelets*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1992.
6. E. Eskin, A. Arnold, M. Prerau, L. Portnoy, and S. Stolfo. A geometric framework for unsupervised anomaly detection: Detecting intrusions in unlabeled data. *Applications of Data Mining in Computer Security*, 2002.
7. H. Fan, O. R. Zaiane, A. Foss, and J. Wu. A nonparametric outlier detection for effectively discovering top-n outliers from engineering data. In *PAKDD*, 2006.
8. R. Fujimaki, T. Yairi, and K. Machida. An approach to spacecraft anomaly detection problem using kernel feature space. In *11th ACM SIGKDD*, 2005.
9. J. Joshua Oldmeadow, S. Ravinutala, and C. Leckie. Adaptive clustering for network intrusion detection. In *8th PAKDD*, 2004.
10. E. M. Knorr and R. T. Ng. Algorithms for mining distance-based outliers in large datasets. In *24rd International Conference on Very Large Data Bases*, pages 392–403, 1998.
11. R. Kwitt and U. Hofmann. Unsupervised anomaly detection in network traffic by means of robust pca. In *International Multi-Conference on Computing in the Global Information Technology*, 2007.
12. A. Lazarevic, L. Ertöz, V. Kumar, A. Ozgur, and J. Srivastava. A comparative study of anomaly detection schemes in network intrusion detection. In *3rd SIAM DM*, 2003.
13. W. Lee and S. J. Stolfo. Data mining approaches for intrusion detection. In *7th conference on USENIX Security Symposium*, 1998.
14. W. Lee and D. Xiang. Information-theoretic measures for anomaly detection. In *IEEE Symposium on Security and Privacy*, 2001.
15. M. Markou and S. Singh. Novelty detection: a review - part I: statistical approaches. *Signal Processing*, 83, 2003.
16. S. Papadimitriou, H. Kitagawa, P.B. Gibbons, and C. Faloutsos. LOCI: fast outlier detection using the local correlation integral. In *19th ICDE*, 2003.

17. A. Patcha and J.-M. Park. An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Comput. Networks*, 51, 2007.
18. L. Portnoy, E. Eskin, and S. Stolfo. Intrusion detection with unlabeled data using clustering. In *ACM CSS Workshop on Data Mining Applied to Security*, 2001.
19. M. Roesch. SNORT, 1998.
20. A. Valdes and K. Skinner. Probabilistic alert correlation. In *Recent Advances in Intrusion Detection*, pages 54–68, 2001.
21. N. Wu and J. Zhang. Factor analysis based anomaly detection. In *IEEE Workshop on Information Assurance*, 2003.
22. V. Yegneswaran, P. Barford, and S. Jha. Global intrusion detection in the domino overlay system. In *Network and Distributed Security Symposium*, 2004.
23. Randy K. Young. *Wavelet Theory and Its Applications*. Kluwer Academic Publishers Group, 1995.
24. S. Zhong, T. M. Khoshgoftaar, and N. Seliya. Clustering-based network intrusion detection. *International Journal of Reliability, Quality and Safety Engineering*, 14, 2007.