

Anomaly Detection in Network Traffic Using Unsupervised Machine Learning

Dipali Paradhi, Mehjabeen Naghma Ansari, Sharmila More

MIT Arts, Commerce and Science College, Alandi, Pune, Maharashtra, India

Abstract: *With the increasing complexity and volume of network traffic, the detection of anomalies has become crucial for maintaining the security and efficiency of computer networks. Traditional rule-based methods often struggle to keep pace with the evolving nature of Cyber threats. In this paper, we propose utilizing unsupervised machine learning techniques for anomaly detection in network traffic. We explore various algorithms including k-means clustering, Isolation Forest, and auto encoders to identify abnormal patterns within network data without the need for labeled examples. Our experiments demonstrate the effectiveness of these approaches in detecting anomalies accurately and efficiently. Furthermore, we discuss the challenges and opportunities in deploying unsupervised machine learning for network anomaly detection in real-world scenarios. This research contributes to the advancement of Cyber security by providing novel methodologies for detecting suspicious activities within network traffic data, thereby enhancing the resilience of computer networks against emerging threats. Unsupervised methods, such as clustering algorithms like k-means or density-based techniques like DB-SCAN, can detect deviations from normal patterns in network traffic, indicating potential intrusions or anomalies. These systems analyze various features of network traffic, such as packet headers, traffic volume, and protocol behavior, to identify suspicious activity. However, they may also generate false positives and require careful tuning to balance detection accuracy and performance.*

Keywords: Cyber threats

I. INTRODUCTION

In the dynamic landscape of the Internet of Things (IoT), where interconnected devices facilitate seamless communication, the integrity of network traffic is paramount. As the IoT ecosystem burgeons, encompassing a multitude of applications across diverse domains, the susceptibility to network attacks becomes a pressing concern. The imperative to secure these networks necessitates advanced intrusion detection mechanisms, particularly those harnessing the power of unsupervised machine learning. [3] This research paper embarks on an exploration of anomaly detection in network traffic within the context of IoT environments. Fueled by the proliferation of wireless networking technologies, the vulnerabilities inherent in these systems expose them to a spectrum of network attacks. Traditional intrusion detection systems (IDS) are challenged by the dynamic nature of IoT, underscoring the need for adaptive, unsupervised machine learning algorithms. [1] Throughout this study, we scrutinize the intricacies of wireless networking protocols, dissect prevalent network attacks targeting IoT devices, and delve into the application of unsupervised machine learning in developing robust intrusion detection systems. By synergizing these elements, our research seeks to contribute to the enhancement of network security paradigms, offering a comprehensive understanding of anomaly detection as a proactive approach in fortifying IoT network traffic against emerging threats. Unsupervised machine learning plays a crucial role in anomaly detection in network traffic. Algorithms like k-means clustering, DBSCAN, and isolation forests can identify unusual patterns or outliers in network data without the need for labeled examples. By analyzing various features such as packet size, protocol type, and timing, these algorithms can flag potentially malicious activity or unusual behavior, helping to enhance network security.

K-means clustering serves as a powerful unsupervised learning technique for identifying anomalies in network traffic by partitioning the data into clusters and flagging data points that deviate significantly from the norm. [2]

Grouping Normal Behavior: Initially, K-means clustering can be applied to partition network traffic data into clusters representing typical or normal behavior. The centroids of these clusters represent the normal patterns of network traffic.

Detecting Anomalies: Once the clusters are established, any data point that does not belong to any of the clusters or is significantly distant from the centroids can be considered an anomaly.

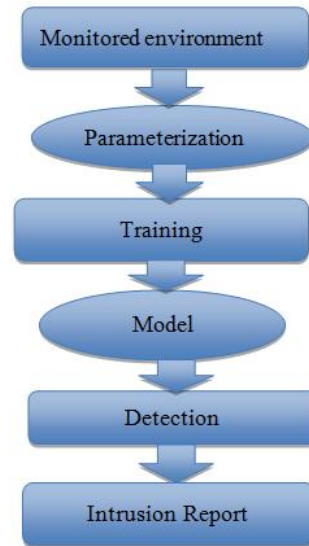
These outliers could indicate potential security threats, network intrusions, or abnormal behavior.

- **Thresholding** : By setting thresholds based on distances from cluster centroids or other statistical measures, anomalies can be detected in real-time as new data comes in. If a data point falls outside these thresholds, it is flagged as an anomaly, triggering alerts or further investigation.
- **Dynamic Updates**: Periodically, the clustering algorithm can be retrained on the latest data to adapt to changes in network behavior. This ensures that the model remains effective in detecting
- **Anomaly detection in network traffic using unsupervised machine learning** involves identifying unusual patterns that deviate from normal behavior. Techniques like clustering or autoencoders can help detect anomalies by learning the typical patterns in network data. Need more details or assistance with a specific aspect.
- **Baseline Establishment**: Create a baseline of normal network behavior by analyzing historical data to understand regular patterns and characteristics.
- **Unsupervised Learning**: Utilize unsupervised machine learning techniques such as clustering or autoencoders to detect anomalies without relying on predefined labels or signatures.
- **Feature Extraction**: Identify relevant features from network traffic data, considering factors like packet size, frequency, protocols, and communication patterns.
- **Real-time Monitoring**: Continuously monitor network traffic in real-time to promptly identify deviations or anomalies as they occur.
- **Thresholds and Alarms**: Set thresholds based on normal behavior and trigger alarms or alerts when network activities surpass these thresholds, indicating potential anomalies.
- **Behavioral Analysis**: Use behavioral analysis to understand the typical interactions between devices, services, or users, making it easier to identify deviations.
- **Adaptive Models**: Implement adaptive models that can adjust to changes in network patterns and adapt to evolving threats over time.
- **Integration with Security Information and Event Management (SIEM)**: Integrate anomaly detection with SIEM solutions for comprehensive security monitoring, analysis, and response .8)Integration with Security Information and Event Management (SIEM): Integrate anomaly detection with SIEM solutions for comprehensive security monitoring, analysis, and response
- **Feedback Loop**: Establish a feedback loop to continuously improve the anomaly detection model based on the evolving nature of network traffic and potential new threats.
- **Incident Response**: Develop a robust incident response plan to address and mitigate security incidents identified through anomaly detection. The effectiveness of anomaly detection relies on a combination of these techniques and continuous refinement based on the evolving nature of network behavior and potential threats.

Anomaly Detection Technique

The basic architecture of all anomaly detection based network intrusion detection system(A-NIDS) is similar. According to [30] and [31], basically all that are consist the following basic module.

The most common unsupervised algorithms are, K-Means, Self-organizing maps (SOM), C- means, Expectation-Maximization Meta algorithm(EM), Adaptive resonance theory (ART), Unsupervised Niche Clustering (UNC) and One-Class Support Vector Machine.



Clustering Techniques

Rawat[3] and many researcher found that Clustering techniques work by grouping the observed data into clusters, according to their similar characters or distance measure. There are least two approaches to clustering based anomaly detection. In the first , the anomaly detection model is trained using unlabeled data that consist of both normal as well as attack traffic. In the second approach, the model is trained using only normal data and a profile of normal activity is created. The idea behind the first approach is that anomalous or attack data forms a small percentage of the total data. If this assumption holds, anomalies and attacks can be detected based on cluster sizes. Large clusters correspond to normal data, and the rest of the data points, which are outliers, correspond to attacks.

Unsupervised Neural Network

The two typical unsupervised neural networks are self-organizing maps and adaptive resonance theory. They used similarity to group objects. They are adequate for intrusion detection tasks where normal behavior is densely concentrated around one or two centers, while anomaly behavior and intrusions spread in space outside of normal clusters. The Self-organizing map (SOM) is trained by an unsupervised competitive learning algorithm [4]. The aim of the SOM is to reduce the dimension of data visualization. That is, SOM outputs are clustered in a low dimensional (usually 2D or 3D) grid. It usually consists of an input layer and the Kohonen layer, which is designed as the two-dimensional arrangement of neurons that maps n dimensional input to two dimensions. Kohonen's SOM associates each of the input vectors to a representative output. The network finds the node nearest to each training case and moves the winning node, which is the closest neuron (i.e. the neuron with minimum distance) in the training course. That is, SOM maps similar input vectors onto the same or similar output units on such a two-dimensional map, which leads to self-organize the output units into an ordered map and the output units of similar weights are also placed nearby after training. SOMs are the most popular neural networks to be trained for anomaly detection tasks. For example Kayacik et al. [15], they have created three layers of employment: First, individual SOM is associated with each basic TCP feature. Second layer integrates the views provided by the first-level SOM into a single view of the problem. The final layer is built for those neurons, which win for both attack and normal behaviors. Oh and Chae [16] proposed an approach to a real-time intrusion- detection system based on SOM that groups similar data and visualizes their clusters. The system labels the map produced by SOM using correlations between features. Jun et al.[17] introduced a novel methodology to analyze the feature attributes of network traffic flow with some new techniques, including a novel quantization model of TCP states. Integrating with data preprocessing, the authors constructed an anomaly detection algorithm with SOFM and applied the detection frame to DARPA Intrusion Detection Evaluation Data. Adaptive Resonance Theory (ART). The adaptive resonance theory embraces a series of neural network models that perform unsupervised or supervised

learning, pattern recognition, and prediction. Unsupervised learning models include ART-1, ART-2, ART-3, and Fuzzy ART. Various supervised networks are named with the suffix “MAP”, such as ARTMAP, Fuzzy ARTMAP, and Gaussian ARTMAP. Amini et al. [18] compared the performance of ART-1 (accepting binary inputs) and ART-2 (accepting continuous inputs) on KDD99 data. Liao et al. [21] deployed Fuzzy ART in an adaptive learning framework which is suitable for dynamic changing environments.

Normal behavior changes are efficiently accommodated while anomalous activities can still be identified

K-Means

The K-means algorithm is a traditional clustering algorithm. It divides the data into k clusters, and guarantees that the data within the same cluster are similar, while the data in various clusters have low similarities. K-means algorithm is first selected K data at random as the initial cluster center, for the rest data add it to the cluster with the highest similarity according to its distance to the cluster center, and then recalculate the cluster center of each cluster. Repeat this process until each cluster center doesn't change.

Thus data are divided into K clusters. Unfortunately, K-means clustering is sensitive to the outliers and a set of objects closer to a centroid may be empty, in which case centroids cannot be updated [23]. [22] proposed K-means algorithms for anomaly detection.

Firstly, a method to reduce the noise and isolated points in the data set was advanced. By dividing and merging clusters and using the density radius of a super sphere, an algorithm to calculate the number of the cluster centroid was given. By a more accurate method of finding the clustering center, an anomaly detection model was presented to get a better detection effect.

Cuixiao et al. [24] proposed a mixed intrusion detection system (IDS) model. Data are examined by the misuse detection module and then the detection of abnormal data is performed by the anomaly detection module. In this model, an unsupervised clustering method is used to build the anomaly detection module. The algorithm used is an improved algorithm of K-means clustering algorithm and it is demonstrated to have a high detection rate in the anomaly detection module.

Fuzzy C-Means (FCM)

Fuzzy C-means is a clustering method, which grants one piece of data to belong to two or more clusters. It was developed by Dunn [25] and improved later by Bezdek [26], it is used in applications for which hard classification of data is not meaningful or difficult to achieve (e.g. pattern recognition). C-means algorithm is similar to K-Means except that membership of each point is defined based on a fuzzy function and all the points contribute to the relocation of a cluster centroid based on their fuzzy membership to that cluster. Shingo et al. [28] proposed a new approach called FC-ANN, based on ANN and fuzzy clustering to solve the problem and help IDS achieving higher detection rate, less false positive rate and stronger stability. Yu and Jian [29] proposed an approach integrating several soft computing techniques to build a hierarchical neuro-fuzzy inference intrusion detection system. In this approach, principal component analysis

neural network is used to reduce feature space dimensions. The preprocessed data were clustered by applying an enhanced fuzzy C-means clustering algorithm to extract and manage fuzzy rules. Another approach that uses a fuzzy approach for unsupervised clustering is presented by Shah et al. [20]. They employed the Fuzzy C-Medoids (FCMdd) in order to index cluster streams of system call, low level Kernel data and network data.

Unsupervised Niche Clustering (UNC)

(UNC) is a robust clustering algorithm, which uses an evolutionary algorithm with a niching strategy (Nasraoui et al. [5]). The evolutionary algorithm helps to find clusters using a robust density fitness function, while the niching technique allows it to create and maintain the niches (candidate clusters). Since UNC is based on genetic optimization, it is much less susceptible to suboptimal solutions than traditional techniques. The algorithm's main advantage is the ability to handle noise and to determine clusters number automatically. Elizabeth et al. [6] combined the UNC with fuzzy set theory for anomaly detection and applied it to network intrusion detection. They associated to each cluster

generated by the UNC a member function that follows a Gaussian shape using evolved cluster center and radius. Such cluster membership functions will define the normalcy level of a data sample.

Expectation-Maximization Meta Algorithm(EM)

EM is another soft clustering method based on Expectation- Maximization Meta algorithm Dempster et al. [7]. Expectation-Maximization is an algorithm for finding maximum probability estimates of parameters in probabilistic models. EM clustering algorithm alternates between performing expectation (E) step, by computing an estimation of likelihood using current model parameters (as if they are known), and a maximization (M) step, by computing the maximum probability estimates of model parameters. The model parameters new estimations contribute to an expectation step of next iteration. Hajji [8] used Gaussian mixture models to characterize utilization measurements. Model parameters are estimate using Expectation-Maximization (EM) algorithm and anomalies are detected corresponding to network failure events. Animesh and Jung [9] proposed an anomaly detection scheme, called SCAN to address the threats posed by network-based denial of service attacks in high speed networks. The noteworthy features of SCAN include: (a) it rationally samples the incoming network traffic to reduce the amount of audit data being sampled while retaining the intrinsic characteristics of the network traffic itself; (b) it computes the missing elements of the sampled audit data by using an enhanced Expectation-Maximization (EM) algorithm-based clustering algorithm; and (c) it enhances the convergence speed of the clustering process by employing Bloom filters and data summaries.

One -Class Support Vector Machine(OCSVM)

The one-class support vector machine is a very specified sample of a support vector machine which is geared for anomaly detection. The one-class SVM varies from the SVM generic version in that the resulting problem of quadratic optimization includes an allowance for a specific small predefined outliers percentage, making it proper for anomaly detection. These outliers lie between the origin and the optimal separating hyperplane. All the remaining data fall on the opposite side of the optimal separating hyperplane, belonging to a single nominal class, hence the terminology “one-class” SVM. The SVM outputs a score that represents the distance from the data point being tested to the optimal hyperplane. Positive values for the one-class SVM output represent normal behavior (with higher values representing greater normality) and negative values represent abnormal behavior (with lower values representing greater abnormality) [10]. Eskin et al. [11] and Honig et al. [12] used an SVM in addition to their clustering methods for unsupervised learning. The SVM algorithm had to be modified a little to work in unsupervised learning domain. Once it was, it performs better than both of their clustering methods. Shon and Moon [13] suggested a new SVM approach, named Enhanced SVM, which merges (soft-margin SVM method and one-class SVM) in order to provide unsupervised learning and low false alarm capability, similar to that of a supervised SVM approach. Rui et al. [14] proposed a method for network anomaly detection based on one class support vector machine (OCSVM). The method contains two main steps: first is the detector training, the training data set is used to generate the OCSVM detector, which is capable to learn the data nominal profile, and the second step is to detect the anomalies in the performance data with the trained detector.

Pros and Cons Technique for Anomaly Detection

Technique	Pros	Cons
K -Nearest Neighbor	Simplicity: KNN is a simple algorithm and easy to understand, making. No Training Phase. Non-parametric no Adaptability No Model Building	High computational cost with large datasets and numerous features. Significant memory usage as it memorizes the entire training dataset. Sensitivity to noise and irrelevant features, requiring preprocessing steps. Dependency on selecting an optimal K value, affecting performance. Tendency to favor majority classes in

		imbalanced datasets, leading to biased predictions.
Neural network	Powerful for complex relationships Adaptable to various tasks Automatic feature extraction Parallel processing capability Robust to noise	Computationally intensive Black-box nature Data dependency Prone to overfitting Require hyperparameter tuning
Decision tree	Interpretable No data preprocessing needed Efficient Can handle nonlinear relationships Provide feature importance	Prone to overfitting Instability Limited expressiveness Bias towards features with many levels Difficulty with continuous variables
Support Vector Machine	Effective in high-dimensional spaces Versatile with various kernel functions Robust to overfitting due to margin maximization Works well with small to medium-sized datasets Effective in cases where the number of features exceeds the number of samples	Computationally intensive, especially for large datasets Requires proper selection of kernel and tuning of hyperparameters Doesn't provide probability estimates directly Sensitive to noise and outliers
Self-Organizing map	Unsupervised learning with topological properties preservation Effective for dimensionality reduction and visualization	Initialization sensitivity Tendency to converge to local minima Need for tuning parameters such as learning rate and neighborhood size
	Can handle non-linear relationships in data Robust to noise Can reveal hidden structures in data	Requires careful interpretation of results Computationally intensive for large datasets
K-means	Simple and easy to implement Scalable to large datasets Efficient computational complexity Can handle large feature spaces Clusters can be easily interpreted	Requires the number of clusters (K) to be specified in advance Sensitive to initial cluster centroids May converge to local optima Assumes spherical clusters of similar sizes Doesn't work well with non-linear data distributions
fuzzy C-means	Provides soft clustering assigning membership probabilities to clusters More robust to noise and outliers compared to K-means Can handle overlapping clusters Allows gradual transition between clusters No need to specify the number of clusters precisely	Sensitive to the choice of initial cluster centers Computationally intensive, especially for large datasets Interpretation of cluster membership is more complex than in K-means Requires tuning parameters such as fuzziness coefficient May not perform well with non-convex clusters
Expectation Maximization	General framework for unsupervised learning, applicable to various probabilistic models	Sensitive to initialization of parameters Computationally intensive, especially for large datasets

	Handles missing data well Provides soft clustering with probability distributions Can model complex data distributions Guarantees convergence to local optimum	May converge to local optima Requires assumptions about data distribution (e.g., Gaussian) Interpretation of results can be complex, especially with high- dimensional data
--	---	---

II. CONCLUSION

Machine learning techniques have received considerable attention among the intrusion detection researchers to address the weaknesses of knowledge base detection techniques. anomaly detection by unsupervised techniques. Many algorithms were used to achieve good results for these techniques. propose of this paper an overview of technique of unsupervised machine learning for anomaly detection. Techniques for unsupervised such as K-Means, SOM, and one class SVM achieved better performance over the other techniques although they differ in their capabilities of detecting all attacks classes efficiently.

One of the most well known techniques for anomaly detection in network traffic using unsupervised machine learning is clustering or DBSCAN. these algorithm can identify unusual patterns or outliers in the network traffic data without the need for labeled example.

REFERENCES

- [1]. Machine Learning Technique For Anomaly Detection:An Overview(International journal of computer application (0975-8887)) authors by salimomar,AsriNgadi,HamidH.jebur
- [2]. Unsupervised Clustering approach For Network Anomaly Detection ,authorized by IwanSysrif,AdamPrugel-Bennett and Gary Wills.
- [3]. Rawat,S.2005. Efficient Data Mining Algorithms for Intrusion Detection. in Proceedings of the 4th Conference on Engineering of Intelligent Systems (EIS'04).
- [4]. Kohonen, 1995." Self-Organizing Map". Springer,Berlin,
- [5]. Nasraoui, O., Leon, E. &Krishnapuram, R. 2005.Unsupervised Niche Clustering: Discovering an Unknown Number of Clusters in Noisy Data Sets. In:GHOSH, A. & JAIN, L. (eds.) Evolutionary Computation in Data Mining. Springer Berlin Heidelberg.
- [6]. Lizabeth, L., Olfa, N. and Jonatan,G.2007. Anomaly detection based on unsupervised niche clustering with application to network intrusion detection. Proceedings of the IEEE Conference on Evolutionary Computation.
- [7]. Dempster,A., Laird, N.and Rubin, D. 1977." Maximum likelihood from incomplete Data via the EM algorithm".J. Royal Stat, Soc, Vol. 39, 1977, pp. 1–38.
- [8]. Hajji ,H." Statistical Analysis of Network Traffic for Adaptive Faults Detection". 2005. IEEE Trans. Neural Networks, Vol.16, NO5, 2005, PP. 1053-1063.
- [9]. Animesh, P. and Jung,M. 2007. "Network Anomaly Detection with Incomplete Audit Data". Elsevier Science,12 February, 2007, pp. 5-35.
- [10]. Gilles,C.,Melanie, H. and Christian, P.2004.One-Class Support vector Machines with a Conformal kernel A case study in handling class Imbalance. In: Structural syntactic and Statistical Pattern Recognition.
- [11]. Eskin,E.,Arnold,A ,Preraua,M., Portnoy,L and Stolfo,S." A geometric framework for unsupervised anomaly detection: Detecting intrusions in unlabeled data". In D. Barber and S. Jajodia (Eds.). Data Mining for Security Applications. Boston: Kluwer Academic Publishers.
- [12]. Honig, A. 2002" Adaptive model generation: An architecture for the deployment of data mining based intrusion detection systems". In D. Barbar and S. Jajodia, (Eds.), Data Mining for Security Applications.Boston: Kluwer Academic Publishers May 2002.
- [13]. Shon and Moon. 2007." A hybrid Machine Learning Approach to Network Anomaly Detection". Inf. SCI,Vol.177, NO 18, PP. 3799-3821.

- [14]. Rui, Z., Shaoyan, Z., Yang, L. and Jianmin ,J.2008.Network Anomaly Detection Using One Class Support Vector Machine. Proceedings of the International Multi Conference of Engineers and Computer Scientists.
- [15]. Kayacik, G., Zincir, H. and Heywood, M.2003. On the Capability of an SOM Based Intrusion Detection System.Proc IEEE, IJCNN.
- [16]. Oh and Chae.2008. Real Time Intrusion Detection System Based on Self-Organized Maps and Feature Correlations. The Proceedings of the Third International Conference on Convergence and Hybrid Information.
- [17]. Jun, Z., Ming, H., Hong, Z .2004. A new Method of Data Preprocessing and Anomaly Detection. Pro. of Third Inter. Conf on Machine Learning and cybernetics.
- [18]. Amini and Jalili. 2004. Network-based intrusion detection using unsupervised adaptive resonance theory.in Proceedings of the 4th Conference on Engineering of Intelligent Systems (EIS'04).
- [19]. Honig, A. 2002” Adaptive model generation: An architecture for the deployment of data mining based intrusion detection systems”. In D. Barbar and S.Jajodia, (Eds.), Data Mining for Security Applications.
- [20]. Shah,H., Undercoffer,J. and Joshi, A. 2003. Fuzzy Clustering for Intrusion Detection. the 12th IEEE International Conference on Fuzzy Systems.
- [21]. Liao,Y. , Vemuri,R. and Pasos,A. 2007.” Adaptive anomaly detection with evolving connectionist Systems”. Journal of Network and Computer Applications, Vol.30, NO1, PP. 60–80.
- [22]. LI,H 2010.Research and Implementation of an Anomaly Detection Model Based on Clustering Analysis. International Symposium on Intelligent Information Processing and Trusted Computing.
- [23]. Han, J. and Kamber, M. 2001.” Data mining: Concept and Techniques. (1th Ed) , Morgan Kaufman publishers,
- [24]. Guobing,Z.,Cuixia,Z.and Shanshan,s.2009. A Mixed Unsupervised Clustering-based Intrusion Detection Model. Third International Conference on Genetic and Evolutionary Computing.
- [25]. Dunn, J. 1973.” A fuzzy relative of the ISO data process and its use in detecting compact well-separated clusters”. Journal of Cyber natics, Vol.3(3), pp. 32–57.
- [26]. Bezdek, J. 1981.” Pattern recognition with fuzzy objective function algorithms”. Kluwer Academic Publishers, Norwell, MA, USA (1981).
- [27]. Rawat,S.2005. Efficient Data Mining Algorithms for Intrusion Detection. in Proceedings of the 4th International Journal of Computer Applications (0975 – 8887) Volume 79 – No.2, October 2013
- [28]. Shingo, M., Ci, C. Nannan, L. Kaoru, S. and Kotaro, H.” An Intrusion Detection Model based on fuzzy Class Association Rule Mining Using Genetic Network
- [29]. Yu, Z. and Jian, F. 2009 Intrusion Detection Model Based on Hierarchical Fuzzy Inference System. Second International Conference on Information and Computing Science Ictic.
- [30]. Estevez,J.,Garcya,P. and Dyaz, J. 2004.”Anomaly detection methods in wired networks: a survey and taxonomy”. Computer Networks. Vol .27, No.16, 2004, pp. 1569–84.
- [31]. Garcia,T. Diaz,V. Macia,F. and Vazquezb. 2009. Anomaly-based network intrusion detection”. Computers and security, Vol. 2 8, 2 0 0 9, pp. 1 8 – 2 8.