

Improving the Performance of Intrusion Detection System by Removing the Count Attribute from KDD Cup 1999 Data

Yash Jain¹ Pratik Jain²

^{1,2}Department of Computer Science & Engineering

^{1,2}IPS Academy, Indore, India

Abstract— Intrusion detection encompasses a range of security techniques designed to detect (and report) malicious system and network activity or to record evidence of intrusion. To understand intrusion detection one must fully understand what intrusion is. Webster's dictionary defines an intrusion as "the act of thrusting in, or of entering into a place or state without invitation or welcome". For the purpose of this article, we will define intrusion as any unauthorized system or network activity on one (or more) computer(s) or network(s). This could be an instance of a legitimate user of a system trying to escalate his privileges so that he can gain greater access to the system that he is currently assigned, or a legitimate user trying to connect to a remote port of a server to which he is not authorized. These intrusions can originate from the outside world, a disgruntled ex-employee who was fired recently, or from your trusted staff. In this paper, one scenario of false positive is considered. The false positive is the case in which the normal data is detected as attack. We are focusing on this problem with the help of an example & proposing one solution for the same problem. The KDD CUP 1999 data set is used. The result of experiment shows that if a class has higher number of counts then this class is considered as an anomaly class. But if the true person is crossing the threshold value of count it will be count as anomaly. To detect the true person & to remove false positive, one solution is proposed.

Key words: Data Mining, Anomaly Detection System (ADS), K-Means, Ensemble, Detection Rate, False Alarm Rate, False Positive, Clustering

I. INTRODUCTION

In the last two decades, computer technology have been utilized by many people all over the world in several areas. With the development of computer technology, security of network system is become an important issue, as network attack have been increased day by day over the past few years. It is very essential to find an effective way to protect our data as it contain highly sensitive information. In the present world, we are having very traditional security such as data encryption, fire wall and VPN. They are good within them. Still they are lacking to detect the attacks by a crackers. However, intrusion detection is a dynamic one which can give dynamic protection to the network security in monitoring attacks and counter attacks.

Network intrusion detection systems (NIDS) generally follow one of three design models. Each design model has its own strengths and weaknesses and many devices are a combination of the three models. These general IDS design categories are signature-based, anomaly-based and protocol modelling.

A. Signature-Based NIDS

This technology analyses packets for specific patterns related to known attacks. This is the most common design: almost all

NIDS devices have a strong dependence on signature-based detection at some level. Signature-based detection is relatively easy to understand, deploy, and update, and is good at positively identifying known attacks. However, one drawback to signature-based systems is that they may not detect unknown or modified attacks.

B. Categories of Intrusion Detection Systems

1) Signature Based Detection Systems

Signature based intrusion detection system (SBIDS) based on the known signature. This type of detection is more effective against known attacks, and it depends on the continuance updating signature. The main drawback of SBIDS is, it is unable to detect the unknown attacks and novel attacks, but the detection rate is higher than anomaly intrusion detection rates [9].

2) Anomaly Based Detection System

Anomaly based intrusion detection system (ABIDS) has attracted many researchers due to its capability of detecting novel attack. Novel attack detection is technique for identification of unidentified attack that the machine learning system is not aware during training [10]. ABIDS has two main advantages over SBIDS, First is the ability to detect unknown and "zero day" attack. This is done by comparing the normal activity with that of deviation from them. Second one is the normal activity profile are customized for system, network and therefore making it very difficult for an attacker to know with certainty what activities it can carry out without getting detected [11]. The efficiency of the system depends on how well it is implemented and tested on all protocols. The major drawback of anomaly detection is defining its rule set.

3) Protocol Modeling

Protocol modeling is performed by analyzing network traffic for abnormal protocol activity and alerting on traffic with certain designated protocols or protocols that are unknown to the system. Protocol modeling relies on several different data sources to determine what normal protocol activity is. Common sources for this data may include protocol specification RFCs, popular applications that use that protocol, and thorough analysis of normal network traffic.

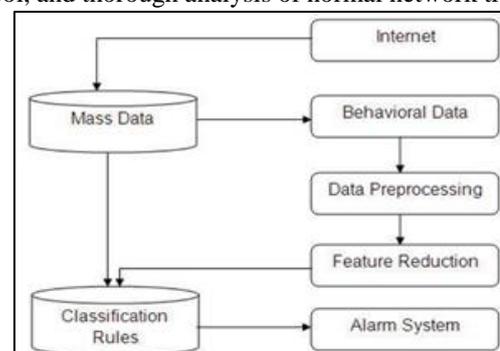


Fig. 1: Process of Intrusion Detection Systems[8]

II. LITERATURE SURVEY

V. Chandola et al, They used Hybrid detection framework that depends on data mining classification and clustering techniques [1]. Francesco Mercaldo, in his work there aim is to use data mining techniques including classification tree and support vector machines for anomaly detection. The result of experiments shows that the algorithm C4.5 has greater capability than SVM in detecting network anomaly and false alarm rate by using 1999 KDD cup data [2]. D. Denning, Algorithm utilizes a feature extraction algorithm called symbolic dynamic filtering (SDF)[3]. In SDF, time-series data are partitioned for generating symbol sequences that then construct probabilistic finite state automata (PFSA) to serve as features for pattern classification [4]. Ugo Fiore et al, in this paper, it is firstly understand the behavior of the leaning method when noise increases because it could alter the capability of extracting correct rules. Effectiveness is evaluated with 3 metrics: Max rule confidence, Precision and Recall [5]. T. Bhavani et al, they uses Cluster Analysis for Anomaly Detection. We used a simple K-mean clustering procedure. K-mean clustering is a simple, well-known algorithm. It is less computer-intensive than many other algorithms, and therefore it is a preferable choice when the dataset is large [6].

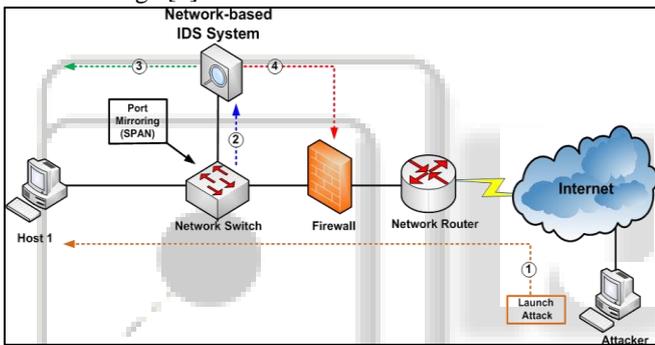


Fig. 2: Architecture of network base IDS system

T. Bhavani et al, they uses Cluster Analysis for Anomaly Detection. We used a simple K-mean clustering procedure. K-mean clustering is a simple, well-known algorithm. It is less computer-intensive than many other algorithms, and therefore it is a preferable choice when the dataset is large [6]. S. Lina et al, for High dimensional dataset these fixed number of cluster given by user are not good estimation, because it leads to inefficient data distribution or its leads to various outlier [7]. B. Thuraisingham, Network intrusion detection systems employ signature-based methods or data mining-based methods which rely on labeled training data. Anomaly network intrusion detection method based on Principal Component Analysis (PCA) for data reduction and Fuzzy Adaptive Resonance Theory (Fuzzy ART) for classifier is presented [8]. S. Wu et al, New hybrid intrusion detection system using intelligent dynamic swarm based rough set (IDS-IR) for feature selection and simplified swarm optimization for intrusion data classification [9]. B. Singh et al, The approach is studied through simulation and applied to an industrial case study. The results suggest potential use for decision making in production management. It uses Algorithm for the creation of a dynamic network based on work order data [10]. M. Xue et al, they used hybrid approach for IDS based on data mining. The main method is clustering analysis with aims of improve detection rate and decreases false alarm rate [11]. K. Wankhade et al, in this paper,

Anomaly traffic detection system based on the Entropy of network features and Support Vector Machine (SVM) are compared. Further, a hybrid technique that is combination of both entropy of network features and support vector machine is compared with individual methods [12]. A. Samad, Focuses on detailed comparative study of several anomaly detection schemas for identifying different network intrusion [13]. J. Jonathan, They present a new density-based and grid-based clustering algorithm that is suitable for unsupervised anomaly detection [14].

III. PROBLEM IDENTIFICATION

Intrusions are controlled by intrusion detection systems. An Intrusion Detection System (IDS) secures the network & protects it. It has the ability to detect anomalous activity automatically. The techniques for the detection of the anomalous activity are classified into two groups:-

A. Predefined intrusion behavior

It first stores the pattern of malicious behavior which is related to intrusion & then judges the intrusion according to the obtained pattern. It has the higher detection accuracy & having low false alarm rate. The main disadvantage of it is that it can only find predefined patterns intrusions.

B. Predefined normal behavior

It first stores the pattern of user's normal behavior into the database & then judges the normal behavior according to the stored pattern. If the deviation is huge enough, we can say that there is anomalous activity [2], [3], [4].

An Intrusion Detection System (IDS) requires high accuracy and detection rate as well as low false alarm rate. In general, the performance of IDS is evaluated in terms of accuracy (AC), detection rate (DR), and false alarm rate (FAR) as in the following formula:

- Accuracy = $(TP+TN) / (TP+TN+FP+FN)$
- Detection Rate = $(TP) / (TP+FP)$
- False Alarm Rate = $(FP) / (FP+TN)$

| Actual | Predicted Normal | Predicted Attack |
|------------|------------------|------------------|
| Normal | TN | FP |
| Intrusions | FN | TP |

Table 1: General Behavior of Intrusion Detection Data

- True positive (TP) means attack data detected as attack.
- True negative (TN) means normal data detected as normal.
- False positive (FP) means normal data detected as attack.
- False negative (FN) means attack data detected as normal.

Now, the problem is related to false positive. In which the normal data is detected as intrusion. For that we have to understand how data is considered as normal or anomaly. For that we will take data of KDD Cup 1999 data. The 1998 DARPA Intrusion Detection Evaluation Program was prepared and managed by MIT Lincoln Labs. The objective was to survey and evaluate research in intrusion detection. A standard set of data to be audited, which includes a wide variety of intrusions simulated in a military network environment, was provided. The 1999 KDD intrusion detection contest uses a version of this dataset. When I observe this data & compare the normal classes & anomaly classes. I found that it takes 41 attributes to check whether the input is of normal class or of anomaly class. The attributes are

- [7] Shih-Wei Lina, Kuo-Ching Yingb, Chou-Yuan Leec, Zne-Jung Leed “An intelligent algorithm with feature selection and decision rules applied to anomaly detection” Elsevier 2011.
- [8] Bhavani Thuraisingham “Data Mining for Malicious Code Detection and Security Applications” 2009 IEEE/WIC/ACM 2009.
- [9] Shu Wu, Member, and Shengrui Wang “Information-Theoretic Outlier Detection for Large-Scale Categorical Data” VOL. 25, NO. 3, MARCH 2013.
- [10] Bharat Singh, Nidhi Kushwaha and OP Vyas “Exploiting Anomaly Detections for high Dimensional data using Descriptive Approach of Data Mining” IEEE (ICCT) 2013.
- [11] M. Xue, C. Zhu, "Applied Research on Data Mining Algorithm in Network Intrusion Detection," jcai , pp.275-277, 2009 International Joint Conference Artificial Intelligence, 2009.
- [12] Kapil Wankhade, Mrudula Gudadhe, Prakash Prasad, “A New Data Mining Based network Intrusion Detection Model”, In Proceedings of ICCCT 2010, IEEE, 2010, pp.731-735.
- [13] Abdul Samad bin Haji Ismail “A Novel Method for Unsupervised Anomaly Detection using Unlabeled Data” IEEE 2008.
- [14] Jonathan J, Davis, Andrew J. Clark “Data preprocessing for anomaly based network intrusion detection: A review” Elsevier 2011.

