

# Enhancing the Efficiency of Detecting Intrusions using Improved PSO GMM

Anidua Bano  
Computer Science  
Sachdeva Institute of Technology Farah  
Dr. Abdul kalam technical University

Pankaj Sharma, PhD  
Computer Science  
Sachdeva Institute of Technology Farah  
Dr. Abdul kalam technical University

## ABSTRACT

Network security is one of the most significant problems in computer network management and intrusion. In recent years, the intrusion has occurred as a major area of security for the network. Each section of the attacks is considered to be a particular problem and IDS are doing well when specialized algorithms are handled. Several surveys show that penetration in the network has been steadily increased and has led to private privacy theft. It is an important platform for recent attacks. A network intrusion is illegal activities in the computer network. It is, therefore, necessary to improve an operative intrusion system. In this paper, we use improved particle swarm optimization Gaussian mixture model (IPSO GMM) to detect infiltrative inspection. This paper shows compatibility between an integrated system using an IGKM algorithm and an interchange control system used by the IPSOGMM algorithm in the KDD-99 dataset. Finding that the test was discovered uses IPSOGMM algorithm is additionally correct when compared to IGKM algorithm.

## Keywords

The intrusion detection system, data mining, KDD Cupp 99, IGKM and IPSOGMM.

## 1. INTRODUCTION

Intrusion detection systems (IDSs) have been added to the safety ratings to prevent harmful activities on a system, focusing on the network Intrusions detection systems (NIDS), because they can find extensive attacks when compared to other types of IDS. Traffic is analyzing network IDS to track upcoming up-to-date attacks. While commercial IDS is mainly used to detect attacks on a system or a host computer, the signature uses an object database. It is utilized to detect the devices used to detect tricks on a computer or a system checking the system.

Intrusion is unofficial and inconsistent activities that were distinct through Christopher Krugel et al. As a sequence of related actions by a malicious opponent, which results in the agreement of the target system "[1] an intrusion detection mechanism is an essential implement for network administrators since without such apparatus It would be impossible to analyze the huge. The amount of packets traces the current networks per second. On systems of intrusion detection after intensive research over nineteen years, the field is still open, especially for additional inquiries about the accuracy of the probe. Often undergo detected without the system.

- IDS' aims to provide IDS policy desires.
- Possible aims comprise:
- Detection of occurrences
- Preclusion of attacks

- Detection of policy destructions
- Implementation of use strategies
- Enforcement of assembly rules
- Collection of indication

In particular, IDS are used to recognize, evaluate, and then report illegal or unauthorized network activity, so that they can take suitable actions for future losses. Constructed on the information sources they use, IDS can be grouped into two categories: Network-based and host-based. The network Intrusion detection system (NIDS) analyzes the network packets seized from the network fragment. It checks the audit tracks or system calls created by individual hosts. TCP dump data to connections containing instances of network sessions.

Because of the increase in network traffic, many NIDS are using To enhance multiple sensors, distribution collaboration, and their processing capability. NIDS can too detect IP-based attacks that involve multiple computers, such as denial-of-service attacks. Host-based IDs make these attacks difficult since only the information collected from the computer system is monitored. NIDS has gained popularity because more and more systems are connected through networks. Discovery methods can be applied when IDS is used [3]. Basically, there are two discovery methods: Abuse detection and unusual discovery. The main difference between the two methods, the misuse of the identifier, is to analyze the typical attack and recognize the intrusions based on the nature of known attacks. [4]

### 1.1 Misuse detection

Discover abuse in the characteristics of a known attack. Watch this pattern and signatures of attack known on network traffic. The regularly updated database is used to keep track of known signals. Any activity that is known or damaged by a known attack is considered to be intracellular. The Fig.1 illustrates the block diagram of the misuse detection system as follows

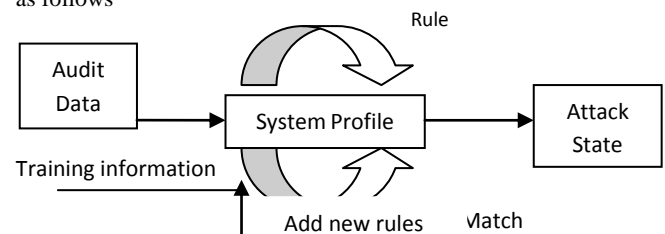
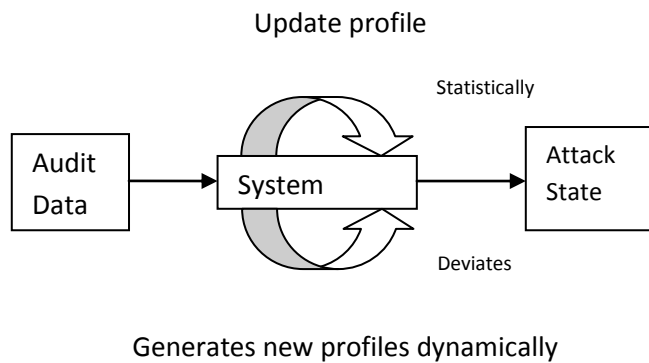


Figure 1 Misuse based detection system

### 1.2 Anomaly detection

This method is based on the discovery of the traffic's impossible. The traffic variation is observed from the normal profile. Different techniques of this technology are based on

the matrix to measure the traffic profile velocity. The following is a discovery of the Block Diagram of Analysis.



**Fig.2 Anomaly-based detection system**

This paper is sorted below. We introduce a review of the literature regarding Section 2. In section 3, it describes a specific K-meaning algorithm and network infiltration. Experimental results are discussed in Chapter 4. The end of this paper ends in Chapter 5.

## 2. LITERATURE REVIEW

ID approaches and Methodology in view of comparison has been discussed in this section.[5] suggested locating malevolent websites. To propagation of signature, the own generated program in JAVA employs for the commonly used webpages to accelerate the analytical process. A web page ends when web pages use Honeypot for browsing. The Microsoft operating system has four modules, which include behavioral record, analyzing the source code of proxy, & behavioral analyze. Even if statistical analyzes of this OS is automatic, resources at any given time are an active detector of the high-level harmful node.

One method has been suggested for IDS. It detects an attack on the network, generates unimportant, untrue and unnecessary alerts. Hence, this is a shortage of the system. The ShahidRajae Port Complex dataset also the DRPP 1999 dataset is also used aimed at an online method. As a result of this approach, the number of alerts decreased by 94.32%. In some cases, this approach provides high deterioration rates with high fake alarm rates. Used to the online analytics, but it is not appropriate, therefore newer model has been designed in order that this will decrease the fake distress rate and upsurge the rate of detection.[6]

Specified SQL injection attack is a technique to steal information from sensitive data like credit card numbers or back end databases. Various types of number and SQL Invention Finder with injection attacks helps to detect query conversion and document identification (IDS-SQLIDDS). Using five honeypot web apps for testing Developed with the MySQL & the PHP. It helps in all kinds of SQL attacks and its models.[7]

In [8] authentication, writer of the advancement persistence threat which usages various types of attack approaches for access to illegal system accessibility later ran across the network. To improve the outcome, IDS at "packet leveled" has designed as a model. This model was done by incident classes, Rules, hypotheses, & seek patterns. If you add log info, the system model will be removed from the distribution system and log-lines will also be available in the network nodes in the absence of knowledge. Using this model form

loglines, we have identified different meaningful subsections. SCADA dataset has also been used for experimental purposes since the correct ratio is 1, followed by false ratio 0.[8]

Instructed to use a method for code infection attack complete XSS (cross-site scripting) attack to adventure Javascript functions and exploit the current damage to the web application in [9]. Various kinds of the XSS attack be introduced & work on both types, finding crossed-sites disabilities intoa web application. The three steps are used for this method, including a Sanitization, Encoding, with normal Express compatibility. The entire HTML tags have been discarded using sanitation to provide malicious code.Java script code delimits possible malicious words. Pre-determined regular expressions are matched to check whether each user is valid.

Suggested [10] as a way to detect malicious JavaScript. Linear regression, 3-level, stack D-noising auto-encoders (SDAs) have a usage in a prescribed manner. In addition, the experimental results compare to other classes that provide a high optimistic rate and second unsurpassed error ratio.

A flexible and effective NIDS has been developed in [11] with the use of a deeper study-related method. This technique termed the Self-trained learning, that is assistances into linking soft maximum deterioration, sparse auto candles. Benchmark uses NSL-KDD data to implement and evaluate a specific method. Classification is done using 5 classes of binary category. The common F Scored is achieved using 75.86% 5-class classification. It is in this way to understand the normal network through Up curse Study. In this way, the drop-down ideas on the deep study, RNN (Recurrent Neural Network) & Auto Encoder have also utilized. Accurateness hasn't entirely malnourished, and the prescribed route may not be accurate. Also, Set an idea to monitor the network data flow. NSL-KDD displays an overview of the data and claims it is 75.75% accurately using six basic features.

A state-of-the-art survey of NIDS in the Health Survey suggested by the State Survey in [12]. Traditional machine learning techniques have been compiled using the Boltzmann Machine, Auto-Encoders, and RNN also with Conversational NN. The consequence is that traditional usual methodology is less and deep learning methods provide high accuracy.

To conflict a specific threat in [13] and use 100 disguised units to be used on a specified deep NN and ADAM Optimization tool worked with the Accelerated Linear Activating function. Uses the KDD data for an estimate furthermore get 99% rate of the accuracy, future long-term memory (LSTM), and the shortening of RNN model. [13]

A survey of the NIDS approach was discussed in [14]. A comprehensive category has been adopted through deep with trivial learning. The largest parts of relevant outcome have collected by these works. The overall comparison of the learning with the technologies of NIDS is displayed in table1.

## 3. PROPOSED METHODOLOGY

In the existing work, a clustering based hybrid approach was used where an optimal number of clusters is generated and later clustering is applied. A genetic algorithm was used for finding the optimal number of clusters and K-means was used for the clustering process.

In the proposed methodology, we are first applying feature selection through info-gain technique. Later particle swarm optimization (PSO) applied The optimal number of clusters is detected by GMM Tech after finding. The key PSO developed

in 1995 in KJ and Eberhurl, suggested the Global Minimum, Optimization Issue, and proposed the algorithmic critical algorithm. The printer of this algorithm is based on the movement of birds. This animal behavior is similar to optimization research. D-dimension is the best value of the position, The value of the location, the value of the location is fast and maximum value is Pg. = (xi1 , xi 2 , ..., xid )The position of the particle.

- Vi = (vi1, vi 2... vid ) velocity of particle.
- Pi = (pi1, pi 2... pid ) best personnel position.
- Pg = (pg1, pg 2... pgd) best global position.

$$pBest_i(t + 1) = \begin{cases} pBest_i(t) & \text{iff}(X_i(t + 1)) \geq f(pBest_i(t)) \\ X_i(t + 1) & \text{iff}(X_i(t + 1)) < f(pBest_i(t)) \end{cases} \quad (1)$$

Each instance of each instance has its speed and position.[6]:

$$vid = wvid + c1r1(Pid - xid) + c2r2(Pgd - xid) \quad (2)$$

$$xid = xid + vid(3)$$

The weight is used for controlling trade between global expansion and local exploration standards. c1 and c2 are scaling cables, serial number r1 and r2 at the same level (0,1)

The algebraic system of the procedure is considered algorithm1.

---

Algorithm 1: Particle Swarm Optimization

---

1. Initialization of the particles
  - a. Start Vi (t), Xi (t), ω, Vmax, q1, q2
  - b. Start the swarm size
2. Each gap calculation is a fitness value
3. Find with pBest and gBest (1)
4. Calculate the velocity of the particle using the equation (2)
5. Refresh the position of the particles using the equation (3)
6. Stop MaxIteration or some other specific criteria.

**GAUSSIAN MIXTURE MODEL**

Assume K clusters (here the number of clusters for simplicity is K). Therefore, μ and Σ estimate are considered each. If there was only one distribution, they would have the maximum possible. However, since these clusters as well as probability density are defined as a linear structure of density in these K-distributions,  $p(X) = \sum_{k=1}^k \pi_k G(X | \mu_k, \Sigma_k)$

Where  $\pi_k$  is the mixing coefficient for k-th distribution?

Calculate how to calculate the parameters using a maximum log-mode scheme  $p(X | \mu, \Sigma, \pi)$

$$\begin{aligned} & \text{Lnp}(X | \mu, \Sigma, \pi) \\ &= \sum_{i=1}^N \ln \sum_{k=1}^K \pi_k G(X_i | \mu_k, \Sigma_k). \end{aligned}$$

**4. IDS DESIGN**

The intrusion detection system enterprise is shown in Fig.3. The flow diagram shows the steps involved in the implementation of this research work. The infiltration detection system may be divided into the following components:

**4.1 Dataset**

The interface was invented using the KDD-99 cups of data. The dataset contains 42 attributes. This data there are 4.8 million incidents. The four main frauds are a database, R2L, U2R, and Probing. The above-mentioned strategy strategies can lead to more than 22 types of attacks. Details Referring Orientation [3]. Find the nature of IGKM algorithm in large datasets Comment this subset is used. In this paper, we use KDD-99 data centers that use a small subset use a pair of context. The IGKM algorithm can be detected in smaller datasets using this diagonal.

**4.2 Feature Selection**

The reason for choosing important and important features is to improve the accuracy of the structured based alert corporation (SAC) on behalf of aggressive actions from the model of alerts. This section describes feature selections of two areas, i.e. feature and more. Feature Ranking Phase Information uses the Gain Algorithm (IG) filtering approach.

Higher Information Technology The objective ranking of the structures, created on low quality, depending on the entropy is targeted. However, identifying the relationship between particular feature stage alerts may mainly involve analyzing the attributes of the alerts, and availing the basic attributes and the alerts cannot be fully detected. Therefore, the aim of this step is to contribute to the relationship between the alerts using the best discretionary ability than the initial sequence features.

**4.3 Training phase**

It is a process in the training phase. We produce algorithms by learning recognized inputs.

In this phase, IGKM train with algorithm attributes

Low data. The number of database clusters and its number

The appropriate value is to produce clusters

From fitness activity. IGKM at the training stage

The algorithm is trained in 60% of the KDD-99 content

Clusters have been created. To get more optimized clusters and the number of generations was ten.

**4.4 Testing phase**

The process is going on during verification

We check that the output is given by unknown inputs

Right? Regular input values from this stepAn input is provided for the remaining 40% KDD-99 dataset. The clusters created Observing the type of attack in the training phase.

**4.5 Classifier**

Uses a classifier to test intrusion detectionWhether the production of the algorithm was correct or not. Check the accuracy of the output Use your classifier ID-mapping. Attribute agents reduced the attribute portals for the data tab at the time of active reloading, which is the attribute ID number. The ID number indicates the official KDD-99 data

corresponding instruction. Run the output ID of the ID number as a reference cord.

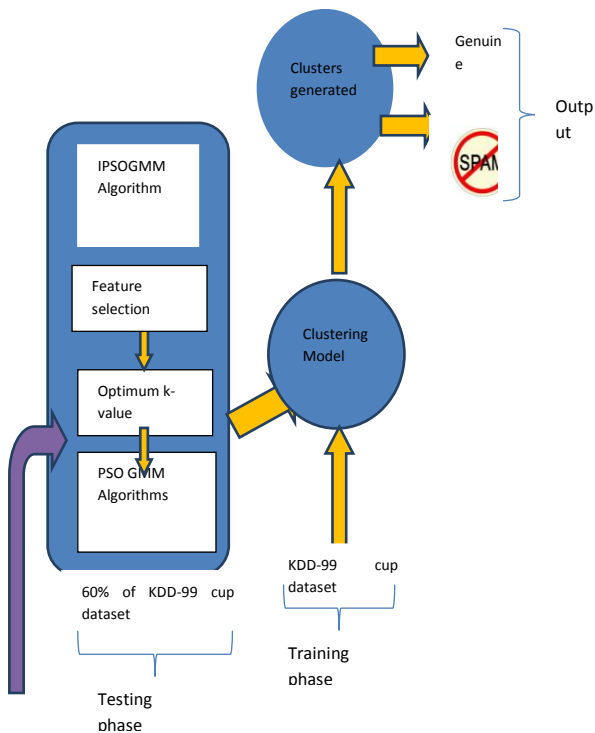


Figure 3 Flow diagram of the intrusion detection system that uses IGKM algorithm.

## 5. RESULT DISCUSSION AND ILLUSTRATIONS SECTIONS

The output evaluation is complete consuming a sequence of factors-

True positive (TP): This specifies the no. of positives

Tuples labeled correctly by a classifier

False positive (FP): This mentions to the no. of negative

Tuples that were incorrectly labeled by the classifier.

False Negative (FN): These are good tuples Uncorrected negativities.

Precision: The number of correct positive ones

The number of false positives.

$$Precision = (TP/FP)$$

Remember: the ratio of the real positive to the number

The number of false positive and false negatives.  $Recall = TP/(FP + FN)$

Accuracy (ACC): This is the total accuracy of classifiers.

$$ACC = (Precision/Recall)$$

Table 1 Comparison of accuracy of IGKM and IPSOGMM algorithm for KDD-99 dataset

Algorithm	Precision	Recall	Accuracy
IGKM	0.831394	0.831394	0.831081
IPSO GMM	0.891959	0.891959	0.884766

Compare current and research activities with IGKM and IPSOGMM in Table 1. In comparison, the suggestion for improving accuracy with better accuracy and pre-orders than previous creations.

Second, the graph shows the fitness values for research activities. The graph shows that IPSOGMM is a higher fitness value than IGKM.

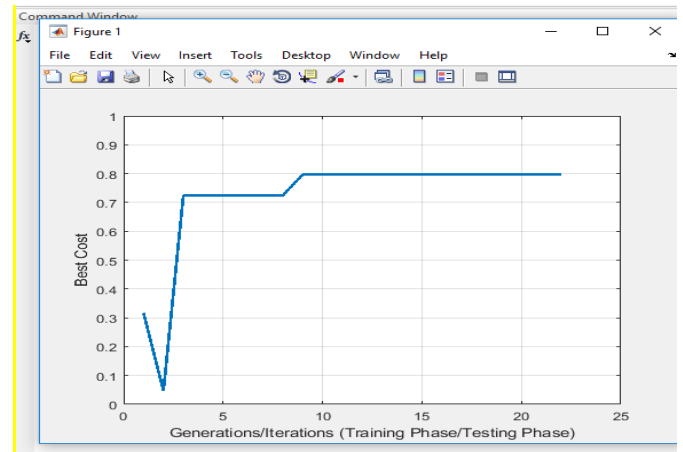


Figure 4 Graph of fitness of KDD-99 dataset of IGKM

The fitness functionality used in these algorithms is the most important priority of the K prior to implementing the cluster stage under training and testing stage.

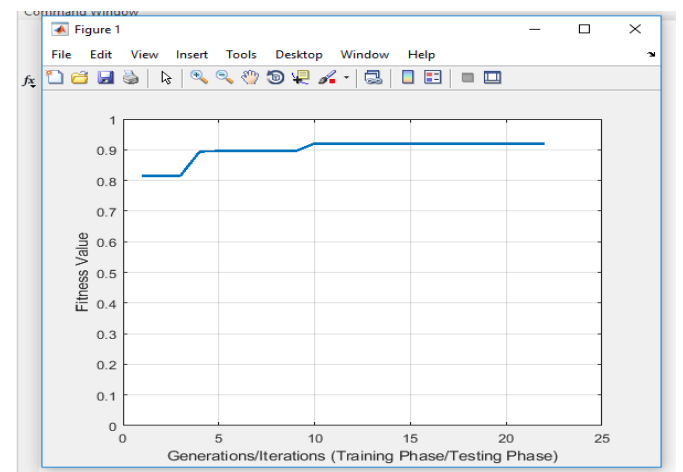


Figure 5 Graph of fitness of KDD-99 dataset of IPSOGMM

Graphs show 6 and 7 on the comparison of time and research activities time combinations. Time complexity is an idea in computer science that measures the amount of time measured by a set of code or algorithm, act or operates as an input operation.



Figure 6 Time complexity of the IGKM

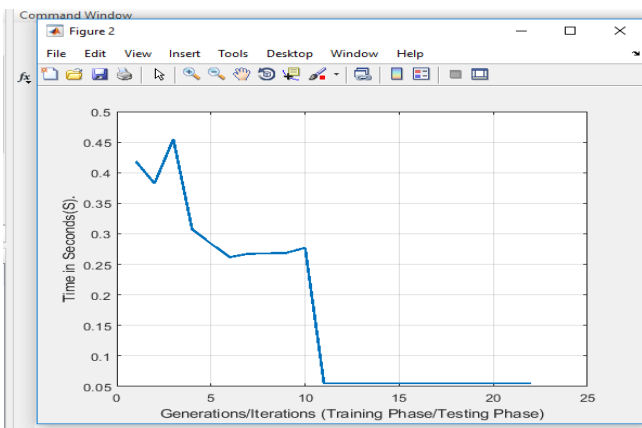


Figure 7 Time complexity of the IPSOGMM

Graph of figure 7 shows the time complexity of the proposed work which stands better than that of the IGKM.

## 6. CONCLUSION

Intrusion crimes increase every day. Therefore, the infiltrative inert detection system detects when the Intrusion explores the query system, the traditional clustered algorithm. In this paper, we have developed an infiltrated query used IPSOGMM. The number of algorithms and clusters (k) to predict the type of intrusion is not predetermined. The most optimal value of K is used by fitness activity. This helps to make optimized clusters efficiently. We use that search detector system from the experiments conducted in this paper. IGKM algorithm shows lesser accuracy on the dataset used but, in the case of IPSOGMM on same dataset intrusion detection system compared to the intrusion scanner using IGKM clustering algorithm shows a relatively high degree of accuracy.

## 7. REFERENCES

- [1] Zhang Z, Shen H. Application of online-training SVMs for real-time intrusion detection with different considerations. *Computer Communications*. 2005; 28(12):1428–42.
- [2] Shyu ML, Chen S, Sarinnapakorn K, Chang L. A novel anomaly detection scheme based on principal component classifier. *Proceedings of the IEEE Foundations and New Directions of Data Mining Workshop, in conjunction with the Third IEEE International Conference on Data Mining (ICDM03)*, 2003. p. 172–79.
- [3] Denning DE. An Intrusion-Detection Model. *IEEE Transactions on Software Engineering*. 2006; SE-13(2):222–32.
- [4] Lee W, Stolfo SJ. A framework for constructing features and models for intrusion detection systems. *ACM Transactions on Information and System Security*. 2000; 3(4):227–61.
- [5] Landgrebe TCW, Pavel P, Duin RPW, Bradley AP. Precision-Recall Operating characteristic (PROC) curves in imprecise environments. *Proceedings of 18th International Conference on Pattern Recognition, ICPR2006, HongKong*. 2006; 4. p. 123–27.
- [6] Wang W, Guan XH, Zhong X. Processing of massive audit data streams for real time anomaly intrusion detection. *Computer communications*. 2008; 31(1):58–72.
- [7] Garcia-Teodoro P, Diaz-Verdejo J, Macia-Fernandez G, Vazquez E. Anomaly-based network intrusion detection: techniques, systems and challenges. *Computer Security*. 2009; 28(1-2):18–28.
- [8] MIT Lincoln Labs. DARPA intrusion detection evaluation [Online]. 2014 Nov. Available from: <http://www.ll.mit.edu/mission/communications/ist/corpora/ideval/index.html>
- [9] Lippmann RP, Fried DJ, Graf I, Haines JW. Evaluating intrusion detection systems: The 1998 DARPA off-line intrusion detection evaluation. *Proceedings of the 2000 DARPA Information Survivability Conference and Exposition (DISCEX'00)*, Hilton Head, SC. 2000; 2. p. 12–26.
- [10] Tavallaee M, Bagheri E, Lu W, Ghorbani AA. Detailed analysis of the KDD CUP 99 Dataset. *Proceedings of the IEEE Symposium on Computational Intelligence in Security and Defense Applications*. 2009; 1–6
- [11] Tsai C-F, Hsu Y-F, Lin C-Y, Lin W-Y. Intrusion detection by machine learning: A Review. *Expert Systems with Applications*. 2009; 36(10):1994–2000.
- [12] Witten IH, Frank E, Hall MA. *Data Mining-Practical Machine Learning Tools and Techniques*. Morgan Kaufmann: San Francisco, CA, 2011.
- [13] Wu X, Kumar V, Quinlan JR, Ghosh J, Yang Q, Motoda H, McLachlan GJ, Ng A, Liu B, Yu PS, Zhou Z, Steinbach M, Hand DJ, Steinberg D. *Top Ten Data Mining Algorithms*. Knowledge and Information Systems Journal, Springer-Verlag London. 2007; 14(1):1–37.
- [14] Gaffney JE, Ulvila JW. Evaluation of intrusion detectors: A decision theory approach. *Proceedings of the IEEE Symposium on Security and Privacy, S&P'01*, Oakland, CA, USA. 2001; 50–61
- [15] Apte C, Weiss S. Data mining with decision trees and decision rules. *Future Generation Computer Systems*. 1997; 13(2-3):197–210.
- [16] AbdJalil K, Kamarudin MH, Masrek MN. Comparison of Machine Learning algorithms performance in detecting network intrusion. *2010 International Conference on Networking and Information Technology (ICNIT)*, Manila, IEEE. 2010. p. 221–26.