# Principle Components Analysis and Multi Layer Perceptron Based Intrusion Detection System

**Najla B. Ibraheem**          **Muna M. T. Jawhar**          **Hana M. Osman**

*College of Computer Sciences and Mathematics*
*University of Mosul, Mosul, Iraq*

## ABSTRACT

Security has become an important issue for networks. Intrusion detection technology is an effective approach in dealing with the problems of network security. In this paper, we present an intrusion detection model based on PCA and MLP. The key idea is to take advantage of different feature of NSL-KDD data set and choose the best feature of data, and using neural network for classification of intrusion detection. The new model has ability to recognize an attack from normal connections. Training and testing data were obtained from the complete NSL-KDD intrusion detection evaluation data set.

**Keywords:** Intrusion Detection, PCA, MLP.

مبدأ تحليل المكونات ونظام كشف التطفل القائم على طبقة متعددة

هناء محمد عصمان          منى جواهر          نجلاء بديع ابراهيم

كلية علوم الحاسوب والرياضيات

جامعة الموصل، الموصل، العراق

الملخص

أصبحت السرية قضية مهمة في الشبكات والاتصالات. تقنيات كشف التطفل هو فرع فعال في التعامل مع مشكلة سرية الشبكات والانترنيت. في هذا البحث تم تقديم نموذج من كشف التطفل بالاعتماد على تقنية (PCA) وعلى شبكة (MLP). الفكرة الأساسية هي في اخذ أفضل حقول من بيانات كشف التطفل المعتمدة (NSL-KDD) من مجموعة حقول مختلفة واستخدام الشبكات العصبية للتصنيف في نظام كشف التطفل. النموذج الجديد له القدرة على تمييز الاتصالات المصابة من الاتصالات السليمة. تم اخذ بيانات التدريب والاختبار من البيانات المعتمدة (NSL-KDD) كاملة.

الكلمات المفتاحية: كشف التطفل، PCA، MLP.

## 1. Introduction

Fast few years have witnessed a growing recognition of intelligent techniques for the construction of efficient and reliable Intrusion Detection Systems (IDS). Due to increasing incidents of cyber-attacks, building effective Intrusion Detection Systems are essential for protecting information system security, and yet it remains an elusive goal and a great challenge.

In general, the techniques for Intrusion Detection (ID) fall into two major categories depending on the modeling methods used: misuse detection and anomaly detection. Misuse detection is based on the knowledge of system vulnerabilities and known attack patterns, while anomaly detection assumes that an intrusion will always reflect some deviation from normal patterns. Many AI techniques have been applied to both misuse detection and anomaly detection. Pattern matching systems like rule-based expert systems, state transition analysis, and genetic algorithms are direct and efficient ways to implement misuse detection. On the other hand, inductive sequential patterns, artificial neural networks, statistical analysis and data mining methods have been used in anomaly detection [1].

Architecturally, an intrusion detection system can be categorized into three types host based IDS, network based IDS and hybrid IDS [2] [3]. A host based intrusion detection system uses the audit trails of the operation system as a primary data source. A network based intrusion detection system, on the other hand, uses network traffic information as its main data source. Hybrid intrusion detection system uses both the methods [4]. However, most available commercial IDS's use only misuse detection because most developed anomaly detector still cannot overcome the limitations (high false positive detection error, the difficulty of handling gradual misbehavior and expensive computation [5]). This trend motivates many research efforts to build anomaly detectors for the purpose of ID [6].

We organize this paper as follows, section 2 provides brief introduction about PCA and Neural Network, section 3 presents previous work, section 4 explain the model designer, section 5 discusses the experiments results followed by conclusion.

## 2. PCA and Neural Network

Principal Component Analysis (PCA) is an effective statistical technique for reducing the dimensions of a given unlabeled high-dimensional dataset while keeping its spatial characteristics as much as possible by performing a covariance analysis between factors. As such, it is suitable for data sets from multiple dimensions field of application, such as image compression, pattern recognition (face recognition in particular), gene expression, data clustering and traffic flow events intrusion detection. One of the main advantages of PCA is that you can compress the data, i.e. by reducing the number of dimensions, without much loss of information.

Now it is mostly used as a tool in exploratory data analysis and for making predictive models. PCA can be done by eigen value decomposition of a data covariance matrix or singular value decomposition of a data matrix. PCA is also known as the discrete Karhunen-Loeve transformation, or the Hotelling transformation[7].

An increasing amount of research in the last few years has investigated the application of Neural Networks to intrusion detection. If properly designed and implemented, Neural Networks have the potential to address many of the problems encountered by rule-based approaches. Neural Networks were specifically proposed to learn the typical characteristics of system's users and identify statistically significant variations from their established behavior. In order to apply this approach to Intrusion Detection, we would have to introduce data representing attacks and non-attacks to the Neural Network to adjust automatically coefficients of this Network during the training phase. In other words, it will be necessary to collect data representing normal and abnormal behavior and train the Neural Network on those data. After training is accomplished, a certain number of performance tests with real network traffic and attacks should be conducted. Instead of processing program instruction sequentially,

Neural Network based models on simultaneously explorer several hypotheses making the use of several computational interconnected elements (neurons), this parallel processing may imply time savings in malicious traffic analysis [8].

## 3. Previous Works

Mrutyunjaya Panda et al. [9] use discriminative multinomial Naïve Bayes with various filtering analysis in order to build a network intrusion detection system, they perform 2 class classifications with 10-fold cross validation for building the model . In [10] Shilpa lakhina et al. propose a new hybrid algorithm PCANNA (principal component analysis neural network algorithm) is used to reduce the number of computer resources, both memory and CPU time required to detect attack. The PCA transform used to reduce the feature and trained neural network is used to identify any kinds of new attacks. The model gives better and robust representation of data as it was able to reduce features resulting in a 80.4% data reduction, approximately 40% reduction in training time and 70% reduction in testing time is achieved. In [11 ] Syed Muhammad Aqil develops intrusion detection system by using principle component analysis and Neural Network the authors use four Multi Layer (MLP) working in parallel for each attack with the normal dataset such as normal vs. probe, normal vs. DoS, normal vs. U2R and normal vs. R2L.

## 4. Experiment Design

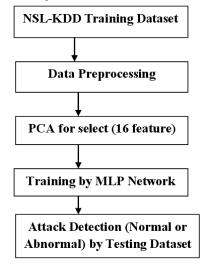The block diagram of the hybrid model is shown in the following figure (1)



**Figure (1).** The Block Diagram of the Model

### A. NSL-KDD Data Set

KDD Cup 1999 intrusion detection benchmark dataset is used by many researchers in order to build an efficient network intrusion detection system [12]. However, recent study shows that there are some inherent problems present in KDD Cup 1999 dataset .The first important limitation in the KDD Cup 1999 dataset is the huge number of redundant records in the sense that almost 78% training and 75% testing records are duplicated, as shown in Tables 1 and 2 [13]; which cause the learning algorithm to be biased towards the most frequent records, thus prevent it from recognizing rare attack records that fall under U2R and R2L categories. At the same time, it causes the evaluation results to be biased by the methods which have better detection rates on the frequent records. This new dataset, NSL-KDD dataset is used for

our experimentation and is now publicly available for research in intrusion detection. It is also stated that though the NSL-KDD dataset still suffers from some of the problems discussed in [14] and may not be a perfect representative of existing real networks, it can be applied an effective benchmark dataset to detect network intrusions. In this NSL-KDD dataset, the simulated attacks can fall in any one of the following four categories [15]:

- DOS (Denial of Service): an attacker tries to prevent legitimate users from using a service e.g. TCP SYN Flood, Smurf.

- Probe: an attacker tries to find information about the target host. For example: scanning victims in order to get knowledge about available services, using Operating System.

- U2R (User to Root): an attacker has local account on victim's host and tries to gain the root privileges.

- R2L (Remote to Local): an attacker does not have local account on the victim host and try to obtain it.

**Table 1.** Statistics of redundant records in the kdd train set [13]

|  | Original Records | Distinct Records | Reduction Rate |
|---|---|---|---|
| Attacks | 3,925,650 | 262,178 | 93.32% |
| Normal | 972,781 | 812,814 | 16.44% |
| Total | 4,898,431 | 1,074,992 | 78.05% |

**Table 2.** Statistics of redundant records in the kdd test set [13]

|  | Original Records | Distinct Records | Reduction Rate |
|---|---|---|---|
| Attacks | 250,436 | 29,378 | 88.26% |
| Normal | 60,591 | 47,911 | 20.92% |
| Total | 311,027 | 77,289 | 75.15% |

### B. Data Preprocessing

Some features have symbolic form (e.g. Protocol type ,Service ,Flag) were converted into numerical ones by assigning a unique number for each feature from the range [1.. no. of the values in the feature] ,lower iteration value takes no.1 and the upper iteration value takes no. equal number of the values within the feature.

### C. Principle Components Analysis (PCA)

The basic knowledge of PCA requires the covariance matrix for the features in the training set. The covariance matrix is defined by

$$Cov_{ij} = \frac{1}{MN} \sum_{k=1}^{M} \sum_{l=1}^{N} \left(X_i(k,l) - M_i\right)\left(X_j(k,l) - M_j\right) \qquad \ldots(1)$$

Where M,N number of the records in training set, number of features in each record, $i$ location of feature in record and j location of the record in dataset, $M_i$, $M_j$ mean of feature i,j.

The mean (μ) is defined by the following Equation:

$$M_i = \frac{1}{MN} \sum_{k=1}^{M} \sum_{l=1}^{N} X_i(k,l) \qquad \qquad …(2)$$

By using Jacobi's Method, we find eigen values as the following steps.

1. Find the largest element in the square matrix that is not in the main Diagonal
2. Find the angle θ

$$\theta = \frac{1}{2} \arctan\left(2a_{ik} / (a_{ii} - a_{kk})\right) \qquad \text{If } a_{ii} \neq a_{kk} \qquad …(3)$$

$$\theta = \begin{cases} \pi/4 & when \ a_{ik} > 0 \\ -\pi/4 & when \ a_{ik} < 0 \end{cases} \qquad \text{If } a_{ii} = a_{kk} \qquad …(4)$$

3. Rotation can be done by the following:

$$d_{ii} = \frac{1}{2}\left(a_{ii} + a_{kk} + \sigma R\right) \qquad …(5)$$

$$d_{kk} = \frac{1}{2}(a_{ii} + a_{kk} - \sigma R) \qquad …(6)$$

$$d_{ik} = d_{ki} = 0 \qquad …(7)$$

Find R by the following Equation

$$R = \sqrt{(a_{ii} - a_{kk})^2 + 4a_{ik}^2} \qquad …(8)$$

Find the value α

$$\sigma = \begin{cases} 1 & if \quad a_{ii} \geq a_{kk} \\ -1 & if \quad a_{ii} < a_{kk} \end{cases} \qquad …(9)$$

Find the other elements of the rotation matrix by the two following equations:

$$d_{ir} = a_{ir} \cos\theta + a_{kr} \sin\theta \qquad …(10)$$

$$d_{kr} = -a_{ir} \sin\theta + a_{kr} \cos\theta \qquad …(11)$$

4. Rearrange the steps from (1-3) until we get the elements of off-diagonal near the zero[16] .

Steps for executing PCA algorithm

1. Reading training NSL-KDD data set.
2. Processing data mentioned above in section B.
3. Calculate Variance/Covariance matrix for the features in every record of the training data.
4. Calculate Eigen vector of Variance/Covariance matrix as follows:
   A. Find the largest element in the matrix.
   B. Find the angle of Rotation.
   C. Find the elements of rotating matrix.
   D. Rearrange the steps from (A - C) until, we get the elements of off-diagonal near the zero.
5. Calculate the values of Eigen vector from the resulted matrix and put it in the Eigen matrix.
6. Arrange the Eigen matrix.

*D. MLP Algorithm*

The anomaly detection is to recognize different authorized system users and identify intruders from that knowledge. Thus, intruders can be recognized from the distortion of normal behavior. Multi-layer feeds forward networks (MLP) is used in this work. The number of hidden layers and the number of nodes in the hidden layers, were also determined based on the process of trial and error. We choose several initial values for the network weight and biases. Generally, theses are chosen to be small random values. The Neural Network was trained with the training data which contain normal and attack records. When the generated output result doesn't satisfy the target output result, adjust the error from the distortion of target output. Retrain or stop training the network depending on this error value. Once, the training is over, the weight value is stored to use in recall stage. In training stage, we used different network architectures with different training algorithms to find the best architecture with a good result. Resilient back propagation and Levenberg- marquardt with two hidden layers were best result from the others. After many experiments to the best features of data which is resulted from PCA algorithm, we take 16 features from 41. The architecture of Multi-layer feeds forward networks consisted from 16 nodes in input layer, 10 nodes, in the first hidden layer, 5 nodes in the second hidden layer, and 1 node in the output layer is illustrated in the following figure.
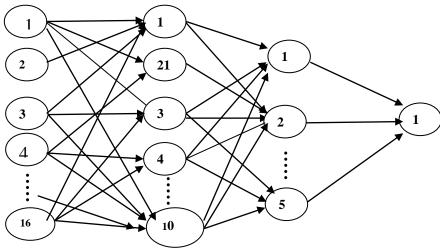


**Figure 2.** The Architecture of the MLP

The goal which is used in the algorithm was 0.001, and the epochs number was 1000. The training time for Resilient back propagation was 50 seconds and the training time for Levenberg- marquardt was 12 minutes. While, the testing time for Resilient back propagation was 17.939403 seconds and the testing time for Levenberg- marquardt was 17.293176 seconds.

The result of recall stage of two algorithms and the previous works is shown in the following table.

**Table (3).** The result of recall stage of two algorithms

| Model | Data set | DR% | FP% |
|---|---|---|---|
| SOM [17] | NSL-KDD | 64% | - |
| ESC-IDS [18] | KDD | 95.3% | - |
| Fuzzy Inference System[18] | KDD | 98% | - |
| PCA- GA[15] | NSL-KDD | 91.6% | 0 |
| Online BPN[19] | KDD | 91.50% | 2.68% |
| Resilient back propagation | NSL-KDD | 99.25% | 0.86% |
| Levenberg-Marquardt | NSL-KDD | 95.34% | 5.23% |

## 5. Conclusions

The main contribution of the present work is to achieve a classification model with a high intrusion detection Rate and with low false negative, this was done through the design of a classification model for the problem using PCA and MLP neural network for the detection of attacks. The first stage of the model is PCA, to find the best filed from the NSL-KDD dataset, we chose 16 features from 41 features. The second stage of the model is MLP neural network which is used for the classification of normal connection from attack connection. After many experiment on the Neural Network by using different training algorithms and object functions, we observe that Resilient back propagation with sigmoid function is the best one for classification. We used two hidden layers, 10 nodes in the first hidden layer and 5 nodes in the second hidden layer. We used the complete NSL_KDD dataset which are  125973 records for the training stage and 22544 records for testing stage.

## *REFERENCES*

[1]     Mahmood Hossain ,2011, "Data Mining Approaches For Intrusion Detection: Issues And Research Directions" , Department of Computer Science, Mississippi State University, MS 39762, USA.

[2]     M. Jawhar, and Mehrotra M., 2009, "Intrusion Detection System: A design perspective", 2rd International Conference on Data Management, IMT Ghaziabad.

[3]     M. Panda, and Patra M., 2009, " Building an efficient network intrusion detection model using Self Organizing Maps", proceeding of world academy of science, engineering and technology, Vol. 38 Feb.

[4]     M. Khattab  Ali, Venus W, and Mamoun Suleiman Al Rababaa, 2009, "The Affect of Fuzzification on Neural Networks Intrusion Detection System", IEEE.

[5]     Mykerjee, Heberlein L.T., and Levitt K.N., "Network Intrusion Detection", IEEE Networks, Vol. 8, No.3,1994. PP.14-26.

[6]     W. Jung K., 2002, "Integration Artificial Immune Algorithms for Intrusion Detection", dissertation in University of London, PP.1-5.

[7]     radi, A. Kartit,et. al, 2011, "On the Three Levels Security Policy Comparison between PCA and SVM", IJCSNS International Journal of Computer Science and Network Security, VOL.11 No.2.

[8]     L.de Sá Silva, Adriana C. Ferrari dos Santos, José Demisio S. da Silva, and Antonio Montes, 2004. "A Neural Network Application for Attack Detection in Computer Networks", Instituto Nacional de Pesquisas Espaciais – INPE, BRAZIL.

[9]     Mrutyunjaya Panda, Ajith Abraham,     Manas Ranjan Patra,     2010, "Discriminative Multino  mial Naïve Bayes for Network Intrusion Detection . http://www.softcomputing.net/ias10_panda.

[10]    Shilpa lakhina, Sini Joseph and Bhupendra Verma, 2010, "Feature Reduction Using Principal Component Analysis for   Effective Anomaly–Based Intrusion Detection on NSL-KDD", International Journal of Engineering Science and Technology, Vol. 2(6), 1790-1799.

[11]    Syed Muhammad Aqil Burney, M. Sadiq Ali Khan and  Dr.Tahseen A. Jilani, 2010. " Feature Deduction and Ensemble Design of Parallel Neural Networks for Intrusion Detection System, (IJCSE) International Journal of Computer Science and Network Security,Vol.10 No.10.

[12]    KDDCup 1999 Dataset. Available at:
http:// kdd.ics.uci.edu/databases/kddcup1999.html

[13]    M. Tavallaee, E. Bagheri, W. Lu, and Ali. A. Ghorbani.,2009, "A detailed analysis of the KDDCup 1999" IEEE.

[14]    J. McHugh, 2000, " Testing Intrusion Detection System: a Critique of the 1998 and 1999 DARPA Intrusion Detection System", Evaluations as performed by Lincoln Laboratory, ACM Transaction on Information and system security, Vol. 3, No. 4, pp.262-294

[15] Hana M. Osman, 2012, " Investigation of Applying Genetic    Algorithm Based - Intrusion   Detection and Classification System to NSL-KDD  Dataset" MSc. Thesis, Computer Science College, University of  Mosul ,Iraq.

[16] Muna J. Alshamdeen, 2011, " The Best Band Selection Using Hybrid Techniques Applied on Remote Sensing Data", MSc. Thesis, Computer Science College, University of  Mosul ,Iraq.

[17] Ritu Ranjani Singh, Prof. Neetesh Gupta, 2010, " To Reduce      the False Alarm in Intrusion Detection System using self Organizing   Map", International journal of Computer Science and its Applications.

[18] Wafa S. Al-Sharafat, Reyadh Sh.Naoum, 2009, "Adaptive Framework   for Network Intrusion Detection by Using Genetic-Based Machine Learning Algorithm", IJCSNS International Journal of Computer Sciences and Network Security, Vol. 9 No.4.

[19] Ibraheem M. Ahmed Al-Haleema, 2011, " Development of Network-Based Intrusion Detection System  Using Artificial Neural  Networks", MSc. Thesis, Computer Science College, University of  Mosul, Iraq.