# Enhanced Intrusion Detection System for Cloud Environment Using Machine Learning Techniques

[1]G.Monika and [2]Y.Kalpana

[1]Research Scholar, Department of Computer science, VISTAS, Chennai, India.
E-mail: aakinomm@gmail.com
[2]Professor, Department of Computer science, VISTAS, Chennai, India.
E-mail: ykalpanaravi@gmail.com

## Abstract

Tremendous growth in Cloud computing environment brings great happiness to the IT sector. Easy access and pay per use are the attractive glimpse towards the cloud environment. Users can't put out the cloud though there may be many threatening security problems. This is because they tasted the cloud easy deployment anywhere. But the security compliance is not yet attained its goal. Cloud has endless boundaries in the field so there is much crucial vulnerability arising every time.IDS being a very good solution to switch off the vulnerabilities. This paper is going to discuss about the cloud IDS which is tested using the feature selection method and eliminates unwanted attributes to gain the effective ids. The feature selection takes only the essential attributes from the total 42 features to bring enhanced model to suit the Cloud computing. The effective ids is found using NSL KDD dataset applying on the methods that is MLP and J48 algorithm to bring higher accuracy rate to improve the intrusion detection in the Cloud environment.

**Key words:** IDS, Feature selection, security, cloud computing, dataset.

*Alpha Publishers*

# 1 Introduction

Nowdays the usage of the cloud computing has increased which also brings the different kind of attacks, security issues the other side. The IDS will become the trump card to reduce these security issues and resolve the attacks [1].The virtualization environment in cloud computing bring exhastic growth in a shared environment as well as the way to the security issues The cloud is used to access our already saved data from various sources through the cloud. The cloud makes the user to deploy anywhere at any time. Lots of data are transformed among the VMs in cloud and between the cloud service providers.

There are more chances for the hijackers and internal and external intruders to handle the data packets.The IDS would be the right technique to deal with these security issues in cloud environment so here in this paper we discussed and propose the technique to overcome the security issue in order to bring availability, integrity and confidentiality. [1]There are many IDS techniques available but can only apply to the ordinary networks but it is not flexible to vital environment like cloud computing. So it is essential for researcher to focus on a proper technique to evolve a solution for the cloud security issues.

IDS can't play its role easier in the Cloud environment as it has many attacks which are predictable and unpredictable. The vast area of cloud in its virtualized environment that brings many loopholes for intruders and many challenges for the cloud service providers.The IDS is a productive solution to identify and withstand these attacks.IDS can be either hardware or software systems that capture intrusion detection, alerts the system, log the captures info.In this paper we are going to discuss about the Machine learning algorithm and they are applied to create an efficient IDS.

Then the IDS is tested with cloud dataset to find out the evaluation of the IDS.Then the IDS are compared to find out which would be the efficient cloud IDS.This paper includes Related literature review, Background of the paper with dataset details, features techniques, about the algorithm used, then Result and Discussion that has the tools used and results obtained in the work, comparing the algorithm's result and finally the conclusion with the future work.

## 2   Related Works

There are many works were done by the Researchers to enhance and reduce the data attacks. There are many importance given to intrusion detection now days due to the global usage of data through internet.This could be clear by viewing the below figure which shows the average cost due to the crime attacks happened globally Figure-1 [2]
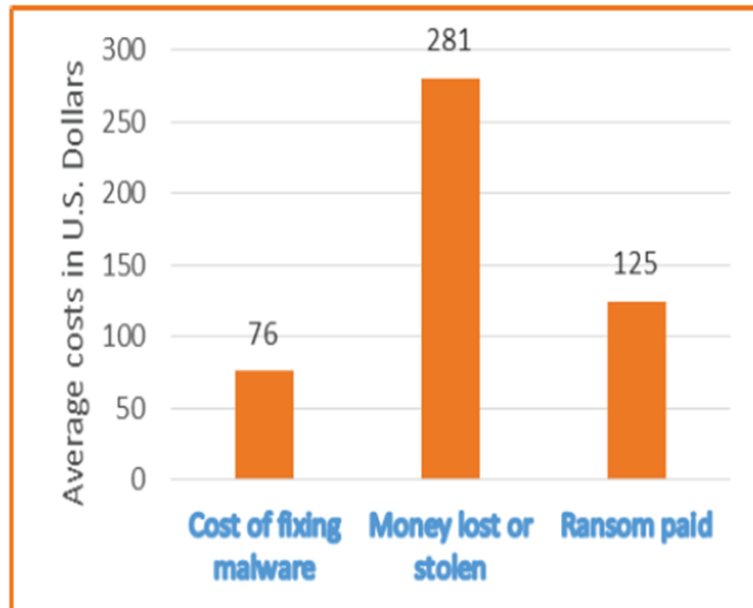
**Figure 1:** [2] Average Costs Due To Intrusion Attacks Happened

Serpil [2] used Feature elimination method with the random classifier and Deep Multilayer perceptron classifier as a deep learning classifier with the CICIDS2017 dataset. This dataset has a total of 225745 records with 80 various elements.This dataset is labelled normal with BENIGN and attacks with that particular name. Dataset is splited into two different parts training and testing the model with 80% for training and 20% for testing. According to the Methods used the whole dataset is reduced and without disturbing the accuracy and not decreasing the speed. The deep learning - DMLP gives a more useful and reduced dataset by achieving an accuracy of 89%.

Priyanka [3] has proposed a model with Neural network based ids and then it is compared with j48 algorithm and multilayered preceptron algorithm to find out the accuracy.Then the model is used with the Kyoto dataset and confusion matrix is created to get the accuracy rate. Then the data is checked with starting from 10% and it shows that even with less data the MLP Classifier performs as 0.9796% precision better than the J48 Classifier 0.9463%.

Imtiaz [4] created a model based on feature selection classifier and implemented on two datasets namely ICSX and NSL-KDD dataset. The introduced model has a subset of four features by ISCX dataset and NSL-KDD dataset takes six attributes from its 42 total attributes .By omitting the unrelated attibutes ISCX and NSL-KDD dataset obtains an accuracy of 99.7% and 99.9%.

Tara Salman [5] has created a model which used two techniques Linear Regression and Random forest. Used to detect the of four stages of classifications which reduces the error rate. So the overall accuracy for the proposed stepwise attack accuracy is attacks and find out the category .In that categorization is 93.6% and detection accuracy is 99%.UNSW was the dataset used in this paper .It has 49 features which could bring out complexity, to reduce the complexity the best-first feature selection technique is being used. It consists 93.35% whereas the single step attack is 80%. The proposed method can be used to analyse the anomaly traffic in cloud computing and in multicloud environment as well.

Jin Kim's paper has done an AI Intrusion detection model using deep neural networks[6]. The preprocessing was done to normalise the data.In this he has used the KDD Cup 99 dataset with 10% train data and full dataset as test data of 4,898,431 records.He witnessed 99% accuracy with the parameters accuracy, detection rate and false alarms.The one drawback is his work was not justified with existing work.

## 3   Background

In this section, the used dataset NSL-KDD attribute selection was performed by Random subset and classification was done by Multilayer perception classifier and j48 algorithm which are described as follows

### *3.1 NSL-KDD Dataset*

This dataset is an improved form of KDD-99. Even though it is a developed one of the KDD data set, it suffers from some of the problems said by McHugh but not proved [7]. It is mostly considered as an effective standard data set which is widely used by researchers to compare between the different intrusion detection methods. [8]This is effective than the KDD-99 dataset and most of the data are reasonable in train and test set.

The duplicate form of data is removed in the train and test NSL-KDD dataset so ease use for classifiers. No need to reduce or select small portion for an effective dataset to apply an algorithm. Already the dataset is very clear and produce the good model using the Machine learning for cloud computing security [8].NSL-KDD Dataset contains train and test records such as 125,973 and 22544 records. The whole dataset can be used fully without the necessity to sample randomly. Figure 2 and Figure 3 states the successful predicted vales of the NSL-KDD dataset. Figure 4 shows the dataset was analysed by 21 classifiers to test the labels and proved it was correctly labelled. [9]

**Figure 2:** NSL-KDD test dataset records [9]



**Figure 3:** NSL-KDD train dataset records [9]



**Figure 4:** Some data from dataset

## 3.2 Feature Selection Using Random Subset Method

Feature selection is a supervised learning algorithm to eliminate the irrelevant features and produce the accurate results by not imposing high risk of unwanted data.

Random subset selection is used to reduce cost and the measures the whole dataset fast. It randomly selects the relevant attributes from the dataset which helps in classification of high dimensional dataset. This algorithm generates q subsets which is created from the data set D each has same size M. The simple method of random subset selection is likely random selection with replacing in which each of the M elements of the subset Si is randomly taken from the size N data set separately with probability 1/N [10]. This feature selection method random subset is used to eliminate the unwanted features and it selects only the 22 attributes from the total 42 attributes shown below in the Figure 5.This brings a very good resultant model with the 22 relevant attributes alone.

The proposed method pseudo code as follows:

Step 1:    The full feature dataset is an input to normalise the data.

Step 2:   Then the Normalized dataset is preprocessed using the Random subset method to select the useful attributes

Step3:    Preprocessed dataset with selected attributes develops a model by J48 and MLP using the dataset

Step 4:   Calculate the precision, accuracy, recall and F-value. Proposed model J48 and MLP

Step 5:    Compare the precision, accuracy, recall and F-value.

Step 6:    Select the efficient model

| protocol_type | flag | src_bytes | wrong_fragment | urgent |
|---|---|---|---|---|
| num_failed_logins | urgent | num_failed_log ins | logged_in | root_shell |
| num_root | num_file_creations | num_outbound_ cmds count | srv_count | srv_serror_rate |
| srv_rerror_rate | same_srv_rate | dst_host_count | dst_host_srv_ count | dst_host_same_ srv_rate |
| dst_host_srv_diff_h ost_rate | class | | | |

**Figure 5:** Selected Attributes :22

### 3.3 Multilayer Perceptron Classifier

Multilayer Perceptron is similar as the structure of the neural network with exclusive feedforward neural networks. All nodes in the layer in a Multilayer Perceptron form a fully connected structure. [11] The three layers of MLP input neurons ,few hidden neurons and output layer. Every layer of neuron is not considered as actual layer. This paper uses the MLP to find out the efficient model for Cloud IDS. In Figure 6 The structure of MLP which has single input layer and output layer, then the hidden layer may be single or so on. In this every layer contains multiple neurons [12].
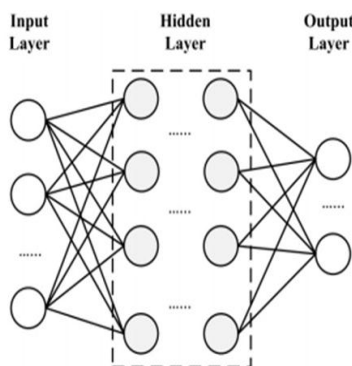


**Figure 6:** Structure of Multilayer Preceptron (MLP) [12]

### 3.4 J48 Algorithm

The J48 algorithm is used mostly to get a decision tree based results for the corresponding dataset. It can be used as a technique for classification and forecast with a representation using nodes and internodes. Split value is the essential part of the decision tree algorithm which separates the data into two parts one is root node and the other is lead node [13]. Root and internal nodes are represented as the test cases whereas the Leaf nodes are considered as class variables. This constructs a tree which gives the ranks starting from root towards the leaves. This picturizes the result and builds the perfect model. The information gain is the important to select the attributes with the highest value and selects the perfect attributes very faster which is explained in the equation (1). Entropy (S,A) A means attribute with the relation to the number of occurrence S.  The Sv v represents value of A attributes. The right side value of Entropy(S,A) is exact value of S but the left part is post value of after partitioning the S with the A attributes. The information gain is explained in equation (2). This act as the abnormal function for attribute selection, split information is found by the equation (3). The Si via Sk which are the k subsets by the occurrence S by k-number of attributes. [14]

$$\text{Gain(S,A)=Entropy(S,A)} - \sum_{v=Value} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \tag{1}$$

$$Gain\ Ratio\ (S,A) = \frac{\text{Gain(S, A)}}{\text{Split Information(S,A)}} \tag{2}$$

$$\text{Split Information(S,A)} = \sum_{j=1}^{[k]} \frac{S_i}{S} log_2 \frac{S_i}{S} \tag{3}$$

## 4    Result and Discussion

The proposed test was performed in the WEKA tool which was being used for Machine learning algorithms. This has been used by the Researchers to enhance their work without worrying about the coding part. Many new algorithms can be downloaded in weka tool which is also an Open source software created by University of Waikato in New Zealand.

### *4.1 Discussion of the MLP Model*

This Model with the Multilayer perceptron classifier has TP rate as 93.2% and FP rate as0.87% for normal .For anomaly it detects 91.3% as TP and FP rate as 0.68%.Then precision is 89% for normal and 94.5% for anomaly.The f-measure is 92.9% for anomaly and 91.1% for normal attacks. Figure 8 represents the accuracy parameters that are visualized through the graph of the proposed MLP model.

### *4.2 Discussion of the J48 Algorithm*

The TP rate is 94.7% and 98.2 %, the FP rate is0.18% and0.53%, the precision is 97.6% and 96%, the f-measure is 96.1%  and 97.1% for normal and anomaly classes. Figure 7 gives the J48 algorithm result tree structure. The time taken to test the split is 0.09 seconds and time taken to build the whole model is 5.31 seconds. This tree has 171 number of leaves and size of the tree is 294.
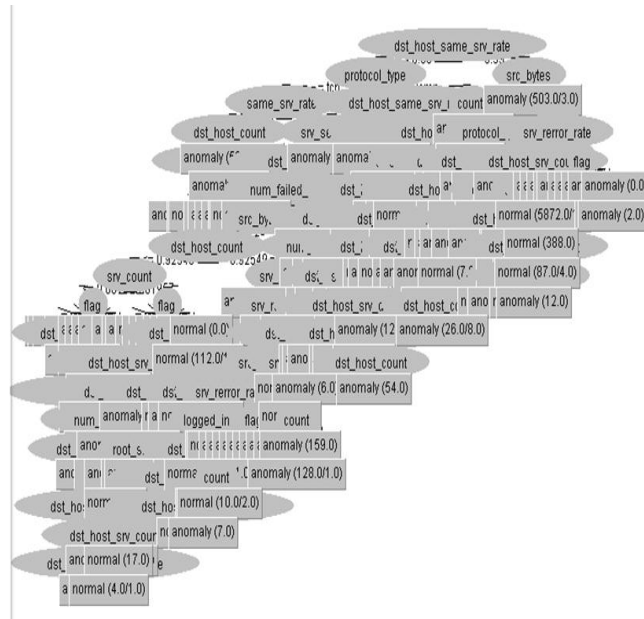
**Figure 7:** Classification tree of J48 algorithm

## *4.3 Comparison of MLP and J48 Algorithm*

Let's see the MLP and J48 algorithm rate to bring out a effective model for the cloud computing data security. The MLP FP rate is very high than the J48.Then the TP rated for MLP is less than the J48 algorithm. The Figure 8 gives the accuracy rate of MLP model and the Figure 9 explains the clear results of J48 algorithm.The MLP TP rate is 97.7 and J48 is 96.7. Precision rate for MLP is 92.2% and J48 is 96.7% so j48 precision rate is high than the MLP. These models are compared on the parameter (TP) True positive, true negatives (TNs), false positives (FPs), the false negatives (FNs) through which the Accuracy is calculated [15]. The most appropriate predictions found by the model are the accuracy. (i.e.)Accuracy is the sum of TP and TN is divided by sum of TP,TN,FP and FN.Table 1 states the comparison of MLP and J48 algorithm.The accuracy is calculated according to the equation (4) as follows

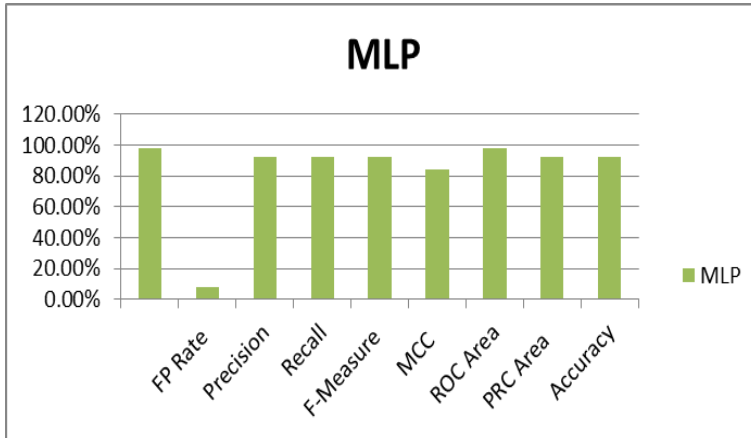$$\text{Accuracy} = \frac{TP+TN}{TN+FP+FN+TP} \qquad (4)$$

**Figure 8:** Accuracy rate of MLP model



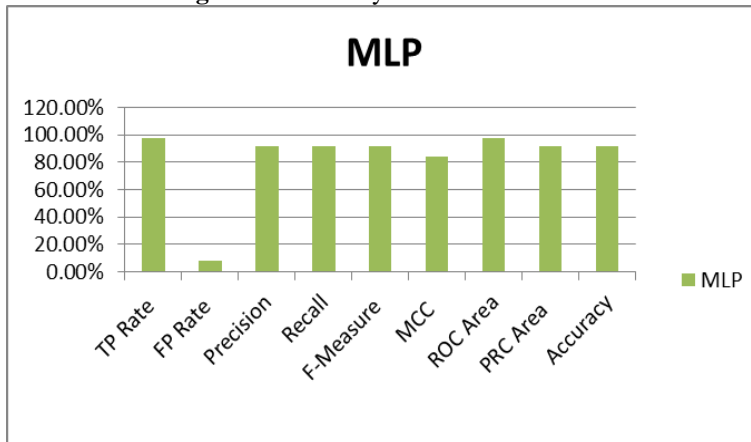**Figure 9:** Accuracy of J48 Algorithm

**Table 1:** Accuracy comparison of MLP and J48

| Class | MLP algorithm | J48 algorithm |
|---|---|---|
| **TP Rate** | 92.10% | 96.70% |
| **FP Rate** | 7.70% | 3.80% |
| **Precision** | 92.20% | 96.70% |
| **Recall** | 92.10% | 96.70% |
| **F-Measure** | 92.10% | 96.70% |
| **MCC** | 84.10% | 93.20% |
| **ROC Area** | 97.90% | 98.30% |
| **PRC Area** | 92.10% | 97.70% |
| **Accuracy** | 92.10% | 96.66% |

The confusion matrix is one of the parameter for analysis of the model which gives the comparison result of predicted class and the original class [15]. This is as follows shown in the Table 2.

**Table 2:** Confusion matrix [5]

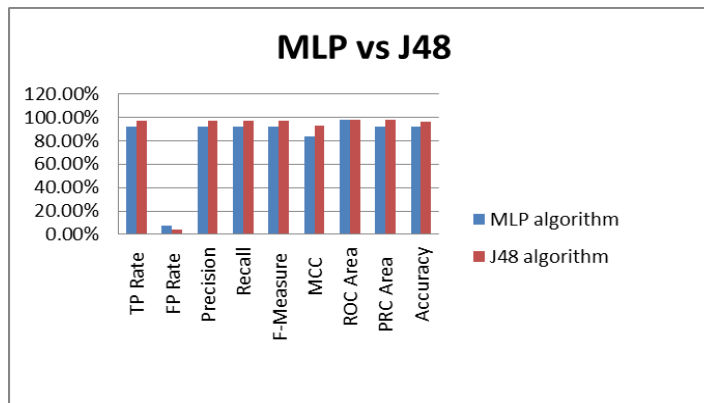| | | Predicted Class | |
|---|---|---|---|
| | | Normal | Anomaly |
| Actual Class | Normal | TN | FP |
| | Anomaly | FN | TP |



**Figure 10:** Accuracy comparison rates of MLP and J48

The above figure 10 shows the clear idea about the effective model the accuracy rate tells us that J48 is very effective than the MLP algorithm .The correctly classified instances are higher in J48 than the MLP. J48 has 3.3% incorrectly classified instances whereas MLP has 7.8%.

## 5  Conclusion

The cloud computing is fast growing field and need lots of help in the data security so this paper has shown the data security method as IDS and effective way to find out a good model to reduce the Intrusion in the cloud computing data. The IDS has various way to solve the data security here we discussed a two ids methods to bring an effective Cloud IDS. The NSL KDD data set is used to test the model and find out the improved model to solve the cloud data security issues. The Random subset is used for feature selection which played a vital role .The total 42 attributes was reduced in the NSL KDD dataset and only 22 attributes was selected using the feature selection and which brings a good result accuracy.

The 22 attributes in the J48 model brings 96.66% accuracy which is higher than the MLP rate i.e. 92.10%. The error rate is also very less in the J48 0.04% than the MLP 0.09%.The attribute selection is very important which makes the accuracy rate higher because the effective and needy attributes have been selected. So the J48 is effective model than the MLP according to the results. Therefore this model can be effective in usage in cloud computing security issues. This ML model could be applied in the Cloud computing as IDS and achieve the intrusion in the cloud environment. The next work will continue to find a very efficient model with different cloud dataset and real cloud data by creating cloud environment test bed.

## References

[1]Omar Achbarou, My Ahmed El Kiram, Outmane Bourkoukou,Salim Elbouanani,"A Multi-agent System-Based-Distributed Intrusion Detection System for a Cloud Computing" ,International Conference on Model and Data Engineering, pp.98-107, 2018.

[2]Serpil Ustebay,Zeynep Turgut,Istanbul, Turkey, "Intrusion Detection System with Recursive Feature Elimination by using Random Forest and Deep Learning Classifier", International Congress on Big Data, Deep Learning and Fighting Cyber Terrorism,IEEE, pp.71-76, 2018

[3]Joshi P.,Prasad R.,Mewada P.,Saurabh P,"A New Neural Network-Based IDS for Cloud Computing",Progress in Computing, Analytics and Networking Advances in Intelligent Systems and Computing, vol.710, pp 161-170, 2018.

[4]Imtiaz Ullah and Qusay H. Mahmoud," A Filter-based Feature Selection Model for Anomaly-based Intrusion Detection Systems", International Conference on Big Data (BIGDATA), IEEE, pp. 2151-2159,2017.

[5]Tara Salman, Deval Bhamare, Aiman Erbad, Raj Jain, Mohammed Samaka, "Machine Learning for Anomaly Detection and Categorization in Multi-cloud Environments",4th International Conference on Cyber Security and Cloud Computing, IEEE,pp.97-103,2017.

[6]J. Kim, N. Shin, S. Y. Jo, and S. H. Kim, "Method of Intrusion Detection using Deep Neural Network",International Conference on Big Data and Smart Computing, IEEE, pp. 313–316, 2017.

[7] Chetna Vaid, Harsh K Verma, "Anomaly-based IDS Implementation in Cloud Environment using BOAT Algorithm", Proceedings of 3rd International Conference on Reliability, Infocom Technologies and Optimization, IEEE,pp.1-6, 2014.

[8]Mostapha Derfouf and Mohsine Eleuldj, "Implementations of Intrusion Detection Architectures in Cloud Computing", International Conference of Cloud Computing Technologies and Applications, vol.49, pp. 100–124, 2019.

[9]M. Tavallaee, E. Bagheri, W. Lu, and A. Ghorbani, "A Detailed Analysis of the KDD CUP 99 Data Set", Symposium on Computational Intelligence for Security and Defense Applications, IEEE,pp.1-6, 2009.

[10]Padmaja Dhyaram, Dr. B. Vishnuvardhan,"Random subset feature selection for classification", International Journal of Advanced Research in Computer Science, vol.9, pp.317-319, 2018.

[11]Ansam Khraisat*, Iqbal Gondal, Peter Vamplew and Joarder Kamruzzaman, "Survey of intrusion detection systems: techniques, datasets and challenges", Cybersecurity, Springer, Vol.2, 2019.

[12]Yidong Liu , Siting Liu , Yanzhi Wang, Fabrizio Lombardi , Jie Han," A Stochastic Computational Multi-Layer Perceptron with Backward Propagation", IEEE Transactions on Computers,Vol.67, pp.1273 – 1286, 2018

[13]Shadi Aljawarneh1, Muneer Bani Yassein1, Mohammed Aljundi1,"An enhanced J48 classification algorithm for the anomaly intrusion detection systems", Journal of Cluster Computing, vol 22, pp.1-17, 2019.

[14]Yasin, Waheed , Hamidah Ibrahim, Nur Izura Udzir, Nor Asilah Wati Abdul Hamid," Intelligent Cooperative Least Recently Used Web Caching Policy based on J48 Classifier", International conference on information integration and web based applications & services, pp.262-269, 2014.

[15]Satendra kumar ,Anamika Yadav," Increasing Performance Of Intrusion Detection System Using Neural Network", International Conference on Advanced Communication Control and Computing Technologies (ICACCCT),IEEE,pp.546-550 ,2014.

## Biographies

**G Monika** is a PhD Full time Research Scholar in Vels institute of science, technology and advanced studies (VISTAS) from Department of Computer science in Chennai, TN, India. She has received her MSc (cs) from University of Madras in the Year 2013. She has 2.5 yrs of Teaching and 5 yrs Research experiences. She has greater thirst on the Research in the field of Computer science which made her to register for PhD. Her Research Area includes Cloud Computing, Security Issues, Intrusion Detection System for Data Security She also worked as teaching faculty in the APL Global School (Cambridge International School) in Tamilnadu, India. She feels that Education makes many differences in one's life so she love teaching and took it as her profession. She has published three journals in her Research area in International Journals.

**Y.Kalpana** is currently working as a professor in the department of Information Technology, Vels Institute of science Technology and Advanced Studies, Chennai, India. Her experience includes, as a teaching faculty for more than two decades in various institutions and a researcher for more than a decade. As part of the research work she has published more than 35 articles in indexed national and international journals and presented more than 20 papers in national and international conferences and authored a book titled "PHP Programming step by step Approach". Dr.Kalpana is currently guiding 8 Ph. D scholars and has produced 15 M.Phil scholars and two Ph.D scholars. She has conducted an overseas seminar "Mastering Algorithmic Techniques" at Putra Intellect International College, Malaysia. Her areas of specialization are Computer algorithms, Neural Networks and Cloud computing. She has received appreciation award for research article publication for two consecutive years(2015 and 2016) and has received best paper award in the year 2017 for her research contribution in cloud computing.