# A Secure Data Classification Model in Cloud Computing Using Machine Learning Approach

Kulwinder Kaur and Vikas Zandu

*Department of Computer Engineering & Technology*
*SVIET, Banur, Punjab, India kulwinder.kaur0117@gmail.com*
*vckz@8037@gmail.com*

## Abstract

*Cloud computing offers numerous benefits including scalability, availability and many services. But with its wide acceptance all over the globe, new risks and vulnerabilities have appeared too. Cloud computing provides facility of storing and accessing information and programs over the web without bothering the storage space on system. Storing the data on cloud eliminates one's worries about space considerations, buying new storage equipment or managing their data, rather they are able to access their data any time from any place provided they have internet access. But the rising security problems have resisted the organizations from connecting with cloud computing completely. Hence security risks have appeared as the main disadvantage of cloud computing. This paper involves the efforts to analyze the security issues and then proposes a framework to address these security issues at the authentication and storage level in cloud computing. While addressing the security issues the first and the foremost thing is to classify what data needs security and what data needn't bother with security and hence data gets classified into two classes sensitive and non-sensitive. To achieve data classification, a data classification approach based on the confidentiality of data is proposed in this paper. Following that an efficient security mechanism has to be deployed by means of encryption, authentication, and authorization or by some other method to ensure the privacy of consumer's data on cloud storage.*

*Keywords: Cloud Computing; security issues, privacy preserving, Integrity, confidentiality, availability, graphical passwords*

## 1. Introduction

Cloud Computing is recognized as a hottest technology which has a significant impact on IT field in the nearby future. In today's time Cloud computing is a fast growing field in computational research and industry. Cloud computing is the internet based computing that provides "IT resources as a service" on demand of the user following the "pay-per-use". It employs parallel processing, distributed processing, grid computing and distributed database to enhance processing, in virtualization technology and the Internet broadband technology and based network [1]. The offered Cloud Service Models can be classified as 3 categories Infrastructure as a service (IaaS), Platform as a service (PaaS), Software as a service (SaaS). The best examples of cloud service are Google App Engine, Gmail, Google Docs, Microsoft Windows Azure, and Amazon Elastic Compute Cloud (EC2).

As cloud computing helps organizations to sharpen their growth and performance. Besides this, it also hosts many users to provide access to shared resources with less effort. But security problems or threats are still a stumbling block in the success path of cloud computing. Numbers of reasons are the matter. First reason is that users and many

organizations store their data on cloud storage, so the primary focus is the data must be secure, and the data are not being lost and tampered while travelling from one place to another over the network. So it is essential that confidentiality, availability and integrity of data should be ensured. Secondly, unauthorized access where an attacker tries to be the impersonator of the legal user. [2]

Security is the number one issue when it comes to any upcoming technology and cloud computing is no exception. Cloud computing poses numerous security risks in distributed cache model. Information security is the primary developing risk for the nature of administrations that prevents the clients to embrace the cloud administrations. In distributed storage, the information is put away on the separates through two cache techniques. The previous is to encode the information and store on the server while the last is to store the information without encryption. These functions can often face confidentiality issue. The data is regularly not of the same sort and may have distinctive properties. As the consumer's data is stored on the remote servers and the consumer has no idea about its physical location, so there is always a risk of confidentiality leakage. This paper concentrates on privacy issue in cloud computing. At whatever point the information is exchanged to the cloud server it experiences a security system i.e. encryption without comprehension the level of sensitivity of the data or the data is essentially put away on cloud server without securing it. All information has diverse sensitivity levels so it is improper to store the information without comprehension its sensitivity level and security necessities. To direct the security requirements of data, we have proposed a data classification model to classify the data according to its sensitivity level and then encrypting the only data which is required to secure using an encryption technique in cloud environment.

Classification of objects is an indispensable field of research and of practical applications in numerous fields like pattern recognition and artificial intelligence, statistics, vision analysis and medicine. A very intelligent technique to secure the data would be to first classify the data into sensitive and non-sensitive data and then secure the sensitive data only. This will help to reduce the overhead in encrypting the entire data which will be exceptionally costly in connection of both time and memory. For encrypting the data many encryption techniques can be used and for classifying the data numerous classification algorithms are available in the field of data mining.

Data classification is a machine learning strategy used to predict the class of the unclassified information. Data mining uses unique instruments to grasp the unknown, legitimate patterns and relationships within the dataset. These tools are numerical calculations, factual models and prediction and evaluation of the data. Consequently, data mining consists of management, collection, prediction and analysis of the data. ML algorithms are described in to 2 categories: supervised and unsupervised. In supervised studying, classes are already outlined. For supervised studying, first, an experiment dataset is defined which belongs to distinctive classes. These lessons are competently labelled with a certain identify. Lots of the data mining algorithms are supervised studying with a detailed goal variable. In unsupervised learning classes are not effectively characterized but rather arrangement of the information is performed automatically. The unsupervised algorithm looks for similarity between two gadgets in order to find whether they are able to be characterised as forming a cluster. In simple words, in unsupervised learning, "no goal variable is identified". The classification of information within the context of confidentiality is the classification of knowledge headquartered on its sensitivity level and the have an impact on to the organization that knowledge be disclosed only licensed users. The data classification helps determine what baseline security requirements/controls are appropriate for safeguarding that data. The information is categorized into two lessons, confidential and non-confidential (non-exclusive) information. The classification of the information depends on the attributes of the

information. The values of the sensitive attributes are classified as "confidential" and values of the non-sensitive attributes are categorized as "non-confidential".

## 2. Prior Work

Tawalbeh L *et al.* [4] proposes a secure cloud computing model based on data classification. The proposed cloud model minimizes the overhead and processing time needed to secure data through using different security mechanisms with variable key sizes to provide the appropriate confidentiality level required for the data. They have stored data using three level- Basic, confidential and highly confidential level and providing different encryption algorithms at each level to secure the data. The proposed model was tested with different encryption algorithms, and the simulation results showed the reliability and efficiency of the proposed framework.

Moghaddam F. *et al.* [5] proposes a hybrid encryption model using classification indexing, attributes and time based procedures. Data classification is mainly based on attributes. A hybrid ring was used to establish the security between the rings. These securely protected rings perform the re-encryption in order to protect themselves from un-authorized access, time based, data owner request and user revocation. The result analysis shows that the hybrid ring model enhances the reliability and the efficiency of the data protection applications.

Dhamija Ankit *et al.* [6] proposes cloud architecture which ensures secure data transmission from the client's organization to the servers of the Cloud Service provider (CSP). In this, combined approach of cryptography and steganography is used because it will provide a two-way security to the data being transmitted on the network. First, the data gets converted into a coded format through the use of encryption algorithm and then this coded format data is again converted into a rough image through the use of steganography. Steganography also hides the existence of the message, thereby ensuring that the chances of data being tampered are minimal.

Singh Amritpal *et al.* [7] proposes an enhanced LSB based Steganography procedure for images bestowing better data security. It exhibits an embedding algorithm for hiding ciphered messages in nonadjacent and irregular pixel areas in edges and smooth regions of images. The edges in the cover-image are detected using improved edge detection filter. The encrypted message bits are then embedded in the least significant byte of randomly selected edge pixels and some specific LSBs of red, green, blue components respectively. Such type of steganography technique ensures least chances of suspicion about message bits hidden in the image and it gets hard to estimate the true message length by standard steganography detection methods. The Proposed approach shows better results in PSNR value and Capacity as compared to other existing techniques.

Mishra R *et al.* [8] describes a secure data transmission using LZW compression before hiding the secret data inside the image. Compression is done in order to reduce the size of the data so that it can be fit inside any media. For data hiding, encryption algorithms are used in order to provide more secure to the secret data using some keys. Now the edges can be detected using canny edge detection method and the encrypted data is embedding inside those edges. Performance analysis shows that the proposed technique is efficient and more secure with less image distortion.

Gaurav S *et al.* [9] describes all graphical methods for password authentication system and also proposed an approach which describes that first calculation has been done by server based on user entered username and according to result one set of images will be transferred on user screen, each set contains hundreds of images, and then user has to

select two images from given set, whereas server also add two images by its own to form complete password.

## 3. The Proposed Work

The cloud services are overtly offered for every kind of organization. As different companies and government/ non- government organizations prefer to store their sensitive knowledge on cloud, there arises the issue of confidentiality and integrity. Knowledge being shared and held at centralized place, makes it easy for malicious users to access, delete and alter the sensitive knowledge hence leading to confidentiality threat. In addition to this, it is important to classify the data to understand what data needs to be secured and what not. To achieve this, a data classification model has been proposed which classifies the data according to its sensitivity level further encrypting the only data which is required to secure using an encryption technique in cloud environment. Proposed work. The research involves exploring various security issues in cloud environment at with respect to three aspects confidentiality, integrity and availability and analyzes their impacts.

The research involves exploring various data classification algorithms in machine learning like KNN, Naïve Bayes and AdaBoost and analyzes their performance.

The paper proposes a secure data classification model using novel boosting supervised machine learning approach. In this, data is classified according to its sensitivity level. Then encrypting only, the data which is required to be secure using a hybrid privacy preserving based image steganography technique in cloud environment.

### 3.1 Image Sequencing Password Authentication

This password is based on the sequences of some images. It is much secure because sequence of images will change each time. This image sequencing password is use for cloud authentication purpose. Only legitimate user will allow entering in cloud, if they enter the correct sequence of image. After authentication, during access of data operations this interface will again ask the user sequence, this time images gets shuffle, based on sequence of images password will also be change.

### 3.2 Data Classification

Classifying the dataset by using Novel Boosting Technique.

The figure below shows that flowchart of the proposed boosting algorithm. In the figure, 'D' is the training data set. 'i' is a variable used for iterating the dataset. It is used to keep track of items. 'd' is a number used to divide data sets so that classifiers can be defined. 'k' is transaction set which contains all the instances. Train the classifier model using this set of transactions 'k'. Then calculate the weight of classifier's vote and update the weights of correctly classified and misclassified transactions. 'M (i) is the model built after each iteration. 'j' is the variable assigned to iterate the error rate for each model

Training set is then passed by means of the Classifier to instruct it for decision making. Replace the weight of correctly classified transaction or tuple and normalize the entire transactions such that sum of weights of the entire transactions stays identical. So Weights 'w' are assigned to the classifier's vote. Alternatively select a collection of transactions in keeping with their weights and coach weak classifiers. And generate T classifiers in each round. At definite stage, whole all the weights of votes of classifiers for class c. which class has higher weight that class would be predictive class.
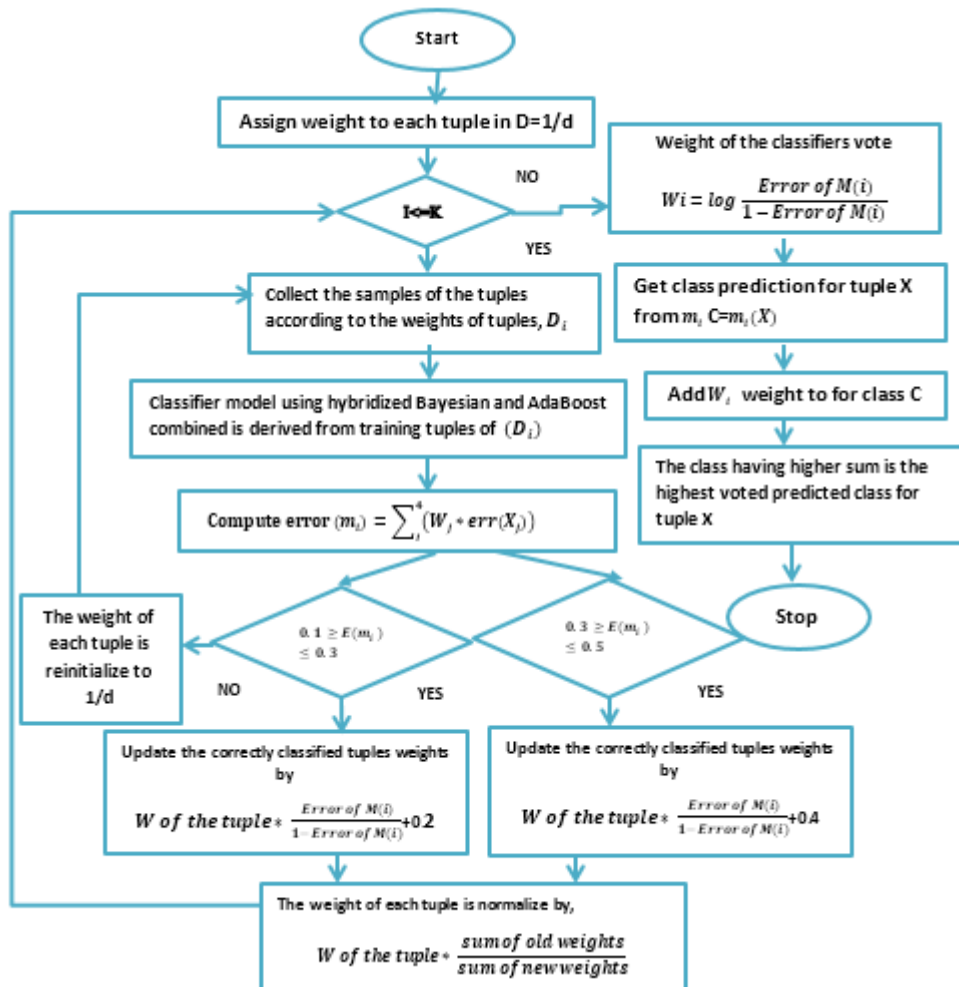
**Figure 1. Flowchart of Proposed Improved Boosting Algorithm**

### 3.3 Data hiding Architecture

To keep the sensitive data, secure from attackers on the network, a new technique has been proposed which provide the privacy to the individual's data. Instead of sending the actual data in anonymized or encrypted form on to the cloud, in this approach the actual data is not sent on to the cloud. In this an image is used to mask the sensitive data and send the randomized index values in the form of text file on to the cloud.

**Hybridized PPDM Technique based on Image Steganography**

Store all the edge pixel calculated from canny algorithm values with their positions i.e. its row and column into an array.

For example:      P (i, j, x)

Where P is the pixel with i= ith row

                 j= jth column

and x= value of the pixel on that address.

### 3.3.1 Randomization

After finding the edges and pixel values of that edges and store them into an array, randomly select the (pixel value or index value) of that array by using random function generator.

### 3.3.2 Masking

In this phase masking of one bit of message with pixel is done. This is done as follows: After pixel value is selected using randomization, LSB (Least Significant Bit) Embedding is done. This is described as:

- If the LSB bit of the image pixel values to I (i, j) which is equal to the message bit m, embed that message bit onto the LSB.
  Else
- Again find another pixel value using randomization.

$$Ls(I,j)= \quad LSB \ (I(I,j)=1) \ and \ m=0,$$
ignore that LSB
$$LSB \ (I(I,j)=0) \ and \ m=1,$$
ignore that LSB
$$LSB \ (I(I,j)=0) \ and \ m=0, \ then,$$
$$((binary \ value \ of \ Pixel \ \& \ 0xFE) \ | \ m)$$
$$LSB \ (I(I,j)=1) \ and \ m=1, \ then,$$
$$((binary \ value \ of \ Pixel \ \& \ 0xFE) \ | \ m)$$

For example: let the pixel value be 111 with binary representation 0110111. This pixel value is having LSB 1 and if message bit is also 1 then mask m with its LSB, otherwise again find another pixel value using randomization. This procedure works until all message bits of the dataset are not masked with LSB.

### 3.3.3 Reduction

Reduction is done in the way that instead of storing the pixel values onto which our message bit is masked, we would store the corresponding index values of that array into a text file and send this text file to cloud. And at our local end store that array containing edge pixels and their positions.
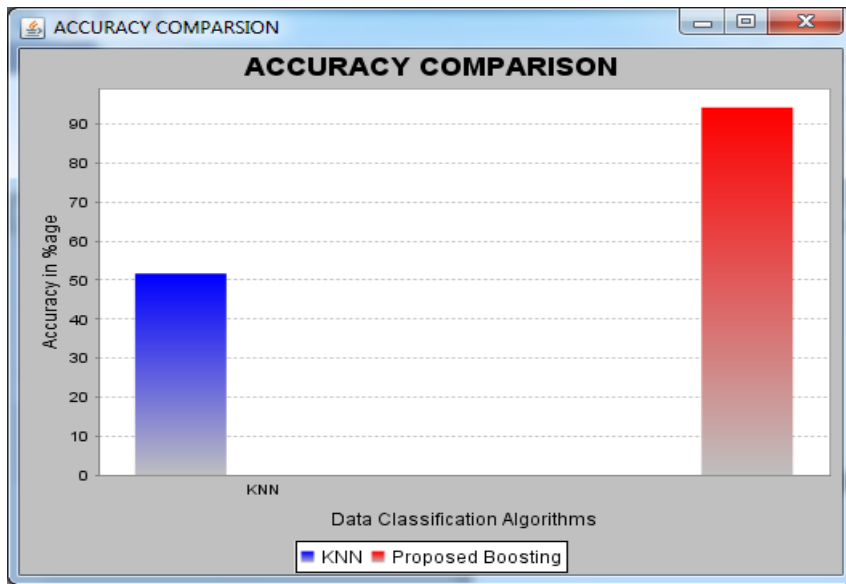
The evaluation parameters considered for evaluating the performance of the proposed system are:
a) Time taken for classifying the data
b) Accuracy of the classified data
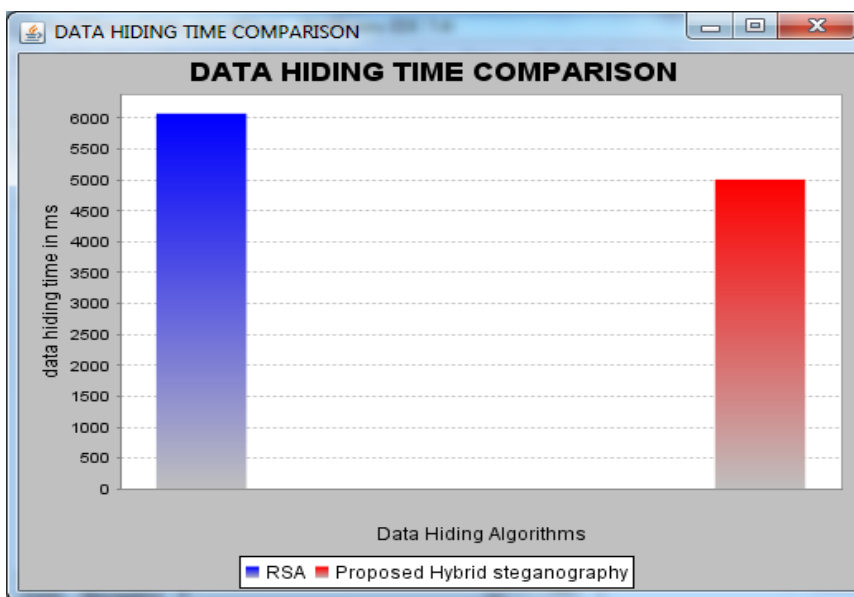c) Encryption time

## 4. Experiments and Results

### 4.1 Simulation and Analysis

The proposed methodology is implemented with the help of Cloudsim and Net beans IDE 8.0. CloudSim is the library that provides the simulation environment of cloud computing and also provide primary classes describing virtual machines, data centers, users and applications. The classification and data hiding time results have been illustrated in the following figures Figure 2 and Figure 3. And comparison between KNN with RSA and Improved Boosting with hybrid PPDM has been made in these figures.

**Figure 2. Accuracy Comparison of KNN with the Proposed Boosting Algorithm**



**Figure 3. Accuracy Comparison of Data hiding time of RSA and Proposed Privacy Preserving based steganography Algorithm**

The above figure shows the performance analysis of the proposed methodology with the previous method. It is clearly analyzed from the performance graphs that the proposed technique is better than the previous approach. Figure 2 shows the accuracy comparison of data classification algorithms KNN and proposed Improved Algorithm. KNN algorithm is having accuracy 51.7241% and improved boosting is having 94.2529% i.e. proposed algorithm has classified data more correctly and performs 42.5288% better than the KNN algorithm. Similarly, Figure 3 shows the data hiding time comparison between the proposed and previous RSA approach. Proposed Hybrid steganography technique takes 5005 milliseconds and the RSA algorithm takes 6073

milliseconds to hide the sensitive data. Therefore, in order to reduce the encryption time on cloud data is classified according to its security needs using machine learning algorithms. From the above analysis it is shown that the proposed methodology performs betters in respect to Accuracy and data hiding time.

## 5. Conclusion and Future Scope

In this research, a technique for data confidentiality in cloud environment is proposed. The focus of the research was to characterize the data taking into account the security prerequisites of the information that divides the data into sensitive and non-sensitive using improved machine learning algorithm. The fundamental contribution of this security model is data confidentiality and classification of data using machine learning classification approach. The classified confidential information is then encrypted using hybrid steganography based privacy preserving approach and is stored in the cloud server with hash key to maintain the integrity of the data while the non-confidential data is sent to the cloud environment as public data directly. Furthermore, to enhance the security at the authentication level, image sequencing passwords based on different themes has been used in order to avoid un-authorized access to the cloud environment. The proposed system has been simulated in a designed cloud simulation environment using cloud sim simulator. The results depict that the proposed technique is more relevant than storing the data without deciding the security needs of the data. Also the results show that the improved boosting technique works better than the K-NN classification technique in terms of both the accuracy and the classification time.

In future, some more security requirements can be taken in account in order to take the classification decision using machine learning algorithm, also the boosting algorithm can be further enhanced using fuzzy logics based decision rules for the classification of data according to the security needs.

## Acknowledgments

## References

[1] Application Architecture for Cloud Computing. IBM, WTHIE PAPER, November 2009.
[2] Almorsy, M., Grundy, J., & Ibrahim, A. S. (2011) "Collaboration- Based Cloud Computing Security Management Framework" IEEE conference of cloud computing, Washington (DC), pp. 364-371,2011.
[3] Song, D., E. Shi, I. Fischer and U. Shankar,(2012) "Cloud data protection for the masses" , IEEE Comput. Soc., 45(1): 39-45
[4] Lo'ai Tawalbeh, Nour S. Darwazeh, Raad S. Al-Qassas and Fahd AlDosari (2015) "A Secure Cloud Computing Model based on Data Classification", First International Workshop on Mobile Cloud Computing Systems, Management, and Security, Elsevier pp. 1153 – 1158,2015
[5] F. F. Moghaddam, M. Vala, M. Ahmadi, T. Khodadadi, and K. Madadipouya, "A reliable data protection model based on re-encryption concepts in cloud environments," 2015 IEEE 6th Control and System Graduate Research Colloquium (ICSGRC), pp. 11–16, 2015.
[6] A. Dhamija and V. Dhaka, "A novel cryptographic and steganographic approach for secure cloud data migration," 2015 International Conference Green Computing and Internet of Things (ICGCIoT), pp. 346–351, 2015.
[7] A. Singh and H. Singh, "An improved LSB based image steganography technique for RGB images," 2015 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), pp. 1–4, 2015.
[8] R. Mishra, A. Mishra, and P. Bhanodiya, "An Edge Based Image Steganography with Compression and Encryption," 2015 International Conference on Computer, Communication and Control (IC4), pp. 2–5, 2015.

[9]  S. M. Gurav, L. S. Gawade, P. K. Rane, and N. R. Khochare, "Graphical password authentication: Cloud securing scheme," 2014 International Conference on Electronic Systems, Signal Processing and Computing Technologies, pp. 479–483, 2014.Mohammed Faez Al-Jaberi and Anazida Zainal (2014)" Data Integrity and Privacy Model in Cloud Computing" International Symposium on Biometrics and Security Technologies, IEEE,pp.280-284, 2014

[10] Du, Y.; Zhang, R.; Li, M. (2013) "Research on security mechanism for cloud computing based on virtualization" Springer Telecommunication systems, Vol. 53, Issue 1, pp. 19-24, 12 July 2013

[11] Chen, D., & Zhao, H. (2012). "Data Security and Privacy Protection in cloud computing." IEEE, Hangzhou, pp. 647-651,2012.

[12] Guo, M.-H., Yen, C. T., & Hsiao, L.-L. (2012) "Authentication using graphical password in cloud", IEEE, Taipei ,pp. 177-181,2012.

[13] Abdullah, A., Hashim, F., & Al-Haddad, S. (2014) " A review of cloud security based on cryptographic mechanisms",IEEE, Kuala Lumpur, pp. 106-111,2014.

[14] Abuhussein, A., Bedi, H., & Shiva, S. (2012) "Evaluating Security and Privacy in Cloud Computing Services: A Stakeholder's Perspective", IEEE, London, pp.388-395,2012.

[15] Banirostam h., & Hedayati, A.(2013) "A Trust Based Approach for increasing Security in Cloud Computing Infrastructure" International Conference on computer modeling and simulation, IEEE, Cambridge, pp. 717-721,2013.

[16] Kang, A.N.; Barolli, L.; Park, J.H.; Jeong, Y.S. (2013) "A strengthening plan for enterprise information security based on cloud computing", Springer cluster computing, vol. 17, Issue 3, pp. 703-710, 2013

## Authors

**Kulwinder Kaur**, She is a student of Swami Vivekanand Institute of engineering and Technology, Punjab. She completed Bachelor of Technology degree in Computer Science and Engineering from Lovely Professional University in 2012. She is now pursuing Master of Technology in Computer Science at SVIET. Her major areas of interest are Cloud Computing, Data Mining and Software Engineering. She has 3 publications in reputed journals.

**Vikas Zandu**, He is an assistant Professor in SVIET. He has completed his B.Tech and M.Tech in Computer Science in 2011 and 2013 respectively. His main area of interest is networking. He has a numbers of publications in good journals and conferences.