

A Comparative Study of Clustering Algorithms for Building a Network Intrusion Detection Model

Mrutyunjaya Panda¹ Manas Ranjan Patra²

ABSTRACT

K-means is a popular clustering algorithm that requires a huge initial set to start the clustering. K-means is an unsupervised clustering method which does not guarantee convergence. Numerous improvements to K-means have been done to make its performance better. Now fuzzy set theory has been applied to many fields including data mining. Fuzzy clustering method is more precise in dealing with data simulation, and the results are easier to be understood and used. Therefore, research into fuzzy clustering method for knowledge is significant not only to theory, but also to application. Expectation Maximization is a statistical technique for maximum likelihood estimation using mixture models. It searches for a local maxima and generally converges very well. In this paper, we propose some clustering algorithms such as K-Means, Fuzzy c-Means, and EM (Expectation and maximization) for network intrusion detection. We have used KDDCup'1999 data set for our experimentation. The simulation results show that EM algorithm is a more statistically formalized method, which accounts for partial membership in classes. It has better convergence properties and is in general preferred to K-Means and Fuzzy c-Means algorithms in building a network intrusion detection model.

Keywords : Intrusion Detection, K-Means, Fuzzy c-Means, EM, MF plot, ROC, log likelihood.

1. INTRODUCTION

An intrusion detection system (IDS) is a component of the information security framework. Its main goal is to differentiate between normal activities of the system and behaviour that can be classified as suspicious or intrusive [1]. The goal of intrusion detection is to build a system which would automatically scan network activity and detect such intrusion attacks. Once an attack is detected, the system administrator can be informed who can take appropriate action to deal with the intrusion.

IDS can be host-based (HIDS), network-based (NIDS) or a combination of both types (Hybrid Intrusion Detection System). HIDS usually observes logs or system -calls on a single host, while a NIDS typically monitors traffic flows and network packets on a network segment, and thus observes multiple hosts simultaneously. Generally, one deal with very large volumes of network data, and thus it is difficult and tiresome to classify them manually in order to detect a possible intrusion. One can obtain labelled data by simulating intrusions, but this will be limited only to the set of known attacks. Therefore, new types of attacks that may occur in future cannot be handled, if those were not part of the training data. Even with manual classification, we are still limited to identifying only the known (at classification time) types of attacks, thus restricting our detection system to identifying only those types.

¹Dept. of ECE, GIET, Gunupur , Orissa , India. E-mail: mrutyunjaya.2007@rediffmail.com

²Dept. of Computer Science , Berhampur University, Orissa, India. E-mail : mrpatra12@gmail.com

To solve these difficulties, we need a technique for detecting intrusions when our training data is unlabeled, as well as for detecting new and un-known types of intrusions. A method that offers promise in this task is anomaly detection. Anomaly detection detects anomalies in the data (i.e. data instances in the data that deviate from normal or regular ones). It also allows us to detect new types of intrusions, because these new types will, by assumption, be deviations from the normal network usage.

It is very difficult, if not impossible, to detect malicious intent of someone who is authorized to use the network and who uses it in a seemingly legitimate way. For example, there is probably no highly reliable way to know whether someone who correctly logged into a system is the intended user of that system, or if the password was stolen.

Under these assumptions we built a system which created clusters from its input data, then automatically labelled clusters as containing either normal or anomalous data instances, and finally used these clusters to classify network data instances as either normal or anomalous. Both the training and testing was done using 10% KDDCup'99 data [2], which is a very popular and widely used intrusion attack dataset.

A popular algorithm is the K-means where, based on a given number of clusters, the algorithm iterates to find best clusters for the objects. Most clustering techniques assume a well defined distinction between the clusters so that each pattern can only belong to one cluster at a time. This supposition can neglect the natural ability of objects existing in multiple clusters. For this reason and with the aid of fuzzy logic, fuzzy clustering can be employed to overcome the weakness. The membership of a pattern in a given cluster can vary between 0 and 1.

In this model a data object belongs to the cluster where it has the highest membership value. In this paper, we aim to propose a fuzzy c-means clustering technique which is capable of clustering the most appropriate number of clusters based on objective function. This, as the name implies, draws the fuzzy boundary, thereby proving efficient when compared with that of its counterpart. Another approach is to use the Expectation Maximization algorithm (EM). The Expectation Maximization algorithm may be a very efficient technique for estimating class conditional probability density Functions (PDF) in both univariate and multivariate cases [3]. This paper discusses about the K-Means, Fuzzy c-Means and EM clustering algorithms and compares their suitability in building an efficient network intrusion detection model.

The rest of the paper is organised as follows. In section 2, we discuss Clustering Methods; followed by Data Clustering algorithms in section 3. In Section 4, various cluster evaluation schemes are discussed. Section 5 describes about the experimental set-up and results obtained. Some discussion is made in Section 6. Finally, Section 7 provides some related works followed by conclusion in Section 8.

2. CLUSTERING METHODS

Clustering may be found under different names in different contexts, such as unsupervised learning (in pattern recognition), numerical taxonomy (in biology, ecology), typology (in social sciences) and partition (in graph theory) [4]. By definition, "cluster analysis is the art of finding groups in data", or from Wikipedia [5], "clustering is the classification of similar objects into different groups, or more precisely, the partitioning of a data into subsets (clusters), so that the data in each subset (ideally) share some common trait-often proximity

according to some defined distance measure. Clustering is a challenging field of research as it can be used as a stand-alone tool to gain insight into the distribution of data, to observe the characteristics of each cluster, and to focus on a particular set of clusters for further analysis. Alternatively, cluster analysis serves as a pre-processing step for other algorithms, such as classification which would then operate on detected clusters.

The process of grouping a set of physical or abstract objects into classes of similar objects is called clustering. A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters. By clustering, one can identify dense and sparse regions and therefore, discover overall distribution patterns and interesting correlations among data attributes. Clustering does not rely on predefined classes and class labelled training examples. For this reason, clustering is a form of learning by observation.

In intrusion detection, an object is a single observation of audit data and/or network packets after the values from selected features have been extracted. Hence, values from selected features, and one observation, define one object (or vector). If we have values from n number of features, the vector (or object plot) fits into an n -dimensional coordinate system (Euclidean space R^n).

In order to derive the objective function and other relevant mathematics for fuzzy c -means and the remaining of its variations, it is better to see the same for the hard (crisp) partitioning technique, so that we may be able to understand the difference between the two approaches. (If we look into these issues all of them appears to be objective functional minimization problems. If the constraints are relaxed we get the possibilistic partition scheme. So, the clustering

algorithm is nothing but a minimization problem which may be constrained or unconstrained.)

2.1. Hard Partitioning

These kind of methods are based on classical set theory and defines the presence or absence of a data point in a partition subset on strict logic, that is the object either belong to a subset or not. So, such kind of methods divides a dataset strictly into disjoint subsets.

Conventional clustering algorithms find a "hard partition" of a given data set based on certain criteria that evaluate the goodness of a partition. By hard partition we mean that each datum belongs to exactly one cluster of the partition. More formally, we can define the concept of "hard partition" as follows:

1) Let X be a data set of data, and x_i be an element of X . A partition $p = \{C_1, C_2, \dots, C_j\}$ of X is "hard" if and only if

$$i) \forall x_i \in X \quad \exists C_j \in P \text{ such that } x_i \in C_j$$

$$ii) \forall x_i \in X \quad x_i \in C_j \Rightarrow x_i \notin C_k \text{ where } k \neq j, C_j \in P.$$

The first condition in the definition assures that the partition covers all data points in X ; the second condition assures that all the clusters in partition are mutually exclusive.

2) Let x be a data set of data and x_i be an element of X . A partition $p = \{C_1, C_2, \dots, C_j\}$ of X is "soft" if and only if the following condition holds:

$$i) \forall x_i \in X \quad \forall C_j \in P \quad \text{for } 0 \leq \mu_c(x_i) \leq 1;$$

$$ii) \forall x_i \in X \quad \forall C_j \in P \text{ such that } \mu_{c_j}(x_i) > 0.$$

2.2. Soft Partitioning

A soft clustering algorithm partitions a given data set not an input space. Theoretically speaking, a soft partition not necessarily a fuzzy partition, since the input space can be larger than the dataset. In practice however most

soft clustering algorithms do generate a soft partition that also forms the fuzzy partition.

A type of soft clustering of special interest is one that ensures the membership degree of a point x in all clusters adding up to one, i.e.

$$\sum_j \mu_{c_j}(x_i) = 1, \forall x_i \in X \text{----- (1)}$$

A soft partition that satisfies this additional condition is called a constrained soft partition. The fuzzy c-means algorithm produces a constrained soft partition. The fuzzy c-means algorithm is best known algorithm that produces constrained soft partition.

The biggest drawback of a hard partitioning is the concept that it either includes a data point in a partition or strictly excludes it; there is no other chance for the data elements to be part of more than one partition at the same time. However, in natural clusters it is always the case that some of the data elements partially belong to one set and partially to one or more other sets. In order to overcome this limitation, the notion of fuzzy partitioning was introduced [6].

3. DATA CLUSTERING ALGORITHMS

The following are the algorithms used for clustering the datasets:

- K-means Algorithm
- Fuzzy c-means Algorithm.
- Expectation-Maximization (EM) Algorithm.

3.1. K-means Clustering

The k-means clustering is a classical clustering algorithm. After an initial random assignment of example to k clusters, the centres of clusters are computed and the examples are assigned to the clusters with the closest centres. The process is repeated until the cluster centres do not significantly change. Once the cluster assignment is fixed, the mean distance of an example to cluster

centres is used as the score. Using the K-means clustering algorithm, different clusters were specified and generated for each output class [7].

K-means clustering is a well known Data Mining algorithm that has been used in an attempt to detect anomalous user behaviour, as well as unusual behaviour in network traffic. There are two problems that are inherent to k-means clustering algorithms. The first is determining the initial partition and the second is determining the optimal number of clusters [8]. In figure 1 depicted K-means algorithm.

K-MEANS ALGORITHM

Input: The number of clusters K and a dataset for intrusion detection

Output: A set of K -clusters that minimizes the squared-error criterion.

Algorithm:

1. Initialize K clusters (randomly select K elements from the data)
2. While cluster structure changes, repeat from 2.
3. Determine the cluster to which source data belongs

Use Euclidean distance formula.

Add element to cluster with min

(distance (x_i, y_j)).

4. Calculate the means of the clusters.
5. Change cluster centroids to means obtained using Step 3.

Figure 1 : K-means Clustering

As the algorithm iterates through the training data, each cluster's architecture is updated. In updating clusters, elements are removed from one cluster to another. The updating of clusters cause the values of the centroids to change, which is a reflection of the current cluster elements. Once there are no changes to any cluster, the training of the K-Means algorithm is complete.

At the end of the K-Means training, the K cluster centroids are created and the algorithm is ready for

classifying traffic. For each element to be clustered, the cluster centroids with the minimal Euclidean distance from the element will be the cluster for which the element will be a member. After training, the cluster centroids remains the same, like the SOM (Self organise Map) can be useful for anomaly detection tool that requires the input to remain static. The k-Means algorithm may take a large number of iterations through dense data sets before it can converge to produce the optimal set of centroids. This can be inefficient on large data sets due to its unbounded convergence of cluster centroids.

3.2. Fuzzy c-Means (FCM) Clustering

Fuzzy c-Means (FCM) algorithm, also known as fuzzy ISODATA, was introduced by Bezdek [9] as extension to Dunn's [10] algorithm to generate fuzzy sets for every observed feature. The fuzzy c-means clustering algorithm is based on the minimization of an objective function called c-means functional.

Fuzzy c-means algorithm is one of the well known relational clustering algorithms. It partitions the sample data for each explanatory (input) variable into a number of clusters. These clusters have "fuzzy" boundaries, in the sense that each data value belongs to each cluster to some degree or other. Membership is not certain, or "crisp". Having decided upon the number of such clusters to be used, some procedure is then needed to location their centres (or more generally, mid-points) and to determine the associated membership functions and the degree of membership for the data points.

Fuzzy clustering methods allow for uncertainty in the cluster assignments. FCM is an iterative algorithm to find cluster centres (centroids) that minimize a dissimilarity function. Rather than partitioning the data into a collection of distinct sets by fuzzy partitioning, the membership

matrix (U) is randomly initialized according to equation 2.

$$\sum_{i=1}^c u_{ij} = 1, \forall j=1,2,\dots,n. \quad (2)$$

The dissimilarity function (or more generally the objective function), which is used in FCM in given equation 3.

$$J(U, c_1, c_2, \dots, c_c) = \sum_{i=1}^c J_i = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d_{ij}^2 \quad (3)$$

Where, U_{ij} is between 0 and 1;

c_i is the centroids of cluster I;

d_{ij} is the Euclidean Distance between i^{th} . Centroids c_i and j^{th} . Data point.

$m \in [1, \infty]$ is a weighting exponent. There is no prescribed manner for choosing the exponent parameter, "m". In practice, $m=2$ is common choice, which is equivalent to normalizing the coefficients linearly to make their sum equal to 1. When m is close to 1, then the cluster centre closest to the point is given much larger weight than the others and the algorithm is similar to k-Means.

To reach a minimum of dissimilarity function there are two conditions. These are given in (4) and (5).

$$c_i = \frac{\sum_{j=1}^n u_{ij}^m x_j}{\sum_{j=1}^n u_{ij}^m} \quad (4)$$

$$u_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{d_{ij}}{d_{kj}} \right)^{\frac{2}{m-1}}} \quad (5)$$

This algorithm determines the following steps in Figure2.

FCM ALGORITHM:

Input: n data objects, number of clusters

Output: membership value of each object in each cluster

Algorithm:

1. Select the initial location for the cluster centres
2. Generate a new partition of the data by assigning each data point to its closest centre.
3. Calculate the membership value of each object in each cluster.
4. Calculate new cluster centers as the centroids of the clusters.
5. If the cluster partition is stable then stop, otherwise go to step2 above.

Figure 2 : Fuzzy C-Means Clustering Algorithm

By iteratively updating the cluster centres and the membership grades for each data point, FCM iteratively moves the cluster centres to the "right" location within a data -set. FCM does not ensure that it converges to an optimal solution, because the cluster centers are randomly initialised. Though, the performance depends on initial centroids, there are two ways as described below for a robust approach in this regard.

- 1) Using an algorithm to determine all of the centroids.
- 2) Run FCM several times each starting with different initial centroids.

More mathematical details about the objective function based clustering algorithms can be found in [11].

3.3. EM Clustering

Finite mixture distributions provide a flexible and mathematical-based approach to the modelling and clustering of data observed on random phenomena. We focus here on the use of the normal mixture models, which can be used to cluster continuous data and to estimate the underlying density function. These mixture

models can be fitted by maximum likelihood via the EM (Expectation-maximisation) algorithm. The main assumption is that data points are generated by, first randomly picking a model j with probability $\tau_j, j=1: k$, and, second, by drawing a point x from a corresponding distribution. The area around the mean of each (supposedly unimodal) distribution constitutes a natural cluster. So, we associate the cluster with the corresponding distribution's parameters such as mean, variance, etc. Each data point carries not only its (observable) attributes, but also a (hidden) cluster ID (class in pattern recognition). Each point x is assumed to belong to one cluster, and we can estimate the probabilities of the assignment $\Pr(C_j|x)$ to j^{th} model. The overall likelihood of the training data is its probability to be drawn from a given mixture model

$$L(X|C) = \prod_{i=1:N, j=1:k} \tau_j \Pr(x_i | C_j) \quad (6)$$

Log-likelihood $\log(L(X|C))$ serves as an objective function, which gives rise to the Expectation-Maximisation (EM) method. EM is a two step iterative optimization. Step (E) estimates probabilities $\Pr(x | C_j)$, which is equivalent to a soft (fuzzy) reassignment. Step (M) finds an approximation to a mixture model, given current soft assignments. This boils down to finding mixture model parameters that maximise log-likelihood. The process continues until log-likelihood converges is achieved.

Because the mixture model has clear probabilistic foundation, the determination of the most suitable number of clusters k becomes a more tractable task. From a data mining perspective, excessive parameter set causes over fitting, while from a probabilistic perspective, number of parameters can be addressed within the Bayesian framework. More details can be found in [12].

One attractive feature of adopting mixture models with elliptically symmetric components such as the normal or t densities is that it implies clustering is invariant under affine transformations of the data (that is, under operations relating to changes in location, scale, and rotation of the data). Thus, the clustering process does not depend on irrelevant factors such as the units of measurement or the orientation of the clusters in space.

As with k -means, the EM algorithm is only guaranteed to converge to a local maximum, not the global one and so the procedure should be repeated a number of times with different initial guesses for the parameter values. In this case, however, the log-likelihood figure can be used to directly compare the final configurations obtained and so the user just has to choose the largest of the local maxima.

The standard EM algorithm generates a series of parameter estimates, where represents the reaching of the convergence criterion, through the following steps, as shown in figure3.

The major disadvantages for EM algorithm are the sensitivity to the selection of initial parameters, the effect of a singular covariance matrix, the possibility of convergence to a local optimum, and the slow convergence rate [13]. Variants of EM for addressing these problems are discussed in [13] and [14]. A valuable theoretical note is the relation between the EM algorithm and the K -means algorithm. Celeux and Govaert proved that classification EM (CEM) algorithm under a spherical Gaussian mixture is equivalent to the K -means algorithm [15, 16].

4. CLUSTERING EVALUATION SCHEMES

The result of the cluster analysis can be evaluated with the help of any of the following, in order to measure

EM ALGORITHM:

- 1) Initialize θ^0 and set $t=0$;
- 2) **e-step:** Compute the expectation of the complete data log-likelihood

$$Q(\theta, \theta^t) = E[\log p(x^0, x^m | \theta) | x^0, \theta^t];$$
- 3) **m-step:** Select a new parameter estimate that maximizes the Q -function,

$$\theta^{t+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^t);$$
- 4) Increase $t=t+1$; repeat steps 2)-3) until the convergence condition is satisfied.

Figure 3 : Expectation-Maximization (EM) Algorithm

their effectiveness in building an network intrusion detection model.

4.1. Objective Function/Log Likelihood Criteria

Usually for clustering, there are two kinds of measures of cluster "goodness" or quality [17]. One type of measure allows us to compare different sets of clusters without reference to external knowledge and is called an internal quality measure. This type of measure uses the "overall similarity" which is based on the pair wise similarity of documents in a cluster. The second type allows us to evaluate how well the clustering is working by comparing the groups produced by the clustering techniques to known classes. This type of measure is called an external quality measure. One external measure is entropy, which provides a measure of "goodness" for un-nested clusters or for the clusters at one level of a hierarchical clustering. Another external measure is the F -measure, which is more oriented towards measuring the effectiveness of a hierarchical clustering. For EM algorithm, we may use "log likelihood" to measure the "overall similarity", since in EM, in each iteration, we optimize the log likelihood of expected parameters.

The algorithm terminates when a formula that measures cluster quality no longer shows significant increases. One measure of cluster quality is the likelihood that the data came from the dataset determined by the clustering. The likelihood computation is simply the multiplication of the sum of the probabilities for each of the instances. With two clusters A and B containing instances x_1, x_2, \dots, x_n , where $P_A = P_B = 0.5$, the computation is:

$$[0.5P(x_1|A) + 0.5P(x_1|B)] + [0.5P(x_2|A) + 0.5P(x_2|B)] \dots [0.5P(x_n|A) + 0.5P(x_n|B)] \quad (7)$$

The log likelihood score is the logarithm of the likelihood measure defined above in equation (7). There is no upper bound on this value; however, larger scores represent clustering of higher quality.

In case of FCM, the minimum value of the objective function is equal to twice the logarithm of the likelihood. If the likelihood is smaller than 1, the log will be negative. The likelihood, in simple normal problems is a sum of squares. If that sum is >1 , then $-2\log$ likelihood will be negative. The absolute value of the objective function is meaningless, as a likelihood is only defined up to an arbitrary proportionality constraints. Only differences between objective functions of nested models are meaningful.

4.2. Using Number of Clusters

EM allows us to choose the number of clusters to be formed. As an alternative, it can also be possible to instruct EM to determine a best number of clusters. The algorithm will converge to an optimal clustering; however, the optimization may not be global.

4.3. Choice of Seeds

The clustering results can vary based on random seed selection. Some seeds can result in poor convergence

rate or convergence to sub-optimal clustering. So, the choice of seed is of paramount importance in good clustering evaluation in order to build an efficient intrusion detection model.

5. EXPERIMENTAL SETUP AND RESULTS

In this experiment, we have used a standard dataset, the raw data used by the KDD Cup 1999 intrusion detection contest [2]. However, in our experiment; we have used 10% KDD Cup'99 datasets. This database includes a variety of intrusions simulated in a military network environment that is common benchmark for evaluation of intrusion detection techniques. This data set consists of 65525 data instances, with 21 training attack types, each of which is a vector of extracted feature values from a connection record obtained from the raw network data gathered during the simulated intrusion and is labelled as either normal or a certain attack type. The distribution of attacks in the KDD Cup'99 dataset is highly unbalanced. Some attacks are represented with only a few examples, e.g. the phf and ftp_write attacks, whereas the Smurf and Neptune attacks cover many records. In general, the distribution of attacks is dominated by probes and DoS attacks.

We carried out the experiments on 2.8GHz Pentium IV processor, 512 MB RAM running Windows XP system. Weka tool [18] was used for performing the K-Means and EM clustering experimentation, whereas Fuzzy Logic Toolbox [19] of MATLAB 7.0 was used for fuzzy c-Means clustering. The K-Means algorithm finds k clusters by choosing k data points at random as initial Cluster centers. Each data point is then assigned to the cluster with center that is closest to that point. Each cluster center is then replaced by the mean of all the data points that have been assigned to that cluster. This process is iterated until no data point is reassigned to a different

cluster. The simulation result of K-Means algorithm is shown in figure4, in terms of Receiver Operating characteristics, which is a measure between Detection rate and False Alarm Rate. This is a good indicator of performance, since it measures what percentage of intrusions the system is able to detect and how many incorrect classifications are made in the process.

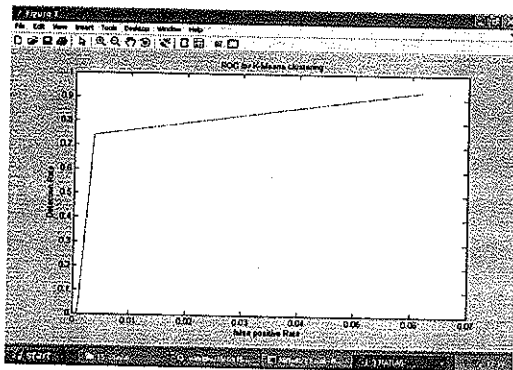


Figure 4 : K-Means Clustering (ROC)

In practice, the number of classes is not always known beforehand. There is no general theoretical solution to find the optimal number of clusters for any given dataset. We choose $k=5$ for the experimentation of the FCM. The simulation results after using FCM are shown in figures 5, and 6. In figure 7, the shape of the membership function for selected values of the fuzzification factor ($m=2$) and cluster number=5 is shown. It can also be seen from these figures that, we are able to group the data by using the objective functions based fuzzy c-means clustering approach. Finally, the relationships of the objective function with the number of iterations are obtained in figure 8.

In case of EM Clustering, the aim is to select the number of clusters automatically by maximizing the logarithm of the likelihood of future data, estimated using cross-validation. Beginning with one cluster, it continues to add clusters until the estimated log-likelihood decreases. The results are illustrated in figures 9, 10, 11, and 12

based on various criteria as illustrated in following section.

6. DISCUSSION

The K-means clustering algorithms are the simplest methods of clustering data. The K-means algorithm uses a set of unlabeled feature vectors and classifies them into k classes, where k is given by the user. From the set of feature vectors k of them are randomly selected as initial seeds. The feature vectors are assigned to the closest seeds depending on its distance from it. The mean of features belonging to a class is taken as the new center. The features are reassigned; this process is repeated until convergence. The effectiveness of the K-Means clustering is shown in figure 4.

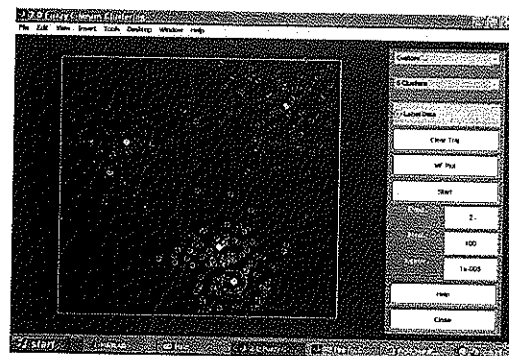


Figure 5 : Five Clusters Of Data After Using FCM

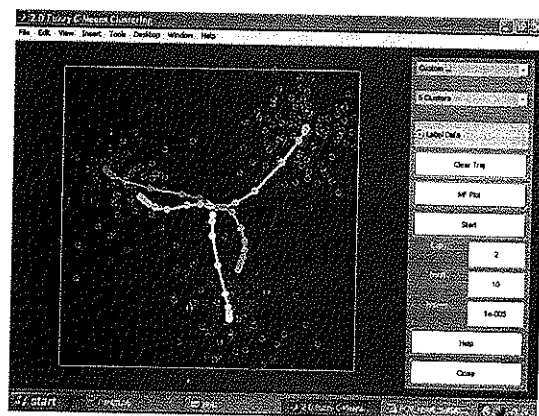


Figure 6: 2-D Fuzzy c-Means Clustering

The fuzzy based clustering methods had shown tremendous achievements in areas of image processing and pattern recognition. The fuzzy c-means is a good choice for circular and spherical clusters, but if the orientation of natural clusters is not spherical, then the algorithm leads to among almost wrong clusters. Another drawback of the algorithm is that it imposes equal size clusters on the data set which is again a deviation from the natural clusters. The performance of any fuzzy based clustering method is the best when the number of clusters is known Apriori. But most of the time, it is not the case and so researchers have devised a number of methods known as cluster validation indices to evaluate the clusters formed [20, 21]. The simulation results of the FCM are shown in figures 5, 6, 7, and 8.

The reason we are using EM is to fit the data better, so that clusters are compact and far from other clusters, since we initially estimate the parameters and iterate to find the maximum likelihood for those parameters. EM uses the Maximum likelihood, in which it assumes that the parameters are fixed; the best estimate of their value is defined to be the one that maximizes the probability of obtaining the samples actually observed. In most cases the observed data could be the samples that are used for training. Based on the number of Seeds and number of Clusters used, the effectiveness of the EM clustering is shown in figures 9 and 10 respectively. Comparison of K-Means and EM clustering is made with the help of ROC is shown in figure11, which shows that EM is a better choice in comparison to K-Means clustering algorithm. Finally, the comparison of K-Means, FCM, and EM algorithms are made with the help of their objective function in figure12. This result shows that, the EM algorithm is more suitable in comparison to all others in order to build a efficient intrusion detection model.

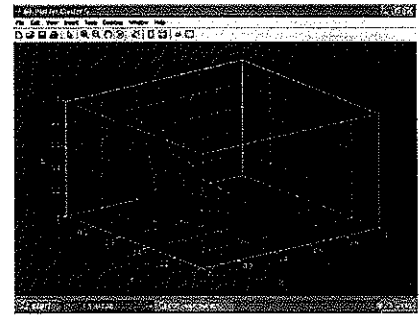


Figure 7 : MF (Membership Function) Plot for FCM

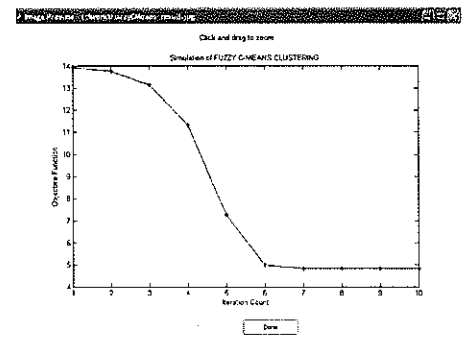


Figure 8 : Simulation of FCM Based on Iteration Count

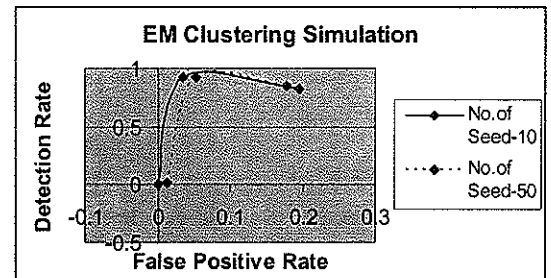


Figure 9 : ROC Comparison Based on Number of Seeds

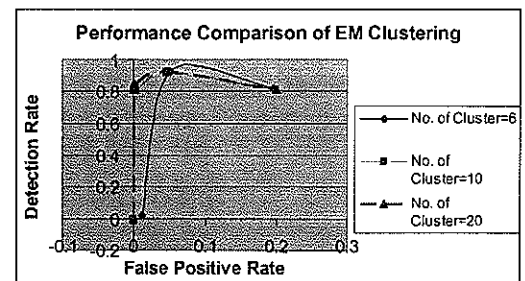


Figure 10 : ROC Comparison Based on Number of Clusters

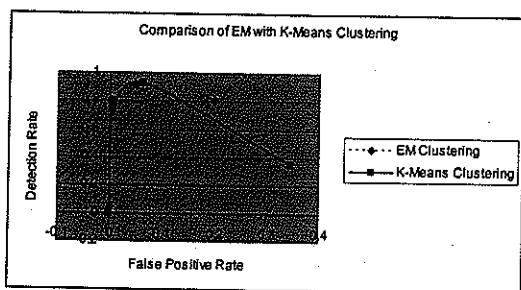


Figure 11 : Comparison of K-Means and EM Clustering Based on ROC

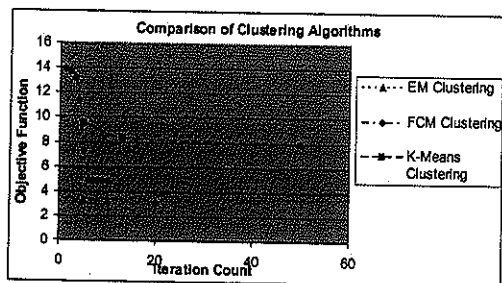


Figure 12 : Objective Function Based Comparison of Clustering Algorithms

7. RELATED WORK

In [22], a speed up technique for image data was proposed. In this method, FCM convergence is obtained by using a data reduction method. Data reduction is done by quantization and speed-up by aggregating similar examples, which were then represented by a single weighted exemplar. The objective function of the FCM algorithm was modified to take into account the weights of the exemplars. However, the presence of similar examples might not be common in all data sets. They showed that it performs well on image data sets. However, the above algorithm does not address the issue of clustering large or very large datasets under the constraints of limited memory.

Recently in [23], a sampling based method has been proposed for extending fuzzy and probabilistic clustering to large or very large data sets. The approach is based on progressive sampling, which can handle the non-image

data. However, the termination criteria for progressive sampling could be complicated as it depends upon the features of the data sets.

In [24], two methods of scaling EM to large data sets have been proposed by reducing time spent in E-step. In the first method, which is referred to as incremental EM, data is partitioned into blocks and then incrementally updating the log-likelihoods. In the second method, lazy EM, at scheduled iterations the algorithm performs partial E and M steps on a subset of data. The methods used to scale EM may not generalize to FCM as they are different algorithms with different objective functions.

8. CONCLUSION

The applications of fuzzy based methods in all fields of engineering and sciences have shown far reaching results and their applications in intrusion detection are also optimistic. In this paper, we have discussed the objective function based fuzzy c-means clustering in detail and their application in detecting anomaly based network intrusions. Fuzzy clustering leads to information granulation in terms of fuzzy sets or fuzzy relations. Membership grades are important indicators of the typicality of patterns or their borderline character. The advantage of using fuzzy logic is that it allows one to represent concepts where objects can fall into more than one category (or from another point of view- it allows representation of overlapping categories). The results obtained in this paper show that FCM works very efficiently in obtaining compact well separated clusters to detect network intrusions. Though we have already seen many examples of successful application of cluster analysis, there still remain many open problems due to the existence of many uncertain factors. These problems have already attracted and will continue to attract intensive efforts from broad discipline.

REFERENCES

- [1] J.Allen, A. Christie, W.Fithen, J.McHugh, J.pickel, and E.Stoner, "State of the practice of Intrusion Detection Technologies", CMU/SEI-99-TR-028, Carnegie Mellon Software Engg. Institute, 2000.
- [2] KDDCup'1999 dataset, <http://kdd.ics.uci.edu/databases/kddcup'99/kddcup99.html>.
- [3] Wael Abd-Elmaged, Aly El-osery, Chistopher E. Smith, "Non-Parametric Expectation Maximization: A Learning Automata Approach", IEEE Conference on Systems, Man and Cybernatics, Washington DC,2003.
- [4] S.theodoridis and K.koutroubas, "pattern Recognition", Academic Press, 1999.
- [5] Wikipedia-Cluster Analysis, http://en.wikipedia.org/wiki/cluster_analysis.
- [6] Johan Zeb Shah and anomie bt Salim, "Fuzzy clustering algorithms and their application to chemical datasets", in Proc. Of the post graduate Annual Research seminar, PP.36-40, 2005.
- [7] Zhengxim Chen, "Data Mining and Uncertain Reasoning-An integrated approach", Willey, 2001.
- [8] Witcha Chimphee, et.al, "Un-supervised clustering methods for identifying Rare Events in Anomaly detection", in Proc. Of World Academy of Science, Engg. and Tech (PWASET), Vol.8, PP.253-258, Oct2005.
- [9] J.Bezkek, "pattern Recognition with fuzzy objective function algorithms", Plenum Press, USA, 1981.
- [10] S.Albayrak and Fatih Amasyali, "Fuzzy C-Means clustering on medical diagnostic systems", International XII Turkish Symposium on Artificial Intelligence and Neural Networks, TAINN-2003.
- [11] Wit old Pedrycz, "Knowledge Based Clustering", John Willey&sons Inc., 2005.ISBN:0-471-46966-1.
- [12] Xindong Wu, Vipin Kumar, J.ross Quinlan,et.al., "Top Ten Algorithms in Data Mining", Knowledge Information System(2008),14, PP.1-37. Published online: 4 December 2007.
- [13] G.McLachlan and T.krishnan, "The EM algorithms and extensions", New York: Wiley, 1997.
- [14] M.Figueiredo and A. Jain, "UnSupervised learning of finite mixture models", IEEE Transaction on Pattern Anal. Mach. Intell., Vol.24, No.3, PP.381-396, March2002.
- [15] G.celeux and G.Govaert, "A classification EM algorithm for Clustering and two stochastic versions", Comut.Statist. Data Anal., Vol.14, PP.315-332, 1992.
- [16] Rui Xu and Donald Wunsch II, "Survey of Clustering Algorithms", IEEE Transaction on Neural Networks, Vol.16,No.3, PP.645-678, May 2005.
- [17] M.Steinbach, G.Karypis and V.Kumar, "A Comparison of Document Clustering Techniques", in KDD workshop on Text Mining, 2000.
- [18] Weka Tool, University of Waikato, New Zealand.
- [19] MATLAB 7.0, Math Works, statistical Toolbox, www.mathworks.org.
- [20] V.Maulik, S.Bandopadhyay, "Performance evaluation of some clustering Algorithms and Validity indices", IEEE transaction on Pattern Analysis and Machine Intelligence, Vol.24, No.12, PP.1650-1654, Dec.2002.
- [21] Steven Eschrich, Jingwei Ke, Lawrence o. Hall and Dmitry B. Goldgof, "Fast accurate fuzzy Clustering

through data reduction", IEEE Transaction on Fuzzy Systems, Vol.11, 2, PP.262-270, 2003.

[22] Richard J. Hathaway and James C. Bezdek, "Extending fuzzy and probabilistic clustering to very large datasets", Journal of Computational Statistics and Data Analysis, Vol.51, Issue 1, PP.215-234, Nov.2006.

[23] Bo Thiesson, Christopher Meek and David Hackerman, "Accelerating EM for large Database", Machine Learning Journal, V.45, PP.279-299, 2001.

[24] A K Jain and R C Dubes, "Algorithm for Clustering Data", Prentice Hall, Engle Wood cliffs, NJ, USA, 1988.

Author's Biography



Mrutyunjaya Panda holds a Master Degree in Engineering and is presently working as an Assistant Professor in the Department of Electronics & Tele Comm. Engg., Gandhi Institute of Engineering and Technology, Gunupur, India. He has 11 years of teaching experience. Currently, he is pursuing Doctoral research in Computer Science in Berhampur University, Orissa,

India. He has about 20 publications to his credit. His research interests include Data Mining, Network Security, Intrusion Detection and Soft Computing. He is also a reviewer of an International journal (IJSDIA). He is a life member of CSI (India), IETE (India), and ISTE (India), ACEEE, MIAENG (Hong Kong).



Dr. Manas Ranjan Patra Holds a Ph.D. Degree in Computer Science from the Central University of Hyderabad and is presently working as a Reader in the Department of Computer Science, Berhampur University, India. He has worked in the International Institute for Software Technology, Macao as a United Nations Fellow during 2000. He has 22 years of experience in teaching and research in different areas of Computer Science. He has about 70 international and national publications to his credit. His research areas include Software Engineering, Data Mining, Intrusion Detection, Artificial Intelligence, and e-business. He has presented papers and chaired technical sessions in many International conferences. He is a member of number of professional bodies. He has executed visiting assignments to many Institutions and Universities.