# Enhanced Expectation–Maximization Clustering through Gaussian Mixture Models

**S.Nagarjuna Reddy, S.Sai Satyanarayana Reddy, M.Babu Reddy**

*Abstract: Clustering is the most important task in data mining. For the intelligent clustering is also the part of the machine learning. Various existing systems are introduced for better clustering. In the past decade so many existing clustering algorithms are introduced to perform better results. These algorithms work on extracting the patterns from the unsupervised decision tree. Binary cuckoo search based decision tree is adopted with Expectation–Maximization (EM) Clustering through Gaussian Mixture Models (GMM) to improve performance of the clustering. Here we are using numerical data set, mushroom and MIST dataset to extract patterns using clustering. The performance will be estimated in terms of various measures like sensitivity, specificity, and accuracy.*

*Keywords: EM-GMM, K-Means, Mushroom, MIST*

## I. INTRODUCTION

The analysis of clustering is one of the learning techniques that establish the foundation of a clever information analyzed process. It is helpful for the investigation between connections among a gathering of examples, by sorting out into homogeneous clusters. It is called unsupervised learning on the grounds that from the earlier marking of certain examples is accessible to use in classifying others and deducing the group structure of the entire information [1]. Density is one of the proportions. A high intra-network implies a decent clustering game plan in light of the fact that the occasions assembled inside a similar group are very subject to one another. Each case in the informational index can be spoken to utilizing a similar arrangement of properties. The properties are clear cut.

To invigorate speculation from a given informational index, a learning framework requires making presumptions about the theory to be educated. These suspicions are called as predispositions. Every learning technique utilizes a few inclinations, it responds well in certain spaces where these are suitable while it performs ineffectively in different areas [2].

The issue with clustering strategies is that the understanding of the clusters might be troublesome. In some cases such as the algorithms always maps the clusters even when there are no clusters in the data. a) Supervised Learning: In controlled adjusting furthermore called facilitate data mining the components under investigation are secluded t into two get-togethers: illustrative elements and (at any rate one) subordinate variables.

b) Unsupervised Learning: In this learning, all of the elements are managed in a similar way, there is no capability contrast to

**S.Nagarjuna Reddy**, Research Scholor, Dept of CSE, JNTUK, Kakinada, India.
**S. Sai Satyanarayana Redy**, Dept of CSE, Vardhaman College of engineering, Hyderabad, India
**M.Babu Reddy**, Dept of CSE, Krishna University, Machilipatnam, India

the name undirected data mining still there is some target to achieve. This goal might be as data diminishment as general or progressively specific like bundling. The dividing line between unsupervised learning and oversaw learning is comparable that perceives discriminate examination from gathering examination.
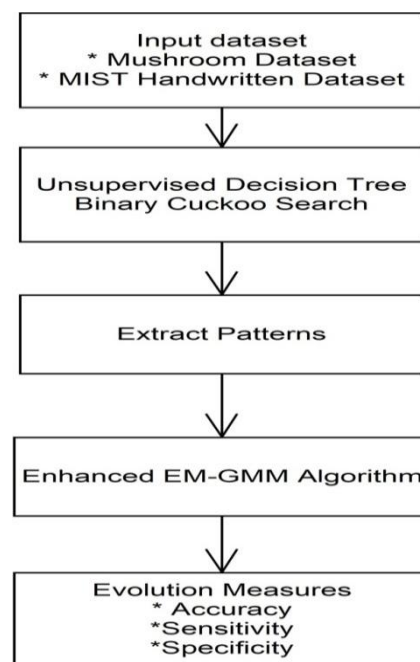


**FIGURE: 1 STEPS OF EM-GMM**

## II. LITERATURE SURVEY

This chapter explains the various clustering algorithms. To reduce the dimensionality of a dataset is used to make the clustering process easy, for this we are using K-means clustering. This is also not fit for huge datasets, PCA is the analysis method adopted to K-means then. For better results, PCA and linear transformation is utilized for data reduction and the basic centroid is calculated and applies the k-means algorithm.

An HLVM used for Data Visualization. Analysts present an HVA which permits the total informational index to be envisioned at the top dimension, with groups and sub clusters of information focuses pictured at more profound dimensions. As per Zhao et al., [3] Fast and top notch grouping calculations are significant for perusing a lot of data. Utilize hierarchal bunching methods like parcelled and agglomerative clustering. Through exploratory assessment, it was seen that divided algorithm is in every case superior to agglomerative calculations because of

relative less computational necessities and bunching execution. In [1] they have introduced another method for clustering calculation it was named as obliged agglomerative calculations which is blend of both apportioned and agglomerative methodology.

## III. DATA SET DESCRIPTION

**Mushroom Dataset: [9]**
It is recorded with 8124 instances, where each instance is represented with 22 attributes with missing values. Each instance describes whether mushroom is poisonous or edible.

**The MNIST Dataset: [10]**
It contains set of handwritten digits, recorded with 60,000 records, and 10,000 records are selected for test set. NIST is a large set will be act super set for MIST. The digit images are normalized b size and generates as a fixed-size image

## IV. EXPECTATION–MAXIMIZATION (EM) CLUSTERING THROUGH GAUSSIAN MIXTURE MODELS (GMM)

The major drawback in the K-means having follows the mean based approach to calculate cluster centre. It can't deal with this on the grounds that the group mean values are extremely near one another. K-Means likewise bombs in situations where the groups are not round, again because of utilizing the mean as a centre of the cluster. The proposed method using GMM gives us more adaptability compared to K-Means. With this we accept to facilitate the information values are Gaussian appropriated; Instead of utilizing the mean it will provide better prohibitive supposition [4]. To depict the shape of the bunches; we have two parameters: Standard deviation and mean. By picking a precedent among the two measurements, the bunches can be formed into a one form of circular shape. In this way, every cluster is assigned by Gaussian distribution. For each bunch we need to discover the parameters of the Gaussian for that we utilize an efficient optimization approach named as Expectation– Maximization (EM). Investigate the realistic graph as an outline of the Gaussians able to fit into the clusters. At that point, we can continue to follow some procedure of proposed clustering [5].

1. We start by choosing the cluster size (number like K-Means does) and haphazardly instating the Gaussian circulation parameters for each group. It can endeavour to give a decent observation time to the underlying parameters by investigating the information as well. Despite the fact that note, seeing that it found in the realistic over, this isn't completely vital as the Gaussians begin our as exceptionally poor however are immediately streamlined. 2. Gaussian distributions are assigned to every cluster for each bunch, figure the likelihood that every datum point has an associated with a specific bunch. The nearest point will be selected as the Gaussian's middle; it has more probable it has a place with that group. This should bode well since with Gaussian dissemination we are expecting that a large portion of the information placed nearer to the focal point of the group.

3. From the obtained probabilities, we process other technique to amplify the probabilities of data points up to what extent it is suitable to a specific cluster. To position a point in the cluster we are computing weighted total of the parameters; based on the probabilities it will be placed in the specific bunch. To clarify this in a graphical way we can

investigate the realistic above, specifically the yellow cluster for instance. The dissemination begins off arbitrarily on the primary emphasis; however, we observe that the majority of the colour yellow focuses are to one side of that conveyance. When we compute a total weight count by the probabilities, despite the fact that there are a few points close to the inside, a large portion of them are on the right. Along these lines normally the dispersion's mean is moved nearer to that arrangement of focuses. We can likewise observe that the majority of the focuses are "upper appropriate to base left". In this manner the parameter (Standard Deviation) changes to make an oval that is increasingly fitted to these focuses, so it leads to boost the total weight count by probabilities.

4. The above steps are rehashed repeatedly until union, where the Gaussian distributions can't be modified much from emphasis to cycle. It contains extremely two notable points of interest in utilizing GMMs. Initially, GMMs are much more adaptable as far as group covariance compared to K-Means; the cluster formed into an oval shape because of parameter called standard deviation; as opposed to being circle shape. Unlike K-Means; In GMM in which each group's covariance along all measurements approaches 0. Also, GMMs follow probabilities; they can form various groups per information point. So if an information point is amidst two covering groups, we can just characterize its group by concluding it has a place X-percent with group 1 and Y-percent to group 2. I.e. GMMs bolster blended enrolment.

**Binary Cuckoo Search-Decision Tree**
Here, we use a double cuckoo seek calculation is utilized to remove designs from a gathering of unsupervised choice trees made through a various levelled strategy. In different investigates in information mining writing has presented various calculations for grouping. In our examination, we present an ideal choice tree system for information bunching dependent on Binary Cuckoo Search Algorithm. The as good as ever clarification is supplanting the dominant part helpful clarification in the home. The accompanying portrayal framework is chosen by Cuckoo Search calculation: Every egg in a home symbolizes a clarification, and a Cuckoo egg symbolizes a novel clarification [7]. The aim is to use the most likely improved egg to re-establish a not all that immense egg from the all the eggs in the homes. The methodology of grouping uses the following steps.

• Only One egg is laid by the cuckoo once. An egg will be placed in an arbitrarily picked home.

• The size of homes to access is fixed, and patch up with the huge calibre of eggs will persist to the following ages.

• If host winged animal found the cuckoo egg; then host animal can discard egg or surrender the home to some other one, and fabricate and moved to new home.

**Processing Steps for EM-GMM**
The idea of the proposed algorithm originates from the GMM. The process adopted by GMM strategy is one approach to progress the thickness of a prearranged arrangement of test information displayed as a component of the likelihood thickness of a solitary thickness estimation technique with numerous Gaussian likelihood thickness capacities to demonstrate the circulation of the information. When all is said in done, to get the assessed parameters of each Gaussian mix segment whenever given an example informational index of the log-probability of the information, the most

extreme is controlled by the EM calculation to gauge the ideal replica. Basically, the proposed grouping technique utilizes the accompanying calculation:

Information: Cluster figure m, a catalogue, ceasing resilience. Yield: A lot of m-groups with a weight augment with probabilities.

1. Expectation advance: For every catalogue record p, process the participation likelihood of every record p in each bunch b = 1,…, m.

2. Maximization advance: Blend model parameter has to update frequently.

3. Stopping criteria: If ceasing criteria are satisfied then stop, else set a = a +1 and then continue from (1).

Expository techniques accessible to accomplish likelihood dispersion parameters, more likely than estimation of the variable are not given. The iterative EM calculation utilizes an irregular variable and, in the end, is a frequent strategy to locate the ideal weighted parameters of the concealed appropriation work from the given information, when the information is inadequate or has some mislaid qualities.

To estimate the parameters of distribution is maximum likelihood estimation.

$$\mathcal{N}(p|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(p-\sigma)^4}{\sigma^2}} \quad Equation - 1$$

Where $\mu, \sigma$ are the mean and covariance,

After calculating our posterior all we need to do is get an estimate of the parameters of each Gaussian defined by the equations below and then evaluate the log-likelihood. These two steps are then repeated until convergence.

$$\mu_a^{new} = \frac{1}{N_a} \sum_{n}^{N} = 1 \, \gamma(z_{na}) a_n \quad Equation - 2$$

Equation 2 considered as mean of Gaussian

$$\sum_{a}^{new} = \frac{1}{N_a} \sum_{n}^{N} = 1 \, \gamma(z_{na})(p_n - \mu_a^{new})(p_n - \mu_a^{new})^T \quad Equation - 3$$

Equation 3 considered as covariance of Gaussian

**Performance Evolution**

The performance of the proposed system can be estimated using False Alaram Rate, False Negative Rate, Sensitivity, Specificity and Accuracy. The count values are computed from the confusion matrix to calculate the above measures.

**False Alarm Rate (FAR)**

The probabilities of cases where an image was falsely classified and rejected the null hypothesise of a specific test. It is also called false positive rate

$$FAR = \frac{FP}{FP + TN}$$

**False Negative Rate (FNR)**

The probabilities of cases where it was recognized to anomalous images, but in fact it did.

$$FNR = \frac{FN}{FN + TN}$$

**Sensitivity**

The proportion of actual positives that live properly known is that the measure of

the sensitivity. It relates to the flexibility of the take a look at to spot positive results.

$$Sesitivity = \frac{No. of \, TP}{No. of \, TP + No. of \, TN}$$

**Specificity**

The proportion of negatives that area unit properly known is that the measure of the specificity. It relates to the power of the check to spot negative results.

$$Specificity = \frac{No. of \, TN}{No. of \, TN + No. of \, FP}$$

## V. RESULTS AND DISCUSSION

Different algorithms can be utilized to isolate information of a comparable sort. In contrast to the arrangement algorithm, clustering has a place with the gathering of unsupervised algorithms. Among these algorithms, the broadly utilized is K-Means and our proposed approach. Linear regression analysis, by utilizing a direct blend of free factors, is a factual system used to anticipate the likelihood of event of an occasion, for example, its probability.[8] However, on the off chance that a lot of information are arranged utilizing strategic relapse investigation just, it is beyond the realm of imagination to expect to ensure the precision of the results. In this paper, Binary Cuckoo Search-Decision Tree subsequent to clustering the exploratory information with the K-means and proposed algorithm to take care of the troublesome issue of guaranteeing the exactness of the acquired outcomes.

The resultant clusters and K-means were functional for assessing the quality of **Parkinson's, spam, based** datasets results are shown in [1]. In this paper, the EM-GMM utilized the mushroom dataset and the MNIST database of handwritten digit datasets is used. Here, we are going to compare the performance of the various discussed approaches based on proposed algorithm.

**Results of Segmentation Evaluation**

The Mushroom dataset is recorded with 8124 instances, where each instance is represented with 22 attributes with missing values. Each instance describes whether mushroom is poisonous or edible.CSV file format.



**Figure: 2, Sample Mushroom dataset**

- Import data, clean, perform exploratory analysis and test the best model fit.
- Take a deep look into important variables and further classify what feature of the variable helps in identifying whether the mushroom is edible or not.

| Cluster-1 | Accuracy | Sensitivity | Specificity | Total records in cluster |
|---|---|---|---|---|
| K-Means Clustering | 0.9777886 | 0.95 | 0.94 | 0.47 |
| EM-GMM | 0.9976876 | 0.99 | 0.96 | 0.51 |

**Table 1: In cluster-1 all the mushrooms are belongs to edible.**

| Cluster-2 | Accuracy | Sensitivity | Specificity | No of records in this cluster |
|---|---|---|---|---|
| K-Means Clustering | 0.9677886 | 0.95 | 0.93 | 0.47 |
| EM-GMM | 0.9976876 | 0.99 | 0.96 | 0.48 |

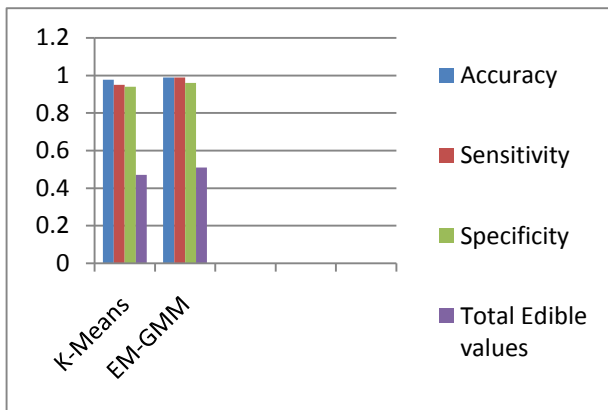**Table 2: In cluster-2 all the mushrooms are belongs to poisonous.**



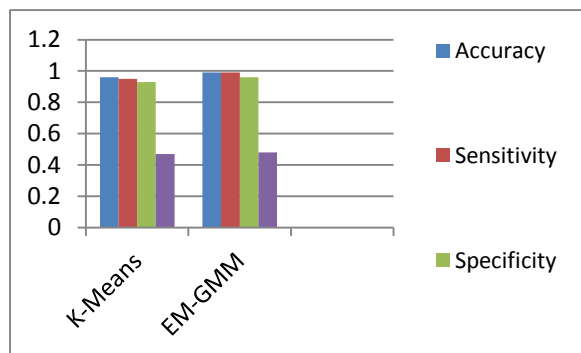**Figure: 3 Performance graph representation for edible (Cluster-1) records in mushroom dataset.**



**Figure: 4 Performance graph representation for poisonous (Cluster-2) records in mushroom dataset.**
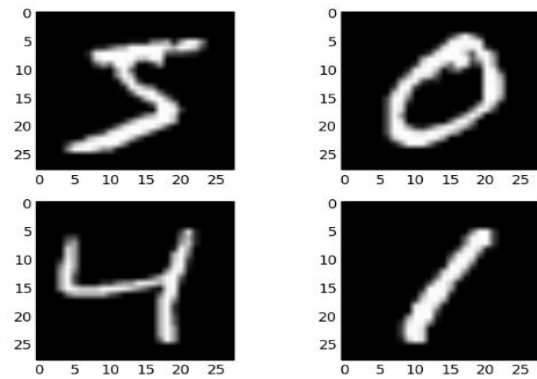


**Figure: 5 Examples of MNIST dataset**

Evolution results on MIST hand written dataset:

| | Accuracy | Sensitivity | Specificity | No of records in this cluster |
|---|---|---|---|---|
| Cluster-1 (Number 0) | 0.99 | 0.96 | 0.95 | 0.3212 |
| Cluster- 2 (Number 1) | 0.99 | 0.97 | 0.95 | 0.3342 |
| Cluster-3 (Number 2) | 0.99 | 0.98 | 0.95 | 0.4321 |
| Cluster- 4 (Number 3) | 0.99 | 0.94 | 0.95 | 0.5433 |
| Cluster-5 (Number 4) | 0.99 | 0.98 | 0.95 | 0.2343 |
| Cluster-6 (Number 5) | 0.99 | 0.9932 | 0.95 | 0.3543 |
| Cluster-7 (Number 6) | 0.99 | 0.9932 | 0.95 | 0.2123 |
| Cluster-8 (Number 7) | 0.99 | 0.9598 | 0.95 | 0.3421 |
| Cluster-9 (Number 8) | 0.99 | 0.9987 | 0.95 | 0.4321 |
| Cluster-10 (Number 9) | 0.99 | 0.96 | 0.95 | 0.4242 |

**Table 3: Performance of the clusters for handwritten dataset 0-9**

The Enhanced component of EM-GMM outcomes show of clustering massive data are group presented by our exposure work. The appraisal between to same like k-means scheme clustering technique to EM-GMM in this evaluation our workpractise will give bepart of high accuracy log values for clustering scheme data. The robin accuracy for the work that exist K-Means is approx.

96% accuracy like for Mushroom dataset low compared to this result to our proposed EM-GMM clustering technique it gives 99% accuracy for Mushroom dataset. However, in the MIST dataset of work, the existing partial technique gives 96% of set exactness which is squat when we associate to this outcome to our anticipated EM_GMM clustering procedure it will give 99% accuracy. From these work, proposed Enhanced EM-GMM based clustering sample give enhanced accuracy endings. Consequently our extension work displays that it is worth save for the clustering.

## VI. CONCLUSION

Extension to this paper, we present small Enhanced biological pattern based on via some clustering methods based return set of patterns. To base scheme performance of the EM-GMM set binary search of cuckoo algorithm in this wild tree will produce in a hide hierarchical procedure. These path algorithms are biased to extract patterns decision tree. The biased decision tree learning school technique will abstract the patterns of previous given data set. The previous data to be clustered based on technique EM-GMM algorithm. The calculated results measures sensitivity, and accurateness were evaluated by our projected method. The bold effectiveness clustering schemes scientific approve data is very high by bestowing very good taster outcomes and also the collected data. From the final outcomes, we have proved some EM-GMM work outperforms the existing algorithms by the simplified very good precision of 99% for both mushroom and MIST database dataset. Therefore by consuming this procedure, our projected Enhanced EM-GMM based practical clustering technique. In the future, we will use another enhanced clustering method for data clustering to make more compatible for any type of datasets fit for the extension to that of mushroom.

## REFERENCES

1.  Bilmes JA. A gentle tutorial of the EM algorithm and its application to parameter estimation for gaussian mixture and hidden Markov models. Berkeley: CA: International Computer Science Institute; 1998.
2.  Jung YG, Lim MJ, Choi YJ. Using PCA and random projections to compare preference of performance . 2012;7(2):469–472.
3.  Y. Zhao, G. Karypis, Evaluation of Hierarchical Clustering Algorithms for Document Datasets, ACM 2002, pp. 515-524.
4.  S. Guha, R. Rastogi, K. Shim, CURE: An efficient Clustering Algorithm for Large Databases, ACM 1998,pp. 73-84.
5.  u ,M.C,Chou,C.H.,"A Modified Version of the KMeans Algorithm with a Distance Based on Cluster Symmetry", IEEE Transactions on Pattern Analysis and Machine, 23 (6), Aug 7, 2002.
6.  Kanungo, T., David M., Piatko C.D.,Silverman,R.,Angela Y.Wu, An Efficient k-Means Clustering Algorithm: Analysis and Implementation, IEEE Transactions on Pattern Analysis and Machine Intelligence, 24,881-892,2002.
7.  Nagarjuna Reddy Seelam,Sai Satyanaryana Reddy Seelam,Babu Reddy Mukkala,Optimal Decision Tree Based Unsupervised Learning Method for Data Clustering,Received: September 21, 2016, IJIES, Vol.10, No.2, 2017
8.  Bilmes JA. A gentle tutorial of the EM algorithm and its application to parameter estimation for gaussian mixture and hidden Markov models. Berkeley: CA: International Computer Science Institute; 1998.
9.  G. H. Lincoff (Pres.), New York: Alfred A. Knopf (1981). The Audubon Society Field Guide to North American Mushrooms [https://archive.ics.uci.edu/ml/datasets/mushroom]
10. Yann LeCun, Corinna Cortes, Chris Burges,(1999), THE MNIST DATABASE of handwritten digits,[ http://yann.lecun.com/exdb/mnist/