

# Enhanced Classification of Incomplete Pattern Using Hierarchical Clustering

Gajanan B. Malekar<sup>1</sup>, Prof. Gurudev Sawarkar<sup>2</sup>

Department of Computer Science & Engineering

<sup>1</sup>ME Student, V. M. Institute Of Engineering & Technology, Nagpur

<sup>2</sup>Assistant Professor, V. M. Institute Of Engineering & Technology, Nagpur

**Abstract-**More often than not values are absent in database, which ought to be managed. Missing qualities are occurred in light of the way that, the data segment individual did not know the right regard or frustration of sensors or leave the space cleanse. The course of action of missing regarded lacking case is a trying errand in machine learning approach. Divided data is not proper for classification handle. Exactly when insufficient cases are masterminded using prototype values, the last class for comparable illustrations may have distinctive results that are variable yields. We can't describe specific class for specific cases. The structure makes a wrong result which also realizes contrasting effects. So to oversee such kind of lacking data, system executes prototype-based credal classification (PCC) strategy. The PCC procedure is intertwined with Hierarchical clustering and evidential reasoning methodology to give correct, time and memory profitable outcomes. This procedure readies the examples and perceives the class prototype. This will be useful for recognizing the missing qualities. By then in the wake of getting each and every missing worth, credal procedure is use for classification. The trial occurs exhibit that the enhanced type of PCC performs better similar to time and memory viability.

**Keywords-**Belief functions, hierarchical clustering, credal classification, evidential reasoning, missing data.

## I. INTRODUCTION

Data mining can be considered as a technique to find fitting information from broad datasets and recognizing outlines. Such cases are further useful for classification handle. The key helpfulness of the data mining method is to find supportive information inside dataset and change over it into an informed association for quite a while later.

In an expansive part of the classification issue, some quality fields of the dissent are empty. There are distinctive clarification for the void attributes including disillusionment of sensors, mixed up qualities field by customer, eventually didn't get the essentialness of field so customer leave that field fumes et cetera. There is a need to find the capable system to portray the challenge which has missing attribute values.

Diverse classification procedures are available in writing to deal with the classification of insufficient cases. Some framework empties the missing regarded cases and just uses complete plans for the classification methodology. In any case, sooner or later inadequate cases contain basic information in like manner this system is not a real course of action. Also this technique is material exactly when insufficient data is under 5% of whole data. Disregarding the divided data may lessen the quality and execution of classification figuring. Next system is simply to fill the missing qualities anyway it is furthermore monotonous process. This paper is based on the classification of divided examples. On the off chance that the missing qualities relate a great deal of data then departure of the data components may come to fruition into a more noticeable loss of the required true blue data. So this paper generally concentrates on the classification of lacking cases.

Progressive Clustering produces a gathering chain of significance or a tree-sub tree structure. Each cluster center point has relatives. Essential gatherings are joined or spilt according to the top down or base up approach. This procedure helps in finding of data at different levels of tree.

Exactly when lacking illustrations are requested using prototype values, the last class for comparative cases may have different results that are variable yields, with the objective that we can't portray specific class for specific cases. While learning prototype regard using ordinary calculation may prompts to inefficient memory and time in results. To overcome these issues, proposed system executes evidential reasoning to process specific class for specific case and Hierarchical Clustering to figure the prototype, which yields successful results with respect to time and memory.

## II. RELATED WORK

Pedro J.Gracia-Laencina, Jose-Luis Sancho-Gomez [2] proposed Pattern classification with accomplishment used as a piece of a couple issue territories, as biometric affirmation, record classification or investigation. Missing information is a standard burden that illustration affirmation frameworks are compelled to change once assurance

certifiable assignments classification. Machine taking in methods and courses outside from associated number-crunching learning theory are most importantly inspected and used in the space.

The essential goal of review is to investigate missing information, plan classification, and to study and take a gander at a portion of the unmistakable courses used for missing data organization.

SatishGajawada and DurgaToshniwal [3] showed a paper; Real application dataset could have missing/cleanse values however a couple classification frameworks require whole datasets. In any case if the articles with divided illustration are in tremendous number then the rest complete inquiries inside dataset square measure slightest. The measure of complete things may be distorted by considering the figured question as aggregate challenge and misuse the registered question for additional tallies by the conceivable complete articles. In this paper they have used the Kmeans and K Nearest neighbor values for the attribution. This technique is associated on clinical datasets from UCI Machine Learning Repository. Cristobal J. Carmona, Julian Luengo proposed a paper [4] Subgroup disclosure may be an expressive data get ready strategy that goes for getting enchanting standards through coordinated learning. All things considered, there are no works separating the results of the closeness of missing qualities in data in the midst of this errand, however less than ideal treatment of this kind of learning inside the examination may familiarize slant and may lead with despicable choices being produced using an investigation consider.

This paper demonstrates an audit on the outcome of manhandle the chief apropos philosophies for pre-treatment of missing qualities in the midst of a chose gathering of computations, the common strategy feathery systems for subgroup disclosure. The trial analyze introduced in the midst of this paper exhibit that, among the methods thought, the KNNI pre-taking care of approach for missing qualities gets the least demanding winds up in natural process fleecy systems for subgroup exposure.

Liu, Z.G.; Pan, Q presented a paper [5] Information blend strategy. It is by and large associated inside data classification to help the execution. A soft conviction K-nearest Neighbor (FBK-NN) classifier is expected maintained basic reasoning for directing unverifiable data. For each dissent which is commitment to amass the question, K fundamental conviction assignments (BBA's) are recognized from the partitions among thing and its K-nearest neighbors under thought the neighbors interests. The KBBA's are joined by new strategy and besides the combinations results decide

the class of the question dissent. FBK-NN framework works with is classification and separate one resolute class, Meta classes and discarded/kept up a vital separation from class. Meta-classes are outlined by blend of various specific classifications. The kept up a key separation from class is utilized for anomaly's recognizable proof.

The handiness of the FBK-NN is elucidated by methods for different examinations and their comparative examination with different customary frameworks. In [6], shown clustering part of data, known as ECM (Evidential c-suggests). It is executed with conviction limits. Methodology focuses on the creedal portion strategy, finishing with hard, feathery and ones. Using a FCM like count a perfect target limit is restricted. System similarly recognizes the right number of bundles authenticity record.

In [7] maker challenge the authenticity of Dempster-Shafer Theory. DS oversees gives contrary to yearning come to fruition. Consider exhibits the strategy for affirmation pooling acts against the typical result of the methodology. Still the researcher assembles working in information blend and article knowledge (AI) is as yet arranged to the DS theory. DS control still can't be used or considered for handling the sensible issues. The main role for this is non-appropriateness to confirmation reasoning. In [9] makers show a detail and relative examination of different systems which are: a Singular Value Decomposition (SVD) based procedure (SVDimpute), weighted K-nearest neighbors (KNNimpute), and push typical. These are used to expect missing qualities in quality microarray data. By testing the three methodologies they exhibit that KNN credit is most correct and generous procedure for assessing missing qualities than remaining two strategies outflank the by and large use draw ordinary technique. They report delayed consequences of the comparative examinations and give recommendations and gadgets to correct estimation of missing microarray data under different conditions.

### III. PROBLEM STATEMENT

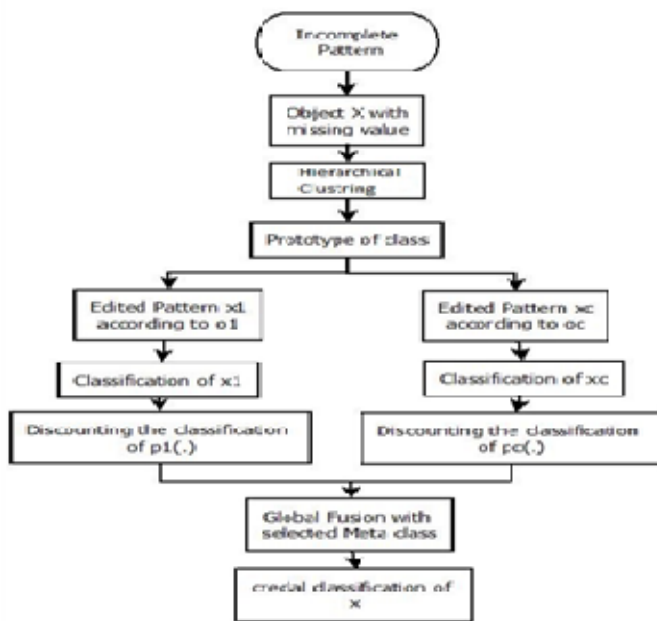
To overcome time, memory and wrong result issues, proposed structure executes evidential reasoning to figure specific class or Meta class for specific case and various leveled clustering to determine the model, which yields capable results similar to time and memory.

### IV. IMPLEMENTATION

#### A. System Architecture

In this structure we are making another strategy to assemble the extraordinary or about hard to sort data with the help of conviction limit Bel (.).In our proposed system we are setting up our structure to tackle missing data from dataset. For this utilization we are using inadequate example dataset as information. For use we can use any standard dataset with missing qualities. Existing system were using mean attribution (MI) approach for registering models in structure. We are using K-Means clustering as starting portion of our use. K-Means clustering gives extra time and memory capable results for our structure than that of mean ascription (MI) framework.

Second some part of our proposed structure is to use dynamic clustering for model calculation. Different various leveled clustering gives more profitable results as diverge from that of K-Means clustering. Henceforth we are focussing on especially dynamic clustering which is used at reason for model creation. After Prototype course of action, we are using the KNN Classifier to describe the examples with the models figured set up of the missing qualities. Since the detachment between the question and the figured model is different we are using the decreasing strategy for the classification. We then wire the classes by using the overall mix control and the as demonstrated by the farthest point regard.



Edge regard gives the amount of the articles that must be fused into the Meta classes. In this manner we augment the exactness by mishitting the question into specific class in case of the vulnerability to describe in one class. We can then apply novel methodology to classifications the challenge into one specific class. In proposed structure we are mostly focussing on time viability in the midst of model improvement.

**B. Algorithms**

**Algorithm 1 Hierarchical Algorithm:**

Input: P objects from dataset
Method:-
1: Amongst the input vector points calculate a distance matrix
2: Every data point must be considered as a cluster.
3: Repeat step 2
4: Combine two nearly similar clusters.
5: Alter distance matrix
6: Go to step 3 until the single cluster remains
7: Stop
Output: Clusters of similar vector.

**Algorithm 2 K means Algorithm:**

Input: N clusters obtained by data set of x objects
Method:-
1: N clusters obtained by data et of x objects.
2: Repeat this 1.
3: Compute distance from centroids to vector.
4: On the basis of mean value of the object in a cluster add every object to the maximum similar cluster.
5: Alter the cluster means.
6: Repeat 3, 4, and 5 until no change.
Output: set of N clusters.

**V. MATHEMATICAL MODEL**

M= (Q, W, P, q0, F) where,

Q is the set of States

W is the set of inputs

P State Transition table q0 is the initial stage

F is the final Stage

1. Q: S1, S2, S3, S4, S5

Where,

S1: Get testing input.

S2: Prototype calculation using hierarchical.

S3: KNN Classification.

S4: Global Fusion using the threshold value and the fusion rule.

S5: Credal classification.

2. W: W1, W2, W3

Where

W1: Incomplete Pattern.

W2: Edited pattern.

W3: Meta Class.

W4: Fusion Data.

3. Q0=S1

4. F: S5

**VI. RESULTS AND DISCUSSION**

*A. Dataset*

Dataset used for proposed system is Breast Cancer and Yeast Data Set that is of Protein Localization Sites. This dataset is assembled from UCI Machine Learning Repository (i.e. <https://archive.ics.uci.edu/ml/datasets/Yeast>). Only 10 to 20 % data or qualities will miss in case of the divided illustrations.

Name	Classes	Attributes	Instances
Cancer	2	9	399
Yeast	3	8	1050

In our use, we use the two veritable enlightening files (malignancy, yeast) open from UCI Machine Learning Repository to test the execution of PCC concerning MI, KNNI, and FCMI. Both EK-NN and ENN are still picked here as standard classifiers. Three classes (CYT, NUC, and ME3) are picked in Yeast educational gathering and two classes (circumspect and hazardous) are picked in Cancer instructive list to our method, since these classes are close and difficult to gathering. The basic information of these instructive files is given in Table.

*B. Result Set*

The outcome set for the paper is for the most part in view of the time and memory examination of the old and the new proposed framework design.

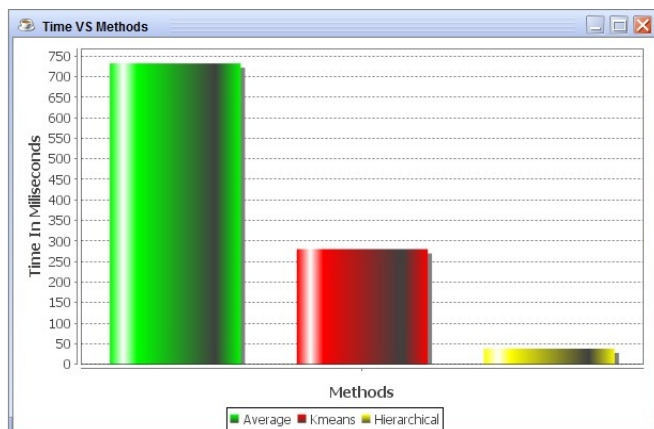


Fig. 2 Time comparison graph

From figure 2 we can see time use of the old structure and proposed system. As ought to be evident that proposed structure puts aside less chance to differentiate and the old or existing system. Proposed structure takes slightest time since it uses different leveled clustering computation for model

figuring and gathering of modified illustrations. Dynamic clustering figuring is more gainful than K-implies computation.

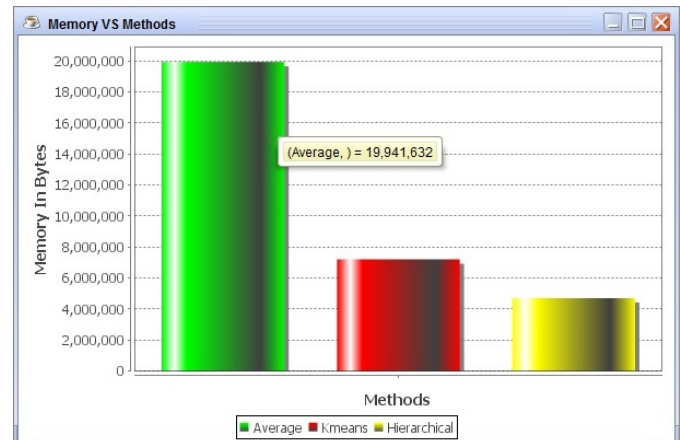


Fig. 3 Memory comparison graph

Figure 3 shows the memory usage by existing framework and proposed framework. As should be obvious that proposed framework devours less memory as contrast and the old or existing framework.

**VI. CONCLUSIONS**

We have proposed a missing example classification for fragmented dissent operation that registers a regard and example by number juggling formula conviction limits. In proposed system evidential intuition describes basic part to miss designs in the dataset. After the discounting system using the conviction work and the edge of the Meta classes the question with deficient example is organized. In case most results square measure tried and true on a classification, the article will be centered on a picked class that is successfully dedicated to the most broadly perceived result. In any case, the high conflict between these results recommends that the classification of the article is somewhat indeterminate or incorrect solely reinforced the far-celebrated far and wide properties data. In such case, the article ends up being horribly hard to classifications really in an exceedingly particular class and it's sensibly conveyed to the benefit meta-class outlined out by the mix of the correct classifications that the article is likely be having a place. By then the conflicting mass of conviction is designated thoroughly to the picked meta-class.

**V. FUTURE SCOPE**

**VI.**

In case the deficient example question is dispersed to a meta-class, it proposes that the correct classifications encased inside the meta-class appear to be ambiguous for this dissent supported the far-celebrated the world over qualities. This structure will be upgraded in taking after ways:

- Client can decide demonstrate a motivating force from manual recognition.
- Diverse clustering estimation can be exchanged for executed different leveled clustering figuring to register the model regard.
- New framework can be used to request last class from meta-classes.

The algorithmic intricacy will be the amount of cycles that are required to organize an inadequate example question suitably to the specific class.

### REFERENCES

- [1] Zhun-Ga Liu, Quan Pan, Grgoire Mercier, and Jean Dezert, "A New Incomplete Pattern Classification Method Based on Evidential Reasoning", Northwestern Polytechnical University, Xian 710072, China, 4, APRIL 2015
- [2] Pedro J. Gracia-Laencina, Jose-Luis Sancho-Gomez, Pattern classification with missing data: a review, Universidad Politecnica de Cartagena, Dpto. Tecnologias de la Information y las Comunicaciones, Plaza del Hospital 1, 30202, Cartagena (Murcia), Spain, 2010.
- [3] SatishGajawada and DurgaToshniwal, "Missing Value Imputation Method Based on Clustering and Nearest Neighbours", The Department of Electronics and Computer Engineering, Indian Institute of Technology Roorkee, Roorkee, India, 2012.
- [4] Cristobal J. Carmona, Julian Luengo, "An analysis on the use of pre-processing methods in evolutionary fuzzy systems for subgroup discovery", Department of Computer Science, University of Jaen, Campus lasLagunillas, 23071 Jaen, Spain, 2012.
- [5] K.Pelckmans, J.D.Brabanter, J. A. K. Suykens, and B.D.Moor, "Handling missing values in support vector machine classifiers, Neural Netw., vol. 18, nos. 5-6, pp. 684-692, 2005.
- [6] P. Chan and O. J. Dunn, "The treatment of missing values in discriminant analysis," J. Amer. Statist. Assoc., vol. 6, no. 338, pp. 473-477, 1972.
- [7] F. Smarandache and J. Dezert, "Information fusion based on new proportional conflict redistribution rules," in Proc. Fusion Int. Conf. Inform. Fusion, Philadelphia, PA, USA, Jul. 2005.
- [8] J. L. Schafer, Analysis of Incomplete Multivariate Data. London, U.K.: Chapman Hall, 1997.
- [9] O. Troyanskaya et al., "Missing value estimation method for DNA microarrays," Bioinformatics, vol. 17, no. 6, pp. 520-525, 2001.
- [10] G. Batista and M. C. Monard, "A study of K-nearest neighbour as an imputation method," in Proc. 2nd Int. Conf. Hybrid Intell. Syst., 2002, pp. 251-260.
- [11] Farhangfar, Alireza, Lukasz Kurgan, "Impact of imputation of missing values on classification error for discrete data", Pattern Recognition, pp. 3692-3705, 2008.
- [12] F. Smarandache and J. Dezert, "On the consistency of PCR6 with the averaging rule and its application to probability estimation", Proceedings of the International Conference on Information Fusion, pp.323-330, July 2013.
- [13] Z.-G. Liu, J. Dezert, G. Mercier, and Q. Pan, "Belief C-means: An extension of fuzzy C-means algorithm in belief functions framework," Pattern Recognition, vol. 33, no. 3, pp. 291-300, 2012.
- [14] P. Garcia-Laencina, J. Sancho-Gomez, A. Figueiras-Vidal, "Pattern classification with missing data: A review", Neural Networks, vol. 19, no. 2, pp. 263-282, 2010.
- [15] A. Tchamova, J. Dezert, "On the Behavior of Dempster's rule of combination and the foundations of Dempster-Shafer theory", In proceedings of Sixth IEEE International Conference on Intelligent Systems, pp. 108-113, 2012.
- [16] Z.-G. Liu, J. Dezert, G. Mercier, and Q. Pan, "Dynamic evidential reasoning for change detection in remote sensing images," IEEE Geosci. Remote Sens., vol. 50, no. 5, pp. 1955-1967, May 2012.
- [17] M.-H. Masson and T. Denoeux, "ECM: An evidential version of the fuzzy C-means algorithm," Pattern Recognit., vol. 41, no. 4, pp. 1384-1397, 2008.
- [18] T. Denoeux and M.-H. Masson, "EVCLUS: Evidential Clustering of proximity data," IEEE Trans. Syst., Man, Cybern. B, Cybern., vol. 34, no. 1, pp. 95-109, Feb. 2004.
- [19] Z.-G. Liu, J. Dezert, G. Mercier, and Q. Pan, "Belief C-means: An extension of fuzzy C-means algorithm in belief functions framework," Pattern Recognit. Lett., vol. 33, no. 3, pp. 291-300, 2012.
- [20] T. Denoeux, "Maximum likelihood estimation from uncertain data in the belief function framework," IEEE Trans. Knowl. Data Eng., vol. 25, no. 1, pp. 119-130, Jan. 2013.
- [21] Z.-G. Liu, J. Dezert, Q. Pan, and G. Mercier, "Combination of sources of evidence with different discounting factors based on a new dissimilarity measure," Decision Support Syst., vol. 52, no. 1, pp. 133-141, Dec. 2011.
- [22] A. Tchamova and J. Dezert, "On the Behavior of Dempster's rule of combination and the foundations of Dempster-Shafer theory," in Proc. 6th IEEE Int. Conf.

- Intell. Syst. (IS'12), Sofia, Bulgaria, Sep. 2012, pp. 108–113.
- [23] D. Dubois and H. Prade, “Representation and combination of uncertainty with belief functions and possibility measures,” *Comput. Intell.*, vol. 4, no. 4, pp. 244–264, 1988.
- [24] F. Smarandache and J. Dezert, “Information fusion based on new proportional conflict redistribution rules,” in *Proc. Fusion Int. Conf. Inform. Fusion*, Philadelphia, PA, USA, Jul. 2005.
- [25] F. Smarandache and J. Dezert, “On the consistency of PCR6 with the averaging rule and its application to probability estimation,” in *Proc. Fusion Int. Conf. Inform. Fusion*, Istanbul, Turkey, Jul. 2013.
- [26] C. M. Bishop, *Neural Networks for Pattern Recognition*. London, U.K.: Oxford Univ. Press, 1995.
- [27] B. W. Silvean, M. C. Jones, E. Fix, and J. L. Hodges, “An important contribution to nonparametric discriminant analysis and density estimation—Commentary on Fix and Hodges (1951),” *Int. Statist. Rev.*, vol. 57, no. 3, pp. 233–227, Dec. 1989.
- [28] L. A. Zadeh, *On the Validity of Dempster’s Rule of Combination of Evidence*, Memo M79/24, Berkeley, CA, USA: Univ. California, 1979.
- [29] J. Lemmer, “Confidence factors, empiricism and the Dempster–Shafer theory of evidence,” in *Proc. 1st Conf. UAI*, 1985, pp. 160–176.
- [30] G. M. Provan, “The validity of Dempster–Shafer belief functions,” *IJAR*, vol. 6, no. 3, pp. 389–399, May 1992.
- [31] P. Wang, “A Defect in Dempster–Shafer Theory,” in *Proc. 10th Conf. Uncertainty AI*, 1994, pp. 560–566.
- [32] 646 *IEEE TRANSACTIONS ON CYBERNETICS*, VOL. 45, NO. 4, APRIL 2015 J. Dezert and A. Tchamova, “On the validity of Dempster’s fusion rule
- [33] and its interpretation as a generalization of Bayesian fusion rule,” *Int. J. Intell. Syst.*, vol. 29, no. 3, pp. 223–252, Mar. 2014.
- [34] L. M. Zouhal and T. Denoeux, “An evidence-theoretic k-NN rule with parameter optimization,” *IEEE Trans. Syst., Man, Cybern., C, Appl. Rev.*, vol. 28, no. 2, pp. 263–271, May 1998.
- [35] A. Frank and A. Asuncion. (2010). *UCI Machine Learning Repository*, University of California, School of Information and Computer Science, Irvine, CA, USA [Online]. Available: <http://archive.ics.uci.edu/ml> S. Geisser, *Predictive Inference: An Introduction*. Boston, MA, USA: Chapman & Hall, 1993.