# M-Denclue for Effective Data Clustering in High Dimensional Non-Linear Data

R.Nandhakumar, Antony Selvadoss Thanamani

Abstract— *Clustering is a data mining task devoted to the automatic grouping of data based on mutual similarity. Clustering in high-dimensional spaces is a recurrent problem in many domains. It affects time complexity, space complexity, scalability and accuracy of clustering methods. High-dimensional non-linear datausually live in different low dimensional subspaces hidden in the original space. As high-dimensional objects appear almost alike, new approaches for clustering are required. This research has focused on developing Mathematical models, techniques and clustering algorithms specifically for high-dimensional data. The innocent growth in the fields of communication and technology, there is tremendous growth in high dimensional data spaces. As the variant of dimensions on high dimensional non-linear data increases, many clustering techniques begin to suffer from the curse of dimensionality, de-grading the quality of the results. In high dimensional non-linear data, the data becomes very sparse and distance measures become increasingly meaningless. The principal challenge for clustering high dimensional data is to overcome the "curse of dimensionality". This research work concentrates on devising an enhanced algorithm for clustering high dimensional non-linear data.*

*Keywords: Clustering, High Dimensional Non Linear Data, curse of dimensionality, Mathematical models.*

## I. INTRODUCTION

Clustering is among the main data exploration jobs and is aimed at group the data items into significant classes (clusters) in a way that your similarity of products inside clusters is obviously maximized and the similarity of products from several clusters is generally minimized. Ton evaluation is generally among the main equipment made for exploring the essential framework of the info collection. Clustering agrees with important applications in a broad collection of professions incorporating handy remote control realizing, routine acceptance, photo application and computer program eyesight. The best objective of any clustering strategy is generally to rupture verified info place comprising N-dimensional elements or perhaps vectors in to a mounted level of Addition clusters. Typically, clustering is going to be thought to be a method that partitioning the information strategies into mutually exclusive complexes or types in a manner that details factors in the same ton are a lot more similar one to the other than to data elements in extra clusters. The dissimilarity between a few facts points is normally measured with a variety metric defined in the differences between your values with the features (dimensions).

**R.Nandhakumar,** Assistant Professor, Department of Computer Science, Nallamuthu Gounder Mahalingam College, Pollachi-642001, Tamilnadu, India (E-mail: nkumarram@gmail.com)

**Dr.Antony Selvadoss Thanamani,** Associate Professor & Head, Department of Computer Science, Nallamuthu Gounder Mahalingam College, Pollachi-642001, Tamilnadu, India

Traditional clustering algorithms utilize all the parts in the details to compute the traces. The bane of dimensionality for neo linear info makes the clustering work very difficult if the information space consists of numerous features. The many attributes helps that end up being computationally infeasible to utilize each one of the attributes to have the clusters. Besides, not all the features happen to be of heap for the clustering work. The fewer relevant features trigger the normal denseness to á ton in practically any society of the details space to be low that means it is difficult to find any essential clusters making use of the initial clustering algorithms in full-dimensional space. This research function focuses primarily on devising an excellent improved duodecimal program with clustering increased dimensional non-linear info.

## II. LITERATURE SURVEY

Clustering high dimensional data is definitely a problem for clustering techniques. This issue offers been studied thoroughly and there are numerous solutions, every befitting various kinds of high dimensional data and data exploration methods. There are numerous potential applications like bioinformatics, text message gold mining with large dimensional info where subspace clustering, forecasted clustering methods may help to discover patterns skipped by current clustering strategy. To be able to gain conceptual clearness of the domain name under research various content articles, books, websites plus some of other personal reviews had been examined. The review provides been carried out by directing on the main element subject of clustering substantial dimensional info.

Maithri. C presented a short a comparison of the prevailing methods that were primarily concentrating in clustering on high dimensional data. The primary objective of the study newspaper is to show the potency of excessive dimensional info evaluation and various algorithm inside the prediction procedure for Data exploration. The overall performance issues of the info clustering in great dimensional data, additionally it is essential to study problems like dimensionality decrease, redundancy elimination, subspace clustering, co-clustering and info Labeling intended for clusters are to analyzed and increased.

SunitaJahirabadkar, presented an assessment of various denseness centered subspace clustering codes as well as a comparative graph concentrating on their particular distinguishing features such as for example overlapping /

non overlapping, axis similar / randomly oriented and so forth. Charles Bouveyron, presents a clustering strategy which estimates the precise subspace and the inbuilt dimension of every class. Their particular strategy gets used to the Gaussian combination unit framework to high-dimensional data and estimations the guidelines which greatest fit the info. We get yourself a robust clustering technique known as Large Dimensional Data Clustering. They used Great Dimensional Data Clustering to find items in organic pictures within a probabilistic platform. Experiments on a lately proposed data source demonstrate the potency of our clustering way for category localization.

E. Kailing, launched a SUBCLU (density- linked Subspace Clustering), a competent method of the subspace clustering issue. Applying the idea of thickness connection fundamental the formula DBSCAN, SUBCLU is founded on an official clustering idea. As opposed to existing grid-based techniques, SUBCLU will be able to detect randomly formed and positioned groupings in subspaces.

## III. EXECUTION PHASES

Stage 1: Applying Denclue, Optical technologies and Fanfare Algorithm upon High-dimensional Non- Linear Dataset.

Stage 2: Predicated on stage 1 effect, two methods (DENCLUE and OPTICS) had been chosen like a greatest formula. To resolve some disadvantages of the two algorithms several mathematical strategies such as for example traguardo heuristics, curse of dimensionality, data redirecting, correlation, regular distribution and Darboux variate had been added with these algorithms. Finally, these modified algorithms had been applied about High dimensional nonlinear info set and result.

## IV. PROPOSED FRAMEWORK

Clustering high dimensional data is definitely challenging because of its dimensionality issue and this affects period complexity, space complexity, scalability and precision of clustering methods. Numerous clustering strategies can be found such as for example hierarchical founded technique, canton based technique, density depending technique, main grid based technique and unit based technique. Among these types of both Denclue algorithm and Optics routine are comes beneath the occurrence centered clustering technique, where as Fanfare algorithm comes beneath the Main grid structured clustering technique. In density centered clustering, the items will be classified predicated on their parts of density. These types of algorithms be capable of discover classes of human judgments designs and omit boisterous items. Upon other submit grid based upon clustering, the info are split into grid of objects. This kind used the algorithm ón the main grid, instead of used it on the data source.

DENCLUE (DENsity-based CLUstEring) is usually a clustering approach predicated on a couple of density passing them out functions. The technique is created on the next ideas: (1) the impact of every info point could be formally patterned utilizing a statistical function, named an impact function, which is the influence of an info stage within just its area; (2) the entire density of the info space

could be made analytically as the amount of money of the impact function put on pretty much all data factors; and (3) clusters may then be identified mathematically by simply determining occurrence attractors, where density attractors are regional maxima of the entire solidity function. However the drawbacks of the algorithm happen to be be; it really is fewer delicate to outliers. It generally does not work very well just for high dimensional data, due to the bane óf dimensionality phenomenon. The density variable and the sound limit have to be picked out carefully since it significantly influences the standard of benefits.
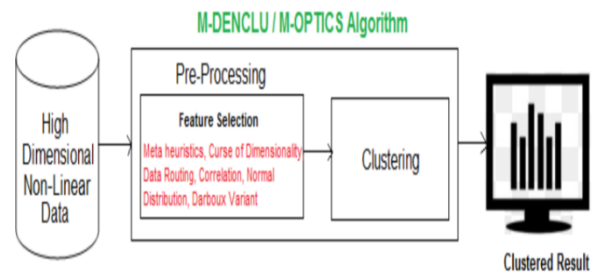


**Figure 1: M-DENCLUE/M-OPTICS Algorithm**

OPTICS technologies algorithm functions in guideline like this prolonged DBSCAN algorithm just for thousands of length details $\epsilon i$ that happen to be smaller sized when compared to a "generating distance" $\epsilon$ ( i actually. vitamin e. $0 \le \epsilon i \le \epsilon$). The sole difference is usually that people you don't need to assign group memberships. Rather, we retail store the purchase in which the items are refined and the info which can be utilized by a protracted DBSCAN routine to designate cluster memberships. However the downsides of the duodecimal system are that expects some type of density refuse to discover cluster limits, in fact it is fewer delicate to erroneous info. Physique one particular displays the M-DENCLUE/M-0PTICS Routine on Great Dimensional Info Set

FANFARE is the initial subspace clustering algorithm. The CLIQUE algorithm discovers the crowed area from the multidimensional data source and discovers the patterns. Any time the machine is going to be dense after that it remains to create a good cluster. The outlier recognition of groupings and be done ? complete about the noisy info also a significant part of superior dimensional info pieces. To take action, clusters happen to be analyzed regarding positive and negative items in CLlQUE by intra-cluster similarity of clusters with regards to the occurrence of negative and positive items through RandIndex. The obsolete items are actually eliminated coming from the spot by simply matrix factorization and canton technique. Some drawbacks of this mechanism are as; Need to tune grid size and density threshold. May fail if clusters are of widely differing densities, since the threshold is fixed. Can still have high mining cost. Same density threshold for low and high dimensionality. DENCLU and OPTICS are density based clustering technique, where as CLIQUE comes under the grid-based clustering technique. Compare to DENCLU

and OPTICS algorithms CLIQUE algorithm provides less performance on result.

So in this research work DENCLU and OPTICS algorithms were selected and modified by adding the mathematics methods such as meta-heuristics, curse of dimensionality in affine sub spaces, data routing, correlation, normal distribution and darboux variate. It performs pre-processing process on data set before implement with our modified algorithm. Multiple sizes will be hard to believe in, difficult to visualize, and, because of the rapid growth of the amount of possible ideals with every dimension, total enumeration of most subspaces turns into intractable with raising dimensionality.

## V. EXPERIMENTAL RESULTS AND ANALYSIS

The Clustering Algorithms DENCLUE, OPTICS and CLIQUE were experimented with the Bio informatics - DNA microarray Datasetwith the implementation of MATLAB R2018b (Version 9.5) - Sep 2018, and the findings yielded that the CLIQUE algorithm did not perform well for Clustering of High Dimensional non-linear data. Then the DENCLUE and OPTICS algorithms were experimented and analysed and the limitations of these algorithms were recorded and swamped with the mathematical models - meta heuristics, curse of dimensionality, data routing, correlation, normal distribution and Darbouxvariate. This process has enhanced the DENCLUE and OPTICS algorithms into M-DENCLUE and M-OPTICS algorithms. These enhanced algorithms were tested for erroneous data, detection of outlier,noisy data, mining performance, computational complexity, execution speed, data quality, scalability and accuracy.

The Research Experiments implemented in MATLAB revealed that M-DENCLUE algorithm suits best for clustering High Dimensional Non-Linear Data experimented with Bio informatics - DNA microarray Dataset.

## VI. CONCLUSION

Clustering High dimensional Non-Linear data sets is a challenging laborious task. The principal challenge for clustering high dimensional data is to overcome the "curse of dimensionality". This research work studied and analyzed various clustering techniques for high dimensional non-linear data clustering, and analyze the limitations of Clustering in High Dimensional Non-Linear data, an effective solution was provided to enhance the performance of clustering on high dimensional non-linear data clustering by overcoming the 'Curse of Dimensionality', by analyzing DENCLUE, OPTICS and CLIQUE algorithms for clustering high dimensional data. The limitations of these algorithms were overcome with incorporating the mathematical concepts of meta-heuristics, curse of dimensionality, sub spaces, data routing, correlation, normal distribution and darboux variate, which has proposed new enhanced algorithms M-DENCLUE and M-OPTICS. The Bio informatics - DNA microarray Dataset which is High Dimensional Non- Linear in nature was used for experimenting the enhanced algorithms and a best fit for clustering High Dimensional Non-Linear data is obtained.

## REFERENCES

1. A. Hinneburg and D. A. Keim, "*Optimal grid clustering: Towards breaking the curse of dimensionality in high dimensional clustering,*" In Proceedings of 25th International Conference on Very Large Data Bases (VLDB-1999), pp. 506-517.
2. Charles Bouveyron, Stephane Girard, et al., "*High Dimensional Data Clustering*", Computational Statistics and Data Analysis 52, 1 (2007) 502-519.
3. E. Müller, S. Günnemann, I. Assent and T. Seidl (2009), "*Evaluating clustering in subspace projections of high dimensional data*", In Proc. of the Very Large Data Bases Endowment, Volume 2 issue 1, pp. 1270-1281.
4. P. Lance, E. Haque, and H. Liu (2004), "*Subspace clustering for high dimensional data: A review*", ACM SIGKDD Explorations Newsletter, Vol. 6 Issue 1, pp 90–105.

## AUTHORS PROFILE

Name: R.NANDHAKUMAR
Assistant Professor ,
Department of Computer Science ,
Nallamuthu Gounder Mahalingam College of Arts and Science,
Pollachi-642001,India.

Name: Dr. ANTONY SELVADOSS THANAMANI
Associate Professor & Head
Department of Computer Science,
Nallamuthu Gounder Mahalingam College of Arts and Science,
Pollachi-642001,India.