

## A Clustering Technique for Reducing Noise in High Dimensional Non-Linear Data Using M-DENCLUE Algorithm

R. Nandhakumar<sup>1</sup>, Dr. Antony Selvadoss Thanamani<sup>2</sup>

<sup>1</sup> Assistant Professor, Department of Computer Science, NGM College, Pollachi-642001, India.  
nkumarram@gmail.com.

<sup>2</sup> Associate Professor & Head, Department of Computer Science, NGM College, Pollachi-642001, India.

### ABSTRACT

Clustering is a method in data mining which deals with huge amount of data. Clustering is intended to assist a consumer in discovering and know-how the herbal structure in a statistics set and abstract the which means of massive dataset. It is the undertaking of partitioning objects of a statistics set into awesome businesses such that two gadgets from one cluster are similar to every other, while objects from wonderful clusters are assorted. Clustering is unsupervised getting to know in which we are not provided with instructions, in which we will area the records items.

With the arrival growth of high dimensional statistics including microarray gene expression facts, and grouping excessive dimensional statistics into clusters will come across the similarity among the items in the full dimensional area is frequently invalid as it consists of exclusive styles of information. The technique of grouping into high dimensional information into clusters is not accurate and possibly not as much as the extent of expectation when the dimension of the dataset is high.

**Key words :** Clustering, Micro array, Noise, Density based, Grid based.

### 1. INTRODUCTION

Clustering is unsupervised learning wherein we aren't supplied with lessons, in which we will location the records items. Clustering is beneficial over category due to the fact fee for labeling is reduced. Clustering has programs in molecular biology, astronomy, geography, client relation control, textual content mining, net mining, etc. Clustering can be used to expect client shopping for patterns based on their profiles to which cluster they belong.

A cluster described as a dense factor, wherein it can grow in any route that density leads [2]. There are processes. The first approach is a density to a schooling facts point like DBSCAN and OPTICS. The 2d method is a density to a data point within the attribute space makes use of a density feature like DENCLUE.

**Curse of Dimensionality** - Dimensionality curse is one of the main issues confronted by excessive dimensional data [4]. In

high dimensional space the points are extra scattered or sparse and all factors are nearly equidistant from every other. Clustering tactics become useless to examine the facts because of this.

**Noise-** The noise present in real packages frequently hides the clusters to be decided on from clustering algorithm and the hassle is worsened in high dimensional data, where the variety of mistakes increases linearly with dimensionality [4].

DENCLUE set of rules and OPTICS algorithm comes underneath the density based clustering approach, in which as CLIQUE algorithm comes beneath the Grid primarily based clustering approach. Clustering in High Dimensional Non-Linear data spaces is a recurrent trouble in many domain names. It influences time complexity, area complexity, Data Size Adaptability and Precision Value of clustering methods.

### 2. RELATED WORK

Clustering is the grouping collectively of comparable records gadgets into clusters. Clustering analysis is one of the fundamental analytical techniques in facts mining; the method of clustering algorithm will have an effect on the clustering effects at once.

Mythili .S and Madhiya mentioned the numerous kinds of algorithms like okay-way clustering algorithms, and so forth. And examine the blessings and shortcomings of the various algorithms. In each type we are able to calculate the distance between each statistics item and all cluster centers in each generation, which makes the Competence Rate of clustering isn't always excessive.

Clustering is the unsupervised class of patterns into companies. The clustering hassle has been addressed in many contexts and by using researchers in many disciplines; this displays its wide attraction and usability as one of the steps in exploratory statistics analysis. M.N. Murty et al., offers a top level view of sample clustering techniques from a statistical pattern recognition attitude, with a aim of imparting useful advice and references to fundamental standards on hand to the broad network of clustering practitioners. We gift taxonomy of clustering techniques, and discover pass-cutting subject matters and current advances.

Sulbha Patil affords reducing as an technique. Data anonymization is a warm studies topic. Slicing can manage high dimensional records and it preserves higher facts utility.

Slicing firsts partitions attributes into columns. Each column incorporates a subset of attributes. Number attributes is identical to wide variety of columns.

XZhang et. Al.. Experimented with actual world large statistics set in cloud from the angle of shielding privacy breaches and to gain excessive degree of Data Size Adaptability and Competence Rate. They proposed a proximity privacy version and a scalable two segment clustering approach based on MapReduce acting data parallel computation in cloud to cope with the difficulty of privateness.

M. Suriyapriya and A. Joicy endorse a scheme that is resilient to replay assaults. In this scheme the usage of Secure Hash set of rules for authentication reason, SHA is one of the several cryptographic hash functions, most usually used to verify that a report has been unaltered. The Paillier crypto machine is a probabilistic uneven set of rules for public key cryptography. Pailier set of rules use for Creation of get admission to policy, report gaining access to and file restoring system. Clouds have large storage space for storing big amount of information. Data owner outsource their information contents on cloud server. Cloud server may have large storage space.

### 3. PROPOSED MODEL

Clustering or facts grouping is the important thing technique of the records mining. It is an unmonitored gaining knowledge of assignment where one seeks to identify a finite set of classes termed clusters to describe the facts. The grouping of records into clusters is based totally on the principle of maximizing the intra class similarity and minimizing the inter magnificence similarity.

#### 3.1. Clustering Technique

The grouping of records into clusters is based totally on the precept of maximizing the intra class similarity and minimizing the inter class similarity. A true clustering technique will produce excessive nice clusters with excessive intra-class similarity - Similar to each other within the identical cluster low inter-elegance similarity [9]. The great of a clustering approach is likewise measured by its ability to find out some or all the hidden styles.

The objective of the clustering technique is to decide the intrinsic grouping in a set of unlabeled facts. The similarity among facts items can be measured with the imposed distance values. Specifying the gap measures for the high dimensional facts is turning into very trivial as it holds distinctive statistics values of their corresponding attributes.

Data mining allows extracting facts from the huge data and changing that statistics into a reasonable and vital shape for additionally utilize [10]. Data mining is a fundamental mission at some stage in the time spent getting to know revelation from big information. Data mining is an enhance mechanism this is very beneficial to mine the understandable expertise, formerly unknown, facts from big amount of information saved in various formats, with the objectives of improving the selection of corporations, agencies where the facts might be collected.

### 3.2. High-Dimensional Data in Knowledge Discovery Database

Clustering excessive dimensional data in a gene expression microarray information set, there can be tens or masses of dimensions, every of which corresponds to an experimental situation. Curse Dimensionality is a loose manner of speaking approximately information separation in high dimensional space. The complexity of many existing facts mining algorithms is exponential with recognize to the variety of dimensions. Each group is a dataset such that the similarity many of the statistics in the institution is maximized and the similarity in outside institution is minimized.

#### 3.3 Density based clustering algorithm

Density primarily based clustering algorithms are very popular inside the applications of facts mining. These procedures use a local cluster criterion and outline clusters because the regions in the data area of better density in comparison to the areas of noise points or border points. Density based totally clustering algorithms using the notion of DBSCAN, can locate clusters of arbitrary length and shape [2]. Density-primarily based clustering may be seen as a non-parametric technique, in which clusters are modeled as regions of high density. CLIQUE is the primary grid based totally subspace clustering approach designed for excessive dimensional statistics. It detects subspaces of the very best dimensionalities.

### 4. RESULT AND DISCUSSION

M-DENCLUE works on two ranges as pre-processing degree and clustering stage. In pre-processing step, it creates a grid for the facts by means of dividing the minimal bounding hyper-rectangle into d-dimensional hyper-rectangles with edge length  $2\sigma$ . In the clustering stage, M-DENCLUE associates an "affect feature" with every facts factor and the general density of the dataset is modeled as the sum of impact functions related to every point [5]. The resulting preferred density feature may have neighborhood peaks, i.E., nearby density maxima, and these nearby peaks can be used to outline clusters.

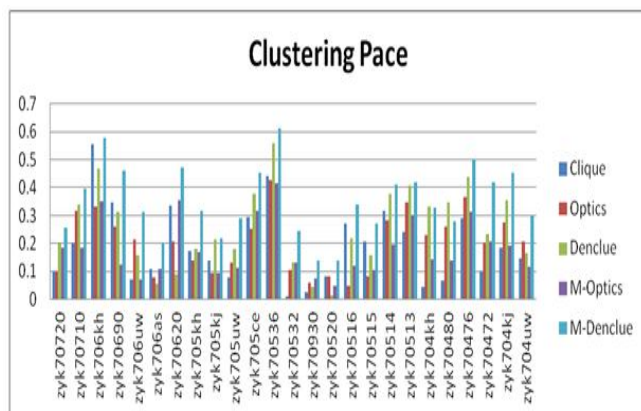
M-DENCLUE makes use of influence features. Influence of every records point may be modeled as mathematical function. The resulting function is referred to as Influence Function. Influence feature illustrates the effect of facts factor inside its neighbourhood.

The M-OPTICS set of rules creates an ordering of the objects in a database, M-OPTICS additionally storing the middle-distance and a appropriate reachability distance for each item. An set of rules turned into proposed to extract clusters based totally on the ordering records produced by M-OPTICS and as soon as the order and the reachability distances are computed, we can extract the clusters for any clustering distance.

The work is illustrated via graphs with the help of - DNA microarray Data.

**Table 1:** Clustering Pace on DNA Microarray data set Clique, Optics, DENCLUE, M-Optics and M-DENCLUE Algorithm

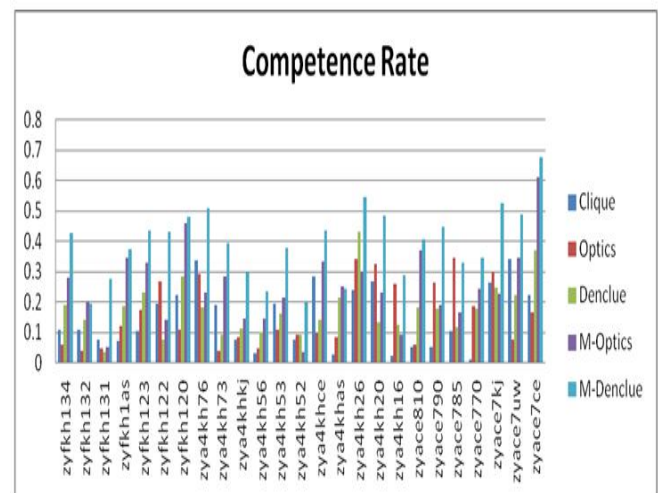
	Clique	Optics	Denclue	M-Optics	M-Denclue
zyk70720	0.101391	0.10273	0.199473	0.183355	0.25517338
zyk70710	0.200161	0.317208	0.337633	0.182342	0.39739786
zyk706kh	0.555159	0.333244	0.469775	0.351924	0.57830717
zyk70690	0.347334	0.260886	0.31413	0.124695	0.4606677
zyk706uw	0.071067	0.214153	0.159508	0.072177	0.31296296
zyk706as	0.108637	0.077021	0.05461	0.109321	0.20170116
zyk70620	0.335354	0.206745	0.090335	0.352806	0.47094082
zyk705kh	0.174004	0.139521	0.181694	0.170664	0.31579885
zyk705kj	0.137659	0.091897	0.213333	0.09523	0.22
zyk705uw	0.079098	0.131335	0.179458	0.112586	0.29094511
zyk705ce	0.29313	0.252728	0.378045	0.317399	0.45156635
zyk70536	0.441155	0.425971	0.558058	0.414736	0.6110013
zyk70532	0.011645	0.106578	0.132544	0.129983	0.24398184
zyk70930	0.025174	0.060603	0.044888	0.075757	0.13996714
zyk70520	0.081876	0.080965	0.015625	0.046639	0.140625
zyk70516	0.269663	0.047389	0.219693	0.121003	0.33900181
zyk70515	0.205245	0.0817	0.159397	0.105742	0.26987437
zyk70514	0.315193	0.281765	0.378375	0.19719	0.410431
zyk70513	0.242647	0.345973	0.40807	0.300967	0.41970008
zyk704kh	0.042641	0.230318	0.330698	0.144264	0.32906407
zyk70480	0.068823	0.258209	0.346408	0.140234	0.27737333
zyk70476	0.291449	0.367389	0.43933	0.311252	0.50195872
zyk70472	0.099653	0.202734	0.234231	0.207511	0.42033378
zyk704kj	0.185863	0.275163	0.355582	0.191243	0.45221191
zyk704uw	0.146397	0.206896	0.166375	0.117798	0.30009238



**Figure 1:** Experimental Clustering Pace

**Table 2:** Experimental Inference-Competence Rate

	Clique	Optics	Denclue	M-Optics	M-Denclue
zyfkh134	0.10726	0.059277	0.190335	0.281155	0.42707339
zyfkh132	0.108917	0.041342	0.140705	0.204966	0.19507669
zyfkh131	0.077543	0.047218	0.036129	0.052473	0.27699428
zyfkh1as	0.070764	0.123056	0.185223	0.346021	0.37236181
zyfkh123	0.104919	0.175736	0.229547	0.331158	0.43615085
zyfkh122	0.192935	0.266635	0.076078	0.143506	0.43048422
zyfkh120	0.222095	0.108282	0.285689	0.461416	0.48158565
zya4kh76	0.337714	0.292634	0.180675	0.230791	0.50999927
zya4kh73	0.191967	0.0418	0.091544	0.284445	0.39554934
zya4khkj	0.078303	0.08539	0.114213	0.146326	0.29810327
zya4kh56	0.030084	0.047161	0.100025	0.14649	0.23644031
zya4kh53	0.193828	0.108925	0.164135	0.217014	0.37879008
zya4kh52	0.077756	0.093267	0.091911	0.034201	0.20067229
zya4khce	0.286268	0.09681	0.140163	0.332161	0.43385884
zya4khas	0.028523	0.083555	0.216592	0.252168	0.24257063
zya4kh26	0.241483	0.342935	0.43169	0.302444	0.54731229
zya4kh20	0.269746	0.325518	0.132449	0.230701	0.48624276
zya4kh16	0.022534	0.261976	0.127218	0.092992	0.28767229
zyace810	0.052984	0.058634	0.182239	0.369891	0.40614054
zyace790	0.053503	0.265784	0.179295	0.189185	0.44763072
zyace785	0.10673	0.344255	0.11689	0.165832	0.32922957
zyace770	0.013127	0.186011	0.180138	0.244415	0.34674772
zyace7kj	0.265722	0.295721	0.246047	0.227704	0.52521232
zyace7uw	0.343682	0.07658	0.224918	0.343935	0.48799682
zyace7ce	0.223578	0.164429	0.370004	0.608936	0.67430078



**Figure 2:** Experimental Competence Rate

**5. CONCLUSION AND FUTURE WORK**

DENCLU and OPTICS are density primarily based clustering technique, wherein as CLIQUE comes beneath the grid-based totally clustering approach. Comparing all those ultimately finish that DENCLUE is the nice one. Since everyday existence adjustments with digital global, to group specific information. DENCLUE helps in reducing noise.

From experimental consequences it has been determined that large and dense records wishes better computational strength. In future the troubles encountered inside the current strategies can be conquering with the aid of growing a hybrid primarily based density algorithm.

## REFERENCES

- [1]. Dr. Anjali B. Raut, "A Hybrid Framework using Fuzzy if-then rules for DBSCAN Algorithm", *Advances in Wireless and Mobile Communications*, ISSN 0973-6972 Volume 10, Number 5 (2017), pp. 933-942.
- [2]. Feng Cao, Weining Qian, "Density-Based Clustering over an Evolving Data Stream with Noise", Department of Computer Science and Engineering, Fudan University.
- [3]. Gaff, B. M., Sussman, H. E., & Geetter, J., "Privacy and big data", *Computer*, 47(6), 7–9. doi:10.1109/mc.2014.161, 2014. <https://doi.org/10.1109/MC.2014.161>
- [4]. Gosain Anjana & Chugh Nikita, "Privacy Preservation in Big Data", *International Journal of Computer Application*, Vol. 100 No.17 August 2014. <https://doi.org/10.5120/17619-8322>
- [5]. Hajar Rehioui, Abdellah Idrissi, "DENCLUE-IM: A New Approach for Big Data Clustering", *The 7th International Conference on Ambient Systems, Networks and Technologies (ANT 2016)*, Science Direct, *Procedia Computer Science* 83 (2016) 560 – 567. <https://doi.org/10.1016/j.procs.2016.04.265>
- [6]. Harsh Shah, Karan Napanda, "Density Based Clustering Algorithms", *International Journal of Computer Sciences and Engineering*, Volume-3, Issue-11, E-ISSN: 2347-2693.
- [7]. Harsh Shah, Karan Napanda, "Density Based Clustering Algorithms", *International Journal of Computer Sciences and Engineering*, Review Paper, Volume-3, Issue-11, E-ISSN: 2347-2693, 2015.
- [8]. Hadi Saboohi, "On Density-Based Data Streams Clustering Algorithms: A Survey", *Journal of Computer Science and Technology* 29(1): 116{141 Jan. 2014. <https://doi.org/10.1007/s11390-014-1416-y>
- [9]. Singh Vijendra, "Efficient Clustering for High Dimensional Data: Subspace Based Clustering and Density Based Clustering", *Information Technology Journal*, Volume 10 (6): 1092-1105, 2011. <https://doi.org/10.3923/itj.2011.1092.1105>
- [10]. Aina Musdholifah And Siti Zaiton Mohd Hashim, "Cluster Analysis On High-Dimensional Data: A Comparison Of Density-Based Clustering Algorithms", *Australian Journal Of Basic And Applied Sciences*, 7(2): 380-389, 2013.