

# Clustering High Dimensional Non-Linear Data with Denclue, Optics and Clique Algorithms

R. Nandhakumar, Antony Selvadoss Thanamani



**Abstract---** Clustering is a technique in data mining which deals with huge amount of data. Clustering is intended to help a user in discovering and understanding the natural structure in a data set and abstract the meaning of large dataset. It is the task of partitioning objects of a data set into distinct groups such that two objects from one cluster are similar to each other, whereas two objects from distinct clusters are dissimilar. Clustering is unsupervised learning in which we are not provided with classes, where we can place the data objects.

With the advent growth of high dimensional data such as microarray gene expression data, and grouping high dimensional data into clusters will encounter the similarity between the objects in the full dimensional space is often invalid because it contains different types of data. The process of grouping into high dimensional data into clusters is not accurate and perhaps not up to the level of expectation when the dimension of the dataset is high.

**Keywords---** Clustering, DNA Micro array, Noise, Density based, Grid based.

## I. INTRODUCTION

Bunching is unaided learning in which we are not furnished with classes, where we can put the information objects. Bunching is advantageous over characterization since expense for naming is decreased. Bunching has applications in atomic science, stargazing, geology, client connection the board, content mining, web mining, and so on. Grouping can be utilized to anticipate client purchasing behaviors dependent on their profiles to which bunch they have a place.

A bunch characterized as a thick segment, where it can develop toward any path that thickness leads. There are two approaches. The main methodology is a thickness to a preparation information point like DBSCAN and OPTICS. The subsequent methodology is a thickness to an information point in the property space utilizes a thickness capacity like DENCLUE. Revile of Dimensionality - Dimensionality revile is one of the serious issues looked by high dimensional information. In high dimensional space the focuses are increasingly dissipated or meager and all focuses are practically equidistant from one another.

Bunching methodologies become incapable to investigate the information because of this. Clamor The commotion present in genuine applications frequently conceals the groups to be chosen from bunching calculation and the issue is declined in high dimensional information, where the quantity of mistakes increments directly with dimensionality.

DENCLUE algorithm and OPTICS algorithm comes under the density based clustering technique, where as CLIQUE algorithm comes under the Grid based clustering technique. Clustering in High Dimensional Non-Linear data spaces is a recurrent problem in many domains. It affects time complexity, space complexity, Data Size Adaptability and Precision Value of clustering methods.

## II. RELATED WORKS

Bunching is the gathering of comparative information things into groups. Grouping investigation is one of the fundamental expository techniques in information mining; the strategy for bunching calculation will impact the grouping results legitimately. Mythili S and Madhiya talked about the different kinds of calculations like k-implies bunching calculations, and so on and investigates the favorable circumstances and deficiencies of the different calculations. In each sort we can ascertain the separation between every datum article and all bunch focuses in every emphasis, which makes the Competence Rate of grouping isn't high.

Bunching is the unaided order of examples into gatherings. The bunching issue has been tended to in numerous unique situations and by scientists in numerous controls; this mirrors its wide intrigue and handiness as one of the means in exploratory information examination. M.N. Murty et al., displays a review of example bunching techniques from a factual example acknowledgment point of view, with an objective of giving valuable counsel and references to crucial ideas available to the expansive network of grouping specialists. We present a scientific classification of grouping strategies, and distinguish cross-cutting subjects and late propels.

Sulbha Patil presents cutting as a methodology. Information anonymization is a hot research point. Cutting can deal with high dimensional information and it jam better information utility. Cutting firsts segments properties into segments. Every section contains a subset of qualities. Number ascribes is equivalent to number of segments.

XZhang et.al.. tried different things with certifiable huge informational index in cloud from the point of view of protecting security breaks and to accomplish high level of Data Size Adaptability and Competence Rate.

Manuscript published on 30 September 2019

\* Correspondence Author

**R.Nandhakumar**, Assistant Professor, Department of Computer Science, Nallamuthu Gounder Mahalingam College, Pollachi-642001, India

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

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](https://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

They proposed a closeness protection model and an adaptable two stage bunching approach dependent on MapReduce performing information parallel calculation in cloud to address the issue of security.

M. Suriyapriya and A. Joicy propose a plan which is versatile to replay assaults. In this plan utilizing Secure Hash calculation for confirmation reason, SHA is one of the few cryptographic hash capacities, regularly used to check that a document has been unaltered. The Paillier crypto framework is a probabilistic hilter kilter calculation for open key cryptography. Paillier calculation use for Creation of access strategy, record getting to and document reestablishing process. Cloud have huge extra room for putting away tremendous measure of information. Information proprietor re-appropriate their information substance on cloud server. Cloud server can have tremendous extra room.

### III. PROPOSED MODEL

Bunching or information gathering is the key procedure of the information mining. It is an unaided learning task where one looks to recognize a limited arrangement of classifications named bunches to portray the information. The gathering of information into bunches depends on the guideline of augmenting the intra class similitude and limiting the entomb class likeness.

#### 3.1. Grouping Technique utilized

The gathering of information into groups depends on the standard of amplifying the intra class likeness and limiting the bury class similitude. A decent bunching strategy will create excellent groups with high intra-class closeness - Similar to each other inside a similar bunch low between class likeness. The nature of a bunching strategy is likewise estimated by its capacity to find a few or the majority of the shrouded examples.

The target of the bunching procedure is to decide the inborn gathering in a lot of unlabeled information. The comparability between information articles can be estimated with the forced separation esteems. Determining the separation measures for the high dimensional information is winding up insignificant in light of the fact that it holds various information esteems in their relating traits.

Information mining permits removing information from the colossal data and changing that information into a sensible and significant structure for moreover use. Information mining is a principal task during the time spent taking in disclosure from huge data. Information mining is a development system that is helpful to mine the understandable learning, beforehand obscure, data from huge measure of information put away in different arrangements, with the targets of improving the choice of organizations, associations where the information would be gathered.

#### 3.3 High-Dimensional Data in Knowledge Discovery Database

Bunching high dimensional information in a quality articulation microarray informational index, there could be tens or many measurements, every one of which relates to

an exploratory condition. Revile Dimensionality is a free method for talking about information partition in high dimensional space. The multifaceted nature of many existing information mining calculations is exponential as for the quantity of measurements. Each gathering is a dataset to such an extent that the likeness among the information inside the gathering is expanded and the similitude in outside gathering is limited.

#### 3.4 Density based bunching calculation

Thickness based grouping calculations are prevalent in the utilizations of information mining. These methodologies utilize a nearby bunch foundation and characterize groups as the areas in the information space of higher thickness contrasted with the locales of clamor focuses or fringe focuses. Thickness based grouping calculations utilizing the idea of DBSCAN, can discover bunches of subjective size and shape. Thickness based grouping can be viewed as a non-parametric methodology, where bunches are displayed as territories of high thickness. Inner circle is the principal matrix based subspace grouping approach intended for high dimensional information. It identifies subspaces of the most elevated dimensionalities..

### IV. RESULT AND DISCUSSION

M-DENCLUE deals with two phases as pre-handling stage and grouping stage. In pre-handling step, it makes a framework for the information by separating the negligible jumping hyper-square shape into d-dimensional hyper-square shapes with edge length  $2\sigma$ . In the grouping stage, M-DENCLUE relates an "impact work" with every datum point and the general thickness of the dataset is demonstrated as the total of impact capacities related with each point. The subsequent general thickness capacity will have nearby crests, i.e., neighborhood thickness maxima, and these neighborhood pinnacles can be utilized to characterize groups.

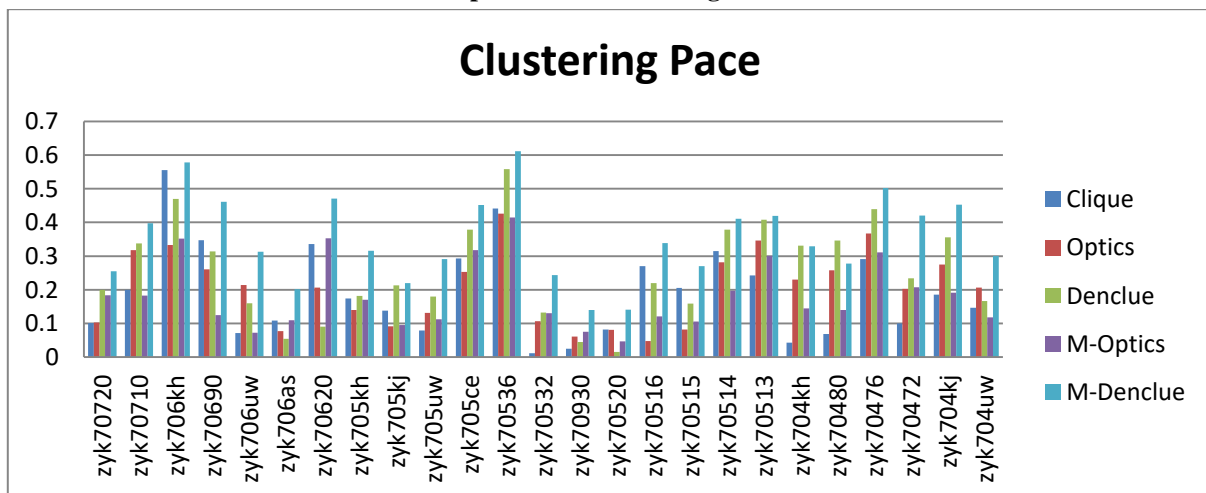
M-DENCLUE uses impact capacities. Impact of every datum point can be demonstrated as scientific capacity. The subsequent capacity is called Influence Function. Impact capacity shows the effect of information point inside its neighborhood.

The M-OPTICS calculation makes a requesting of the articles in a database, M-OPTICS also putting away the center separation and a reasonable reachability separation for each item. A calculation was proposed to concentrate bunches dependent on the requesting data created by M-OPTICS and once the request and the reachability separations are figured, we can extricate the groups for any bunching separation. The work is illustrated through graphs with the help of - DNA microarray Data.

**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

**Experimental Clustering Pace**

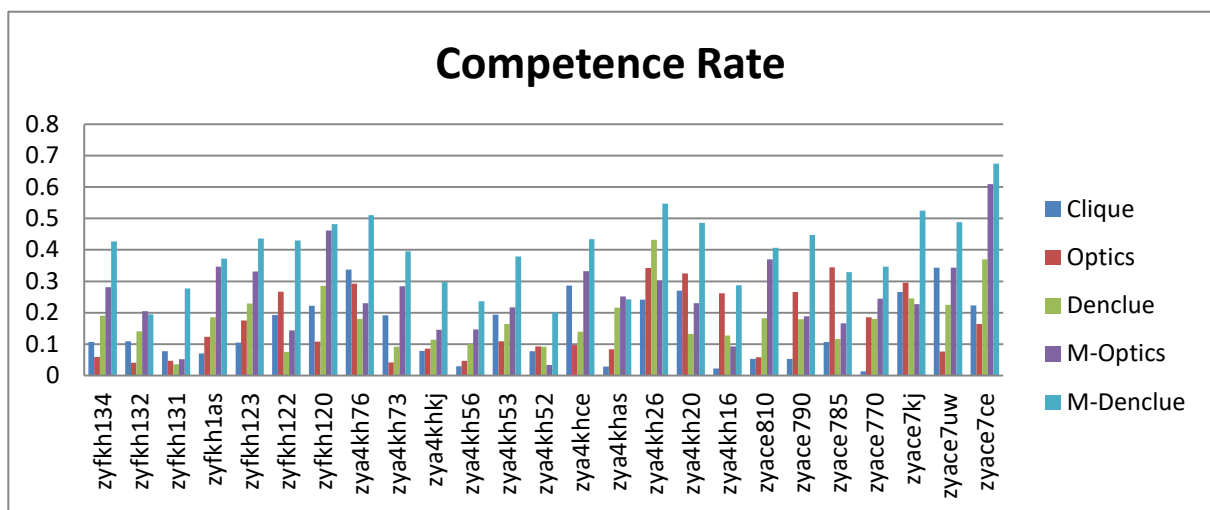


**Experimental Inference-Competence Rate**

## Clustering High Dimensional Non-Linear Data with Denclue, Optics and Clique Algorithms

	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

**Experimental Competence Rate**



### V. CONCLUSION AND FUTURE WORK

DENCLU and OPTICS are density based clustering technique, where as CLIQUE comes under the grid-based

clustering technique. Comparing all those finally conclude that DENCLUE is the best one. Since day to day life changes

with digital world, to group particular data. DENCLUE helps in reducing noise.

From experimental results it has been found that large and dense data needs higher computational power. In future the problems encountered in the existing methods can be overcome by developing a hybrid based density algorithm.

## REFERENCE

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 et al., "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.
4. Gosain Anjana & Chugh Nikita, "Privacy Preservation in Big Data", International Journal of Computer Application, Vol. 100 No.17 August 2014.
5. Hajar Rehioui, Abdellah Idrissi et al., "DENCLUE-IM: A New Approach for Big Data Clustering", The 7th International Conference on Ambient Systems, Networks and Technologies (ANT 2016), ScienceDirect, Procedia Computer Science 83 (2016) 560 – 567.
6. Harsh Shah, Karan Napanda et al., "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 et al., "On Density-Based Data Streams Clustering Algorithms: A Survey", Journal of Computer Science and Technology 29(1): 116{141 Jan. 2014.

## AUTHORS PROFILE



**R. Nandhakumar**, Assistant Professor in Computer Science, Nallamuthu Gounder Mahalingam College, undergoing Ph.D in Data Mining. Published various papers under reputed journals. He is the coordinator in various activates in college like NAAC, ISO, etc..



**Dr. Antony Selvadoss Thanamani**, Associate Professor and Head, Research Department of Computer Science, He is Research supervisor for Ph.D. degree in Computer Science in the Bharathiar University, Dravidian University, etc.. He established Common Research centre, E-content studio, ISBN Nodal Agency at NGM College, Pollachi. He is Advisory Committee Member in National Conference on Advanced Computing. Editor in International Journal of Advanced Scientific Research, India . Editorial Board Member in International Journal of Advanced Research in Computer and Communication Engineering, India.