



# Stream Clustering With Hierarchical Micro-Clusters and Seed Clusters on Density Extraction

Archana<sup>1</sup>, T S Bhagavath Singh<sup>2</sup>

Student, Dept of ISE, RNSIT<sup>1</sup>

Associate Professor, Dept of ISE, RNSIT<sup>2</sup>

**Abstract:** As an ever increasing number of users deliver gushing information, clustering data streams has turned into an essential strategy for information and learning designing. A typical approach is to abridge the data streams progressively with an online procedure into an extensive number of alleged smaller scale bunches (micro-clusters). Micro-clusters are representatives for set of similar data points and are created using a single pass over the data. A conventional clustering algorithm is used in a second offline step to re-cluster the micro-clusters into final clusters sometimes referred to as macro-clusters. This paper depicts Novel Selection, is applied to the medical datasets which has many attributes. In the online stage for the selected disease name in the dataset micro-clusters are formed whereas in offline stage the doctor name is chosen for the selected disease, so that the macro-clusters are formed. This is done by the concept of shared density between the clusters i.e., which are similar to selected attributes, so the large number of smaller clusters will be created. Graph is plotted for both the clusters and also for the accuracy. The clustering quality will be increased by using shared density concept.

**Index Terms:** Data mining, data stream clustering, density-based clustering.

## 1. INTRODUCTION

Clustering data streams has turned into a vital method for information and learning designing. A data stream is a requested and conceivably unbounded succession of data focuses. Such surges of continually arriving data are produced for some sorts of utilizations and incorporate GPS information from advanced mobile phones, web click-stream information, PC arrange observing information, media transmission association information, readings from sensor nets, stock cites, and so forth.

Data stream clustering is commonly done as a two-arrange prepare with an online part which condenses the data into numerous micro-cluster groups or network cells and after that, in a disconnected handle, these micro-clusters (cells) are re-clustered/combined into fewer final clusters. Most papers propose to utilize an (occasionally somewhat adjusted) existing ordinary grouping calculation (e.g., weighted k-implies in CluStream) where the micro-clusters are utilized as pseudo focuses. Another approach utilized as a part of DenStream is to utilize reach ability where every single micro-cluster which are less than a given separation from each other are connected together to shape clusters. Matrix based calculations regularly consolidate nearby thick matrix cells to frame bigger cluster.

Current re-clustering approaches totally overlook the information thickness in the region between the micro-clusters (matrix cells) and in this way may join micro-clusters (cells) which are near one another yet in the meantime isolated by a little territory of low thickness. To address this issue, Tu and Chen [9] acquainted an augmentation with the matrix based D-Stream calculation

[7] in view of the idea of fascination between contiguous matrices cells and demonstrated its adequacy.

In this paper, the micro-clusters are created by selecting the name of the disease once the main clusters are formed, the graph for main clusters are plotted based on index value. Then by selecting the name of the doctor, micro-clusters are re-clustered to form macro-clusters. Then the graph for the dense cluster is also formed. At the end the accuracy graph is also plotted in this project.

## 2. RELATED WORK

Density-based clustering is one of the well-researched areas. The prototypical density-based clustering approach includes DBSCAN [10] and several of its improvements. DBSCAN gauges the thickness around every data point by including the quantity of points a client determined eps-neighborhood and applies client indicated edges to recognize center, outskirts and noise points. In a next step, center points are joined into a group on the off chance that they are thickness reachable (i.e., there is a chain of center points where one falls inside the eps-neighborhood of the following). At long last, clusters are assigned with border points. Different methodologies depend on kernel density estimation (e.g., DENCLUE [8]) or utilize shared closest neighbors (e.g., SNN [6], CHAMELEON [9]).

Be that as it may, these algorithms were not created because of data streams. An data stream is a requested and conceivably unbounded grouping of data points  $X = (x_1,$



$x_2, x_3, \dots, x_i$ ). It is impractical to for all time store every one of the data in the stream which suggests that rehashed arbitrary access to the data is infeasible. Likewise, data streams display idea float after some time where the position as well as state of group changes, and new clusters may show up or existing clusters vanish. This makes the utilization of existing clustering algorithms troublesome. Data stream clustering algorithms constrain data access to a solitary disregard the data and adjust to idea drift. In the course of the most recent 10 years numerous algorithms for clustering data streams have been proposed [5], [6], [8]. Most data stream clustering algorithms utilize a two-stage on the web/disconnected approach [4]:

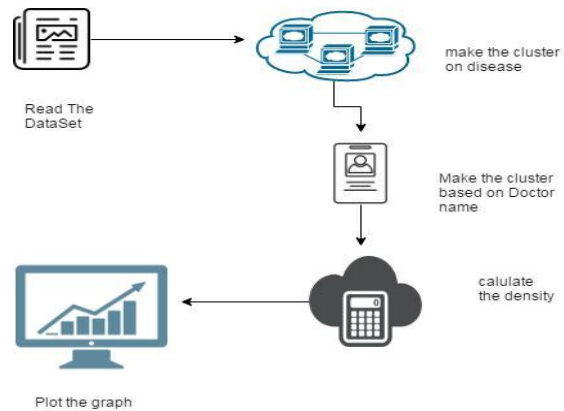


Figure.1. System Architecture

**1) Online:** Summarize the information utilizing an arrangement of  $k'$  micro-clusters sorted out in a space-effective information structure which additionally empowers quick query. Micro-clusters are delegates for sets of comparative information focuses and are made utilizing a solitary disregard the information (regularly continuously when the information stream arrives). Micro-clusters are ordinarily represented by cluster points and extra insights as weight (thickness) and scattering (fluctuation). Each new data indicate is relegated its nearest (regarding a similar function) micro-cluster. A few algorithms utilize a grid rather and non-empty grid cell represents to micro-clusters (e.g., [8], [9]). On the off chance that another data point cannot be allotted a current micro-cluster, another micro-cluster is made. The algorithm may likewise play out some housekeeping (combining or erasing micro-clusters) to keep the quantity of micro-clusters at a sensible size or to remove noise or data obsolete because of concept drift.

**2) Offline:** When the client or the application requires a grouping, the  $k'$  micro-clusters are re-clustered into  $k$  ( $k \ll k'$ ) last groups now and again referred to as macro-clusters. Since the disconnected part is normally not respected time basic, most specialists just express that they utilize a customary bunching algorithm (typically  $k$ -means or a variety of DBSCAN [10]) by in regards to the micro-cluster focus positions as pseudo-foci. The algorithms are regularly altered to consider the weight of micro-clusters.

**3. PROPOSED WORK**

**3.1 Architecture Diagram:**

In this project we have used the medical data sets for clustering. Figure.1. shows the architecture of the project. The data sets contain the medical data; it has many attributes like name of the patient, name of the disease, doctor who handles the particular disease etc. Once the data is read from the dataset, cluster is formed based on the selected disease name. Then make a cluster based on the doctor name. The shared density between two micro-cluster is estimated and based on shared density the shared density graph is plotted. The clusters are formed using novel selection approach.

First based on the decease with respect to doctors (who are treating that disease) will be framed as the first cluster. Second based on the dynamic selection of the attributes micro clusters will be framed. These micro clusters can be in any form but based on the second level of attributes selection. Once first level and second level data points (clusters) framed, sandwich model applied with fine grained grid as shown in Figure.2. The algorithm checks which grid is having highest combinational data points. For example Second level data points = 34 and first level data points are 153 then density would be  $(34/153) * 100$ .

**3.2 Shared Density between micro-clusters**

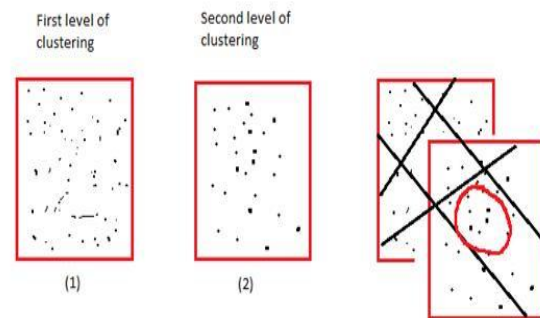


Figure.2. Shared Density between micro-clusters

**3.3 Algorithms**

**Algorithm 1:** Level1 → Online algorithm

**INPUT:** Read data set with Attribute

**OUTPUT:** Reduce data with selective attribute dense cluster

$\lambda = m$  // size of the dataset

1. Data:
  - // Pn is Parameters
  - // Δ is a whole data set
2. For 1 to n of (P)
3. Start: SELECT (D (diseases))
4. = FETCH [1-m-1]
- // m is the dataset
5. If (Δt (m) == disease)
6. UPDATE [MAKE (Cluster)]



7. Cu= ..... (\*)  
**Algorithm 2:** Level 2 -> Offline algorithm  
**INPUT:** Cu (\*) //first level cluster  
**OUTPUT:** // dense cluster  
 // t -> doctor specific attribute  
 1. ≡ Scan (Cu)  
 // equality symbol  
 2. For n in Dc  
 3. Start:  $\Delta\text{Data} = \text{GRT} (\text{Dc } \Delta t_{12})$   
 4. Density  $\approx$  Scatter (Data)  
 5. PLOT

5. CONCLUSION AND FUTURE WORK

In this project the shared density graph together with the algorithms needed to maintain the graph in the online component of a data stream mining algorithm is introduced. Project also shows that shared-density re-clustering already performs extremely well when the online data stream clustering component is set to produce a small number of large MCs. By selecting two attributes that is name of disease and doctor, the micro-clusters are formed using re-clustering concept and graph is plotted for the same. This improves the accuracy and saves the memory.

The algorithms for formation of clusters and re-clustering are written in Algorithm 1 and Algorithm 2 respectively.

In Future work clustering and re-clustering is done using many attributes from the data sets and density clusters with respect to security.

4. RESULTS

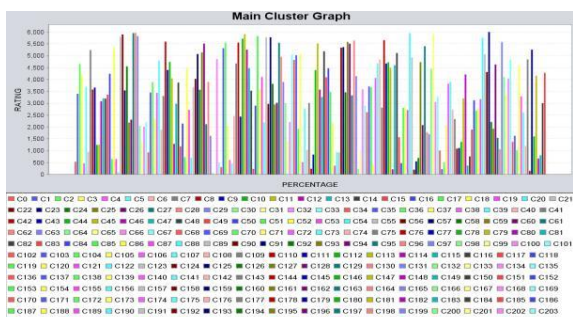


Figure.3. Graph for first level

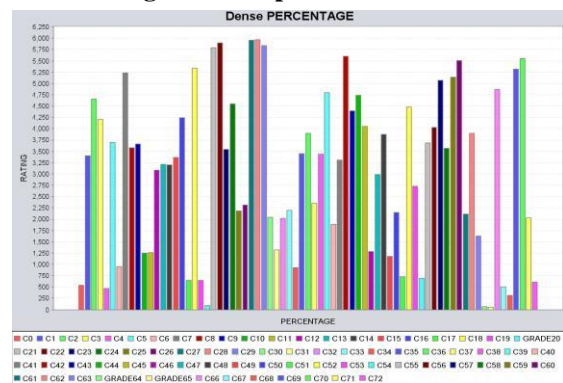


Figure.4. Graph for second level clustering

After formation of clusters the density graph is plotted for each cluster. Graph for first level clustering, second level clustering and accuracy is shown in Figure.3, Figure.4 and Figure.5 respectively

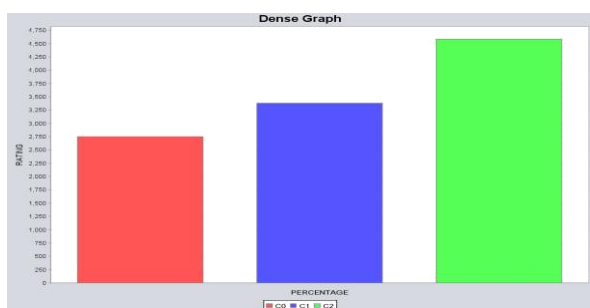


Figure.5. Accuracy graph

REFERENCES

- [1] Michael Hahsler and Matthew Bolanos, "Clustering Data Streams Based on Shared Density Between Micro- Clusters" in IEEE transaction on Knowledge and Data engineering, 2016.
- [2] G. Bifet, G. de Francisci Morales, J. Read, Holmes, and Pfahringer, "Efficient online evaluation of big data stream classifiers," in Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ser. KDD '15.ACM, 2015, pp. 59–68.
- [3] M. Hahsler, M. Bolanos, and J. Forrest, stream: Infrastructure for Data Stream Mining, 2015, R package version 1.2-2.
- [4] A. J. A. Silva, E. R. Faria, R. C. Barros, E. Hruschka, A. C. P. L. F. d. Carvalho, and J. a. Gama, "Data stream clustering: survey," ACM Computing Surveys, vol. 46, no. 1, pp. 13:1–13:31, Jul. 2013.
- [5] J. Gama, R. Sebasti-ao, and P. P. Rodrigues, "On evaluating stream learning algorithms," Mach. Learn., vol. 90, pp. 317–346, 2013.
- [6] Amini and T. Y. Wah, "Leaden-stream: A leader density-based clustering algorithm over evolving data stream," Journal of Computer and Communications, vol. 1, no. 5, pp. 26– 31, 2013.
- [7] Isaksson, M. H. Dunham, and M. Hahsler, "Sostream: Self organizing density-based clustering over data stream," in Machine Learning and Data Mining in Pattern Recognition, ser. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 2012, vol. 7376, pp. 264–278.
- [8] M. Hahsler and M. H. Dunham, "Temporal structure learning for clustering massive data streams in real-time," in SIAM Conference on Data Mining (SDM11). SIAM, April 2011, pp. 664–675.
- [9] P. Kranen, I. Assent, C. Baldauf, and T. Seidl, "The clustree: indexing micro-clusters for anytime stream mining," Knowledge and Information Systems, vol. 29, no. 2, pp. 249–272, 2011.
- [10] H. Kremer, P. Kranen, T. Jansen, T. Seidl, A. Bifet, G. Holmes, and B. Pfahringer, "An effective evaluation measure for clustering on evolving data streams," in Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining . ACM, 2011, pp. 868–876.