# A Multicore Approach for DBCAN Data Clustering Using KD Tree

Ms. Alka Shrivastava
alkashriwastava@gmail.com
Bharti College of Engineering & Technology, Durg

Mr. Suman Swarnkar
sumanswarnkar17@gmail.com
Bharti College of Engineering & Technology, Durg

*Abstract—* **The issue of recognizing classes of data points is challenging when the clusters are of distinctive size, shape and density. Most of these issues turn out to be much more critical when the data is of high dimensionality and when it incorporates outliers and noise. In earlier implementation of DBSCAN algorithms some bottlenecks that need to be resolved, i.e. presence of noise in input dataset affects the algorithm performance and increase in outlier, existing parallel implementation of DBSCAN algorithms are based master slave which imposes I/O overhead henceforth there is need of better implementation of DBSCAN. In this paper we have proposed distributed DBSCAN implementation by means of multicore programming, further we have used k-dimension tree as data structure.**

*Keywords—KD,Eps,DCPGS*

## I. INTRODUCTION

Data clustering is a data mining practice that assemblies data into evocative subclasses, termed as clusters, such that it curtails the intra-differences and get the most out of inter-differences of these subclasses. Renowned algorithms for data clustering include DBSCAN, K-medoids, Kmeans, BIRCH, WaveCluster and STING. These algorithms have been castoff in numerousmethodical areas for example satellite image segmentation, noise sieving and outlier detection, unsupervised text document clustering, and clustering of bioinformatics data. Earlier data clustering algorithms have been unevenlyclassified into four classes: partitioning-based, hierarchybased, grid-based, and density-based.

DBSCAN (Density Based Spatial Clustering of Applications with Noise) is a density centered clustering algorithm. The concept behind DBSCAN algorithm is that for every data point in a cluster, the vicinity within a specified radius (eps) has to comprise at least a minimum number of points (minpts), i.e. the density of the neighborhood has to outstrip some threshold.

DBSCAN algorithm was projected by Martin Ester and others in 1996. DBSCAN algorithm entails only two factors: Eps & MinPts. Eps frights with a randompreliminary point that has not been stayed. This point's neighborhood is salvaged, and if it comprisesadequately many points, a cluster is underway. MinPts is the minutest number of points needs to form a cluster. Eps is aprevalent parameter for DBSCAN algorithm but shrewd the value of Eps is time consuming and difficult. For calculating Eps.

## DBSCAN (D, eps, MinPts)

C = 0
for each unvisited point p in dataset D

mark P as visited
NeighborPts = regionQuery (P, eps)
If sizeof(NeighborPts) < MinPts
mark P as NOISE
else
C= next cluster
expandCluster (P, NeighborPts, C, eps, MinPts)
add P to cluster C
for each point P' in NeighborPts
if P' is not visited
mark P' as visited
NeighborPts' = regionQuery (P', eps)
If sizeof (NeighborPts') >= MinPts
NeighborPts = Neighborpts joined with NeighborPts'
if P' is not yet member of any cluster
add P' to cluster C
regionQuery (P, eps)
return all points within P's eps-neighbourhood

## KD Tree (K Dimensional Trees)

In this paper proposed algorithm uses KD Tree data structure while implementing multicore DBSCAN. Splitting planes perpendicular to Coordinate Axes.Useful in Nearest Neighbor Search.KD tree reduces the Overall Time Complexity to O(log n) which has been used in many clustering algorithms and other domains.

**Input:**
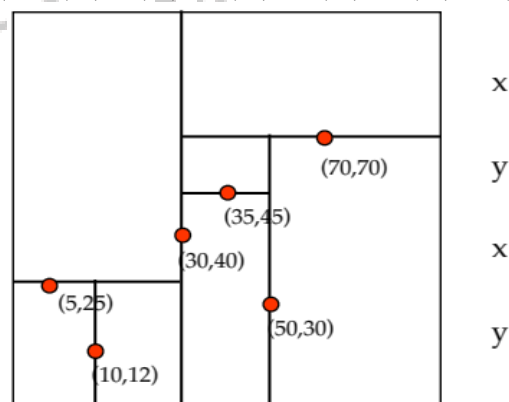Insert: (30,40), (5,25), (10,12), (70,70), (50,30), (35,45)



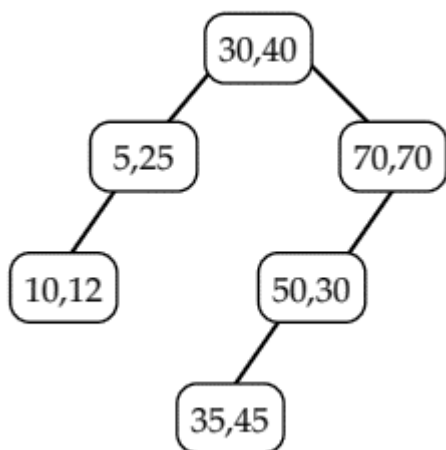Fig.-1 Splitting planes perpendicular to Coordinate Axes

Fig.-2 kD Tree Constructed

Further in this paper in section-II we have discussed different literature and given tabular association among different literature, in section-IIIwe will discuss advantage and disadvantages of different implementation, in section IV proposed methodology discussed, in section V implementation section at last we will conclude our research.

## II. LITERATURE SURVEY

I this paper we have gone through different literature some of them are briefed below:

Xiangliang Zhang et. al. said that proposed algorithm STRAP aims at clustering data streams with evolving data distributions. STRAP confronts the arriving items to the current AP model, storing the outliers in a reservoir and monitoring the ratio of outliers using the PH change point detection test. Upon triggering the PH test, the clustering model is rebuilt from the current one and the reservoir using WAP \. The key issue here was to build the change indicator, monitored by the PH test, in order to preserve the computational cost vs accuracy tradeoff. In this paper, we monitored the ratio of outliers over a sliding window and adapted the change detection threshold in real-time. The proposed approach STRAP was theoretically analyzed on guaranteeing acceptable distortion loss when exemplars slightly drift from the already selected ones, on consuming small amount of memory with little variation, and on requiring acceptable computing time that depends on the complexity of underlying distribution. The performance of STRAP in clustering quality and efficiency is empirically validated on KDD'99 benchmark problem and the URLs stream [IEEE 2013].

Jieming Shi et. al. studied for the first time the problem of Densitybased Clustering Places in Geo-Social Networks (DCPGS). Our clustering model extends the density-based clustering paradigm to consider both the spatial and social distances between places. We defined a new measure for the social distance between places, considering the social ties between users that visit them. Also said that our measure is shown to be more effective and efficient to compute, compared more complex ones based on node-to-node graph proximity. We analyzed the effectiveness of DCPGS via case studies and demonstrated that DCPGS can discover clusters with interesting properties (i.e., barrier-based splitting, spatially loose clusters, clusters with fuzzy

boundaries), which cannot be found by merely using spatial clustering. Besides, we designed two evaluation measures to quantitatively evaluate the social quality of clusters detected by DCPGS or competitors, called social entropy and community score, which also confirm that DCPGS is more effective than alternative approaches.[IEEE 2017].

Kamran Khan et. al. said that Data Mining is all about data analysis techniques. It is useful for extracting hidden and interesting patterns from large datasets. Clustering techniques are important when it comes to extracting knowledge from large amount of spatial data collected from various applications including GIS, satellite images, X-ray crystallography, remote sensing and environmental assessment and planning etc. To extract useful pattern from these complex data sources several popular spatial data clustering techniques have been proposed. DBSCAN (Density Based Spatial Clustering of Applications with Noise) is a pioneer density based algorithm. It can discover clusters of any arbitrary shape and size in databases containing even noise and outliers. DBSCAN however are known to have a number of problems such as: (a) it requires user's input to specify parameter values for executing the algorithm; (b) it is prone to dilemma in deciding meaningful clusters from datasets with varying densities; (c) and it incurs certain computational complexity. Many researchers attempted to enhance the basic DBSCAN algorithm, in order to overcome these drawbacks, such as VDBSCAN, FDBSCAN, DD_DBSCAN, and IDBSCAN. In this study, we survey over different variations of DBSCAN algorithms that were proposed so far. These variations are critically evaluated and their limitations are also listed [IEEE 2014].

Chetan Dharni et. al. Analyzed that the varied density clusters and time complexity are the main problems which has been improved in many DBSCAN variations. The time complexity is reduced to O(N) with the help of indexing in DBSCAN algorithm. The problem of varied density and increasing dimensionality effect the performance of DBSCAN algorithm. Still it needs more improvement [IJCST 2013].

Fang Huang et. al. said that Density-based spatial clustering of applications with noise (DBSCAN) is a density-based clustering algorithm that has the characteristics of being able to discover clusters of any shape, effectively distinguishing noise points and naturally supporting spatial databases. DBSCAN has been widely used in the field of spatial data mining. This paper studies the parallelization design and realization of the DBSCAN algorithm based on the Spark platform, and solves the following problems that arise when computing macro data: the requirement of a great deal of calculation using the single-node algorithm; the low level of resource-utilization with the multi-node algorithm; the large time consumption; and the lack of instantaneity. The experimental results indicate that the proposed parallel algorithm design is able to achieve more stable speedup at an increased involved spatial data scale[MDPI 2017].

| S.No. | Author/Title/Year/Publication | Method Used | Description |
|---|---|---|---|
| | | | |

| 1. | Xiangliang Zhang et. Al./Data Stream Clustering with Affinity Propagation//HAL-Inria2014 | Affinity Propagation (AP) algorithm | Clustering data streams with evolving data distributions. STRAP confronts the arriving items to the current AP model. |
|----|----|----|----|
| 2. | Jieming Shi et.l./Density-based Place Clustering in Geo-Social Networks/2014/SIGMOD | DBSCAN | Studied for the first time the problem of Density based Clustering Places in Geo-Social Networks (DCPGS). Clustering model extends the density-based clustering paradigm to consider both the spatial and social distances between places. |
| 3. | Saif Ur Rehman et. Al./DBSCAN: Past, Present and Future/2016/IEEE | DBSCAN | Presented the summary information of the different enhancement of density-based clustering algorithm called the DBSCAN. |
| 4. | SANJAY CHAKRABORTY et. Al./Analysis and Study of Incremental DBSCAN Clustering Algorithm/2014/ARXIV | DBSCAN | Propose an improve DBSCAN clustering approach which provides better and fastest result compare to the existing DBSCAN clustering algorithm up to some certain point of change in the original database. |
| 5. | Xiangliang Zhang et. Al./ Data Stream Clustering With Affinity Propagation/IEEE 2016 | STRAP | Presented STRAP algorithm combines AP with a statistical change point detection test; the clustering model is rebuilt whenever the test detects a change in the underlying data distribution. Besides the validation on two benchmark data sets, the presented algorithm is validated on a real-world application. |
| 6. | Ali Seyed Shirkhorshidi et. Al/Big Data Clustering: A Review/Springer 2014 | MapReduce, Parallel Clustering | This study is aimed to review the trend and progress of clustering algorithms to cope with big data challenges from very first proposed algorithms until today's novel solutions. |
| 7. | Zexuan Ji et. Al./Interval-valued possibilistic fuzzy C-means clustering algorithm/Elsevier 2013 | Fuzzy C-means clustering | Author compared the proposed algorithm with five fuzzy clustering approaches, including the FCM, PCM, PFCM, IFCM and IPCM, on two-dimensional Gaussian data sets and four multi-dimensional benchmark data sets. Also applied these clustering techniques to segment the brain magnetic resonance images and natural images. Results of this paper show that the proposed IPFCM algorithm is more robust to outliers and initializations and can produce more accurate clustering results. |

## III. FINDINGS

There are several implementations of DBSCAN algorithms.

El-Sonbaty et. al. has proposed a efficient density based algorithm which require preprocessing over input dataset to improve the performance of clustering, dataset partitioning used in pre-processing step. Next stage is merge stage for merging the nearest dense region which can be done by applying the DBSCAN separately on each partition in order to achieve final cluster.

**Advantage**:
1. Algorithm is scalable.
2. In future parallel version of algorithm can be implemented which will improve the performance of algorithm.

**Disadvantage:**
1. Before clustering preprocessing over input dataset is required.
2. I/O load may be augmented.

Adriano Moreira et. al. has proposed DBSCAN implementation which starts by identifying the k nearest neighbours of each point and identify the farthest k nearest neighbour (in terms of Euclidean distance) . The average of all this distance is then calculated. After that, for each point of the

dataset the algorithm identifies the directly density-reachable points (using the Eps threshold provided by the user) and classifies the points into core or border points.

**Advantage:**

1. Does not entail priori specification of number of clusters.
2. Capable to recognize noise data while clustering.

**Disadvantage:**

1. Algorithm miss the mark in case of varying density clusters.
2. Performance degrades in the case of high dimension input dataset.

Duan et.al. has proposed Local Density Based Spatial Clustering Algorithm With Noise having time complexity is O(n). In the proposed algorithm, LOF (local outlier factor) is used to detect the noise and has overawed the problem of global density parameter.

**Advantage**:

1. Easier for the user to prefer the suitable parameters.

**Disadvantage**:

1. Cluster analysis is difficult.Bing Liu has proposed Fast Density Based Clustering Algorithm for Large Databases in which objects are organized by certain dimensional coordinates. This algorithmuses global Eps parameter. If the value of Eps is less, then the few or single cluster containing all objects is formed and if the value of Eps is high, many trifling clusters are produced.

**Advantage:**

1. Time complexity condensed.

**Disadvantage:**

1. The problem of diverse density clusters is not analyzed.

Zhang et. al has proposed A Linear DBSCAN Algorithm Based on LSH (Locality-Sensitive Hashing). The improvement of using LSH is that it decreases the time complexity and the scale of data.The time complexity of algorithm is O(NlogN).

**Advantage:**

1. Comparative original DBSCAN time complexity reduces.
2. Outlierseradicated using LSH and nearest neighbor points are obtained.

**Disadvantage**:

1. Preprocessing over input dataset is required.

Pathway et. al. proposed A New Scalable Parallel DBSCAN Algorithm Using Disjoint-Set Data Structure To build clusters, a tree based bottom-up method is used. The disjoint-set data structure is cast-off to break the data access order and to achieve the merging efficiently.

**Advantage**:

1. The master-slave method speeds up the process.

**Disadvantage**:

1. Upsurges the I/O load and also effects the cost.

## IV. PROPOSED METHODOLOGY

To overcome the different bottlenecks of sequential DBSCAN algorithm we have proposed multicore programming based DBSCAN algorithm which will use K-Dimension tree as data structure for parallel implementation.
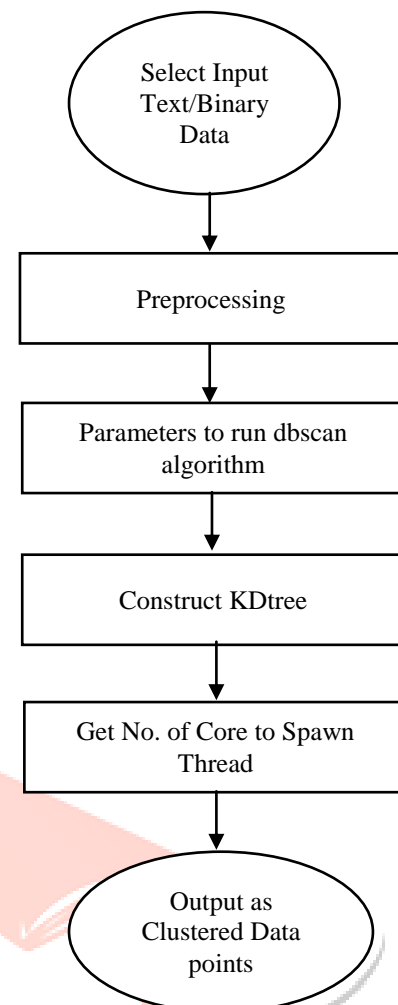


Fig.-3 Proposed System Architecture

| **Proposed Algorithm** |
|---|
| Step-1 Select input text or binary dataset. |
| Step-2 Set DBSCAN parameter eps and minpoints. |
| Step-3 Remove noise. |
| Step-4 Construct KD Tree. |
| Step-5 Construct Minimum spanning tree using REM's algorithm. |
| Step-6 Set #pragma to each node. |
| Step-7 Union of cluster node. |
| Step-8 Merge all visited and non-visited node of KD tree. |

| **KD Tree Construction** |
|---|
| Step-1. if P contains only one point |
| Step-2. then return a leaf storing this point |
| Step-3. else if depth is even |
| Step-4. then Split P with a vertical line ` through the median x-coordinate into P1 (left of or on `) and P2 (right of `) |
| Step-5. else Split P with a horizontal line ` through the median y-coordinate into P1 (below or on `) and P2 (above `) |
| Step-6. vleft ← BuildKdTree(P1,depth+1) |
| Step-7. vright ← BuildKdTree(P2,depth+1) |
| Step-8. Create a node v storing `, make vleft the left child of v, and make vright the right child of v. |
| Step-9. return v |

### V. Result and Discussion

For implementation of proposed algorithm we have used OpenMP, gcc compiler in Ubuntu 16.0 operating system. OpenMp is an open source library for multicore programming. We have implemented sequential DBSCAN and proposed DBSCAN and tested on different dataset. Proposed implementation takes following input:

1. Input file name
2. Input parameter of DBSCAN, min points to form a cluster.
3. Input parameter of DBSCAN, radius or threshold on neighborhoods retrieved.
4. Clustering results file name.
5. Number of threads

**Example**

./multicore_dbscan -i filename -b -m minpts -e epsilon -o output -t threads
./multicore_dbscan -i example.dat –b -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i example1.dat–b -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus5k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus7.5k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus10k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus12.5k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus15k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus20k.txt -m 5 -e 30 -t 8 -o output1
./multicore_dbscan -i clus25k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus30k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus35k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus40k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus45k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus50k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus60k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus70k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus75k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus80k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus90k.txt -m 5 -e 30 -t 8 -o output
./multicore_dbscan -i clus100k.txt -m 5 -e 30 -t 8 -o output

| S. No. | Dataset size | Execution time in sec | |
|--------|--------------|-----------|----------|
| | | Sequential | Proposed |
| 1. | 54 | 0.0011 | 0.000990349 |
| 2. | 106 | 0.052 | 0.000991960 |
| 3. | 5000 | 0.8837 | 0.000993491 |
| 4. | 7500 | 2.2899 | 0.00117969 |
| 5. | 10000 | 5.5712 | 0.001884051 |
| 6. | 12500 | 5.5556 | 0.00372282 |
| 7. | 15000 | 7.9576 | 0.003791021 |
| 8. | 20000 | 26.6375 | 0.0038201 |
| 9. | 25000 | 32.8844 | 0.00388885 |
| 10. | 30000 | 32.1504 | 0.00554149 |
| 11. | 35000 | 45.7200 | 0.0056283 |
| 12. | 40000 | 114.8900 | 0.005866381 |
| 13. | 45000 | 140.6387 | 0.00627656 |
| 14. | 50000 | 143.8052 | 0.006368 |
| 15. | 60000 | 154.0679 | 0.0064895 |
| 16. | 70000 | 209.5745 | 0.006605 |
| 17. | 75000 | 289.8763 | 0.007475 |
| 18. | 80000 | 527.0518 | 0.007840121 |
| 19. | 90000 | 630.1514 | 0.00858002 |
| 20. | 100000 | 702.6091 | 0.124919 |



Fig.-5 Snippet of Input dataset



Fig.-6 Comparison based on Execution Time



Fig.-7 Proposed System Output

Fig.-7 Shows the snippet of our proposed algorithm implementation, which shows number of point noise and execution time I/O time.



Fig.-8 Proposed System Output File

## VI. CONCLUSION

Data clustering is a renowned technique in numerous areas of computer science and associated domains.DBSCAN is the illustrious method density centered data clustering method. Several implementation of DBSCAN has been done which made improvement in different aspects of DBSCAN and improved the performance. Existing parallel implementations of DBSCAN clustering algorithm espouse a master-slave tactic which can simply grounds an unbalanced workload and therefore result in low parallel competence. Earlier techniques are not time efficient. Multicore programming facilitated us to make time efficient DBSCAN algorithm.

## REFERENCES

[1] Xiangliang Zhang et. Al. Data Stream Clustering with Affinity Propagation HAL-Inria 2014

[2] Jieming Shi et.l. Density-based Place Clustering in Geo-Social Networks 2014 SIGMOD

[3] Saif Ur Rehman et. Al. DBSCAN: Past, Present and Future 2016 IEEE

[4] SANJAY CHAKRABORTY et. Al. Analysis and Study of Incremental DBSCAN Clustering Algorithm 2014 ARXIV

[5] Xiangliang Zhang et. Al. Data Stream Clustering With Affinity Propagation IEEE 2016

[6] Ali Seyed Shirkhorshidi et. Al Big Data Clustering: A Review Springer 2014

[7] Zexuan Ji et. Al. Interval-valued possibilistic fuzzy C-means clustering algorithm Elsevier 2013

[8] K.Kameshwaran, K.Malarvizhi Survey on Clustering Techniques in Data Mining IJCSIT 2014

[9] Ilias K. Savvas, and Dimitrios Tselios Parallelizing DBSCAN Algorithm Using MPI IEEE 2016