



A NOVEL ALGORITHM FOR DUAL SIMILARITY CLUSTERS USING MINIMUM SPANNING TREE

S.JOHN PETER¹, S.P.VICTOR²,

1. Assistant Professor 2. Associate Professor

Department of Computer Science and Research Center

St. Xavier's College, Palayamkottai

Tamil Nadu, India.

ABSTRACT

The minimum spanning tree clustering algorithm is capable of detecting clusters with irregular boundaries. In this paper we propose minimum spanning trees based clustering algorithm. The algorithm produces k clusters with center and it also creates a dendrogram for the k clusters. The algorithm works in two phases. The first phase of the algorithm produces subtrees. The second phase converts the subtrees into dendrogram. The key feature of the algorithm is it uses both divisive and agglomerative approaches to find Dual similarity clusters.

Key Words: *Euclidean Minimum Spanning Tree, Clustering, Eccentricity, Center, Hierarchical Clustering, Dendrogram, Subtree.*

1.INTRODUCTION

Given a connected, undirected graph $G = (V, E)$, where V is the set of nodes, E is the set of edges between pairs of nodes, and a weight $w(u, v)$ specifying weight of the edge (u, v) for each edge $(u, v) \in E$. A spanning tree is an acyclic subgraph of a graph G , which contains all vertices from G . The Minimum Spanning Tree (**MST**) of a weighted graph is minimum weight spanning tree of that graph. Several well established **MST** algorithms exist to solve minimum spanning tree problem [17],[12],[13]. The cost of constructing a minimum spanning tree is $O(m \log n)$, where m is the number of edges in the graph and n is the number of vertices. More efficient algorithm for constructing **MSTs** have also been extensively researched [11], [4], [7]. These algorithms promise close to linear time complexity under different assumptions. A Euclidean minimum spanning tree (**EMST**) is a spanning tree of a set of n points in a metric space (\mathbf{E}^n), where the length of an edge is the Euclidean distance between a pair of points in the point set.

The hierarchical clustering approaches are related to graph theoretic clustering. Clustering algorithms using minimal spanning tree takes the advantage of **MST**. The **MST** ignores many possible connections between the data patterns, so

the cost of clustering can be decreased. The **MST** based clustering algorithm is known to be capable of detecting clusters with various shapes and size [20]. Unlike traditional clustering algorithms, the **MST** clustering algorithm does not assume a spherical shapes structure of the underlying data. The **EMST** clustering algorithm [16], [20] uses the Euclidean minimum spanning tree of a graph to produce the structure of point clusters in the n -dimensional Euclidean space. Clusters are detected to achieve some measures of optimality, such as minimum intra-cluster distance or maximum inter-cluster distance [2]. The **EMST** algorithm has been widely used in practice.

Clustering by minimal spanning tree can be viewed as a hierarchical clustering algorithm which follows a divisive approach. Using this method firstly **MST** is constructed for a given input. There are different methods to produce group of clusters. If the number of clusters k is given in advance, the simplest way to obtain k clusters is to sort the edges of minimum spanning tree in descending order of their weights and remove edges with first $k-1$ heaviest weights [2], [19].



Geometric notion of centrality are closely linked to facility location problem. The distance matrix D can be computed rather efficiently using Dijkstra's algorithm with time complexity $O(|V|^2 \ln |V|)$ [18].

The *eccentricity* of a vertex x in G and radius $\rho(G)$, respectively are defined as

$$e(x) = \max_{y \in V} d(x, y) \quad \text{and} \quad \rho(G) = \min_{x \in V} e(x)$$

The *center* of G is the set

$$C(G) = \{x \in V \mid e(x) = \rho(G)\}$$

$C(G)$ is the center to the "emergency facility location problem" which always contains a single block of G . The length of the longest path in the graph is called *diameter* of the graph G . We can define diameter $D(G)$ as

$$D(G) = \max_{x \in V} e(x)$$

The *diameter set* of G is

$$Dia(G) = \{x \in V \mid e(x) = D(G)\}$$

All existing clustering algorithms require a number of parameters as their inputs and these parameters can significantly affect the cluster quality. In this paper we want to avoid experimental methods and advocate the idea of need-specific as opposed to care-specific because users always know the needs of their applications. We believe it is a good idea to allow users to define their desired similarity within a cluster and allow them to have some flexibility to adjust the similarity if the adjustment is needed. Our algorithm produces clusters of n -dimensional points with a given cluster number and a naturally approximate intra-cluster distance.

Hierarchical clustering is a sequence of partitions in which each partition is nested into the next in sequence. An Agglomerative algorithm for hierarchical clustering starts with disjoint clustering, which places each of the n objects in an individual cluster [1]. The hierarchical clustering algorithm being employed dictates how the proximity matrix or proximity graph should be interpreted to merge two or more of these trivial clusters, thus nesting the trivial clusters into a second partition. The process is repeated to form a sequence of nested clustering in which the number of clusters decreases as a sequence progresses until a single cluster containing all n objects, called the *conjoint clustering*, remains [1].

An important objective of hierarchical cluster analysis is to provide a picture of data that can easily be interpreted. A picture of a hierarchical

clustering is much easier for a human being to comprehend than is a list of abstract symbols. A *dendrogram* is a special type of tree structure that provides a convenient way to represent hierarchical clustering. A dendrogram consists of layers of nodes, each representing a cluster.

In this paper we have used **EMST** based clustering algorithm to address the issues of undesired clustering structure and unnecessary large number of clusters. The algorithm assumes the number of clusters is given. The algorithm constructs an **EMST** of a point set and removes the inconsistent edges that satisfy the inconsistency measure. The process is repeated to create a hierarchy of clusters until k clusters are obtained. In section 2 we review some of the existing works on graph based algorithm. In Section 3 we propose **HCEMST** algorithm which produces k clusters with a dendrogram. Hence we named these new clusters as *Dual similarity clusters*. Finally in conclusion we summarize the strength of our methods and possible improvements.

2. RELATED WORK.

Clustering by minimal spanning tree can be viewed as a hierarchical clustering algorithm which follows the divisive approach. Clustering algorithms based on minimum and maximum spanning trees were extensively studied. In the mid of 80's, Avis [3] found an $O(n^2 \log^2 n)$ algorithm for the min-max diameter-2 clustering problem. Asano, Bhattacharya, Keil and Yao [2] later gave an optimal $O(n \log n)$ algorithm using maximum spanning trees for minimizing the maximum diameter of a bipartition. The problem becomes NP-complete when the number of partitions is beyond two [10]. Asano, Bhattacharya, Keil and Yao also considered the clustering problem in which the goal is to maximize the minimum inter-cluster distance. They gave a k -partition of a point set by removing the $k-1$ longest edges from the minimum spanning tree constructed from that point set [2]. The identification of inconsistent edges causes a problem in the **MST** clustering algorithm. There exist numerous ways to divide clusters successively, but there is not a suitable choice for all cases.

Zahn [20] proposes to construct **MST** of a point set and delete inconsistent edges – the edges, whose weights are significantly larger than the average weight of the nearby edges in the tree.



Zahn's inconsistent measure is defined as follows. Let e denote an edge in the **MST** of the point set, v_1 and v_2 be the end nodes of e , w be the weight of e . A *depth neighborhood* N of an end node v of an edge e defined as a set of all edges that belong to all the path of length d originating from v , excluding the path that include the edge e . Let N_1 and N_2 be the depth d neighborhood of the node v_1 and v_2 . Let \hat{W}_{N_1} be the average weight of edges in N_1 and σN_1 be its standard deviation. Similarly, let \hat{W}_{N_2} be the average weight of edges in N_2 and σN_2 be its standard deviation. The inconsistency measure requires one of the three conditions hold:

1. $w > \hat{W}_{N_1} + c x \sigma N_1$ or $w > \hat{W}_{N_2} + c x \sigma N_2$
2. $w > \max(\hat{W}_{N_1} + c x \sigma N_1, \hat{W}_{N_2} + c x \sigma N_2)$
3. $\frac{w}{\max(c x \sigma N_1, c x \sigma N_2)} > f$

where c and f are preset constants. All the edges of a tree that satisfy the inconsistency measure are considered inconsistent and are removed from the tree. This result in set of disjoint subtrees each represents a separate cluster. Paivinen [15] proposed a Scale Free Minimum Spanning Tree (**SFMST**) clustering algorithm which constructs scale free networks and outputs clusters containing highly connected vertices and those connected to them.

The **MST** clustering algorithm has been widely used in practice. Xu (Ying), Olman and Xu (Dong) [19] use **MST** as multidimensional gene expression data. They point out that **MST**-based clustering algorithm does not assume that data points are grouped around centers or separated by regular geometric curve. Thus the shape of the cluster boundary has little impact on the performance of the algorithm. They described three objective functions and the corresponding cluster algorithm for computing k -partition of spanning tree for predefined $k > 0$. The algorithm simply removes $k-1$ longest edges so that the weight of the subtrees is minimized. The second objective function is defined to minimize the total distance between the center and each data point in the cluster. The algorithm removes first $k-1$ edges from the tree, which creates a k -partitions.

Hierarchical clustering algorithm proposed by S.C.Johnson [9] uses proximity matrix as input data. The algorithm is an agglomerative scheme that erases rows and columns in the proximity matrix as old clusters are merged into new ones.

The algorithm is simplified by assuming no ties in the proximity matrix. Graph based Hierarchical Algorithm was proposed by Hubert [6] using single link and complete link methods. He used threshold graph for formation of hierarchical clustering. An algorithm for single-link hierarchical clustering begins with the minimum spanning tree (**MST**) for $G(\infty)$, which is a proximity graph containing $n(n-1)/2$ edge was proposed by Gower and Ross [8]. Later Hansen and DeLattre [5] proposed another hierarchical algorithm from graph coloring.

3. OUR CLUSTERING ALGORITHM

A tree is a simple structure for representing binary relationship, and any connected components of tree is called *subtree*. Through this **MST** representation, we can convert a multi-dimensional clustering problem to a tree partitioning problem, i.e., finding particular set of tree edges and then cutting them. Representing a set of multi-dimensional data points as simple tree structure will clearly lose some of the inter data relationship. However many clustering algorithm proved that no essential information is lost for the purpose of clustering. This is achieved through rigorous proof that each cluster corresponds to one subtree, which does not overlap the representing subtree of any other cluster. Clustering problem is equivalent to a problem of identifying these subtrees through solving a tree partitioning problem. The inherent cluster structure of a point set in a metric space is closely related to how objects or concepts are embedded in the point set. In practice, the approximate number of embedded objects can sometimes be acquired with the help of domain experts. Other times this information is hidden and unavailable to the clustering algorithm. In this section we present clustering algorithm which produce k clusters, with center of each cluster. We also present another algorithm to find the hierarchy of k clusters.

3.1 HCEMST CLUSTERING ALGORITHM:

Given a point set S in \mathbf{E}^n and the desired number of clusters k , the hierarchical method starts by constructing an **MST** from the points in S . The weight of the edge in the tree is Euclidean distance between the two end points. Next the average weight \bar{W} of the edges in the

entire **EMST** and its standard deviation σ are computed; any edge with $W > \hat{W} + \sigma$ or *current longest edge* is removed from the tree. This leads to a set of disjoint subtrees $S_T = \{T_1, T_2 \dots\}$ (*divisive approach*). Each of these subtrees T_i is treated as cluster. Oleksandr Grygorash et al proposed algorithm [14] which generates k clusters. We modified the algorithm in order to generate k clusters with centers and each of the k clusters (subtrees) is again an **EMST**, which is further used for clustering using hierarchical clustering method (*agglomerative approach*) for generating dendrogram. Hence we named the new algorithm as Hierarchical Center Euclidean Minimum Spanning Tree (**HCEMST**). Each center point of k clusters is a representative point for the each subtree T' . A point c_i is assigned to a cluster i if $c_i \in T_i$. The group of center points is represented as $S = \{c_1, c_2, \dots, c_k\}$. This algorithm use both divisive as well as agglomerative approach to find Dual similarity clusters. Since the subtrees are themselves are clusters, are further, classified in order to get more informative similarity clusters

Algorithm: HCEMST(k)

Input : S the point set

Output : k number of clusters with dendrogram and S (set of center points)

Let e be an edge in the **EMST** constructed from S

Let W_e be the weight of e

Let σ be the standard deviation of the edge weights

Let S_T be the set of disjoint subtrees of the **EMST**

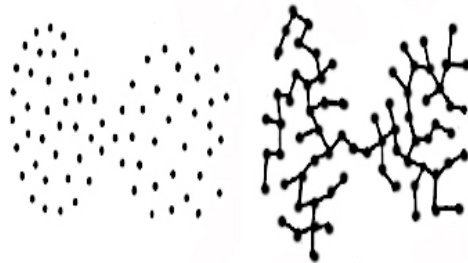
Let n_c be the number of clusters

1. Construct an **EMST** from S
2. Compute the average weight of \hat{W} of all the edges
3. Compute standard deviation σ of the edges
4. $S_T = \emptyset; n_c = 1$
5. **Repeat**
6. **For** each $e \in \mathbf{EMST}$
7. **If** ($W_e > \hat{W} + \sigma$) or (current longest edge e)
8. Remove e from **EMST** which result T' , a is new disjoint subtree
9. $S_T = S_T \cup \{T'\}$ // T' is new disjoint subtree
10. $n_c = n_c + 1$
11. Compute the center C_i of T_i using eccentricity of points
12. $S = \cup_{T_i \in S_T} \{C_i\}$
13. Begin with T' , disjoint clusters with level

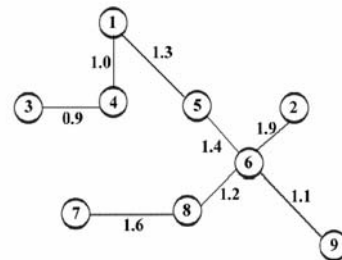
$L_{nc}(0) = 0$ and sequence number $m = 0$

14. **While** (T' has some edge)
15. $e = \text{get-min-edge}(T')$ // for least dissimilar pair of clusters
16. $(i, j) = \text{get-vertices}(e)$
17. Increment the sequence number $m = m + 1$, merge the clusters (i) and (j) , into single cluster to form next clustering m and set the level of this cluster to $L_{nc}(m) = e$;
18. Update T' by forming new vertex by combining the vertices i, j ;
19. **Until** $n_c = k$
20. **Return** k clusters with k dendrogram

Figure 1 illustrate a typical example of cases in which simply remove the $k-1$ longest edges does not necessarily output the desired cluster structure. Our algorithm finds the center of the each cluster, which will be useful in many



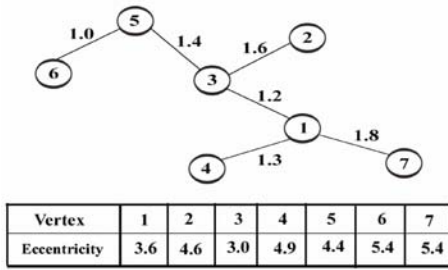
applications. Figure 2 shows the possible distribution of the points in the two cluster structures with their center vertex as 5 and 3.



Vertex	1	2	3	4	5	6	7	8	9
Eccentricity	5.5	6.7	7.4	6.5	4.2	4.6	7.4	5.8	5.7

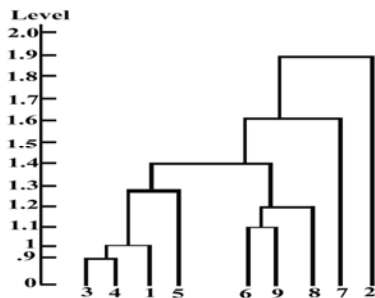
Figure 1. Clusters connected through a point

(a)

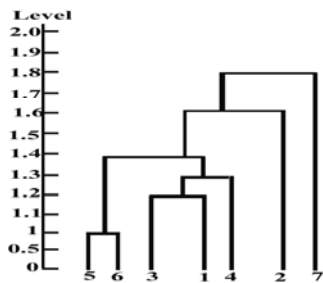


(b)

Figure 2. Two Clusters with Center vertices 5 and 3



(a)



(b)

Figure 3. Dendrogram for two Clusters

Euclidean minimum spanning tree is constructed in line 1. The average and standard deviation of the weighted edges of Euclidean minimum spanning tree are computed to find inconsistent edges are specified in lines 2-3. The inconsistent edges are identified and removed from Euclidean minimum spanning tree in order to generate subtree T' is specified in the lines 7-9. For the newly created subtree T' again further hierarchical clustering is performed is specified in the lines

13-18. It places the entire disjoint cluster at level 0 (line 13). It then checks to see if T' still contains some edge (line 14). If so, it finds minimum edge e (line 15). It then finds the vertices i, j of an edge e (line 16). It then merges the vertices and forms a new vertex (*agglomerative approach*). At the same time the sequence number is increased by one and the level of the new cluster is set to the edge weight (line 17). Finally, the Updation of Euclidean minimum spanning tree is performed at line 18. The lines 15-18 in the algorithm are repeated until the minimum spanning tree T' has no edge. The algorithm takes $O(k|E|^2)$ time, where k is number cluster and E be the number edges in the subtree T' . The outcome of the algorithm is shown in the figure 3. Our **HCEMST** algorithm works in two phases. The first phase produces subtree T' . The second phase converts T' into dendrogram. Our algorithm uses both divisive as well as agglomerative approach to find Dual similarity clusters.

4. CONCLUSION

Our **HCEMST** clustering algorithm assumes a given cluster number. The algorithm gradually finds k clusters with center for each cluster. The **HCEMST** clustering algorithm also generates dendrogram which is used to find the relationship between objects with in a cluster. The intra- cluster distances in the k clusters are shown in the dendrogram (figure3). This will be very useful in many applications. Our algorithm does not require the users to select and try various parameters combinations in order to get the desired output. All of these look nice from theoretical point of view. However from practical point of view, there is still some room for improvement for running time of the clustering algorithm. This could perhaps be accomplished by using some appropriate data structure. In the future we will explore and test our proposed clustering algorithm in various domains. The **HCEMST** algorithm uses both divisive as well agglomerative approaches. In this paper we used both the approaches to find Dual similarity clusters. We will further study the rich properties of EMST-based clustering methods in solving different clustering problems.



REFERENCES

- [1] Anil K. Jain, Richard C. Dubes 'Algorithm for Clustering Data', *Michigan State University, Prentice Hall, Englewood Cliffs, New Jersey* 07632.1988.
- [2] T.Asano, B.Bhattacharya, M.Keil and F.Yao.'Clustering Algorithms based on minimum and maximum spanning trees'. In *Proceedings of the 4th Annual Symposium on Computational Geometry*,Pages 252-257, 1988.
- [3] D.Avis 'Diameter partitioning'. *Discrete and Computational Geometry*, 1:265-276, 1986
- [4] M.Fredman and D.Willard. 'Trans-dichotomous algorithms for minimum spanning trees and shortest paths'. In *Proceedings of the 31st Annual IEEE Symposium on Foundations of Computer Science*,pages 719-725, 1990.
- [5] P. Hansen and M. Delattre 'Complete-link cluster analysis' by graph coloring *Journal of the American Statistical Association* 73, 397-403, 1978.
- [6] Hubert L. J 'Min and max hierarchical clustering using asymmetric similarity measures' *Psychometrika* 38, 63-72, 1973.
- [7] H.Gabow, T.Spencer and R.Rarjan. 'Efficient algorithms for finding minimum spanning trees in undirected and directed graphs'.*Combinatorica*, 6(2):109-122, 1986.
- [8] J.C. Gower and G.J.S. Ross 'Minimum Spanning trees and single-linkage cluster analysis', *Applied Statistics* 18, 54-64, 1969.
- [9] S. C. Johnson 'Hierarchical clustering schemes' *Psychometrika* 32, 241-254, 1967.
- [10] D.Johnson. 'The np-completeness column: An ongoing guide'. *Journal of Algorithms*,3:182-195, 1982.
- [11] D.Karger, P.Klein and R.Tarjan. 'A randomized linear-time algorithm to find minimum spanning trees'. *Journal of the ACM*, 42(2):321-328, 1995.
- [12] J.Kruskal. 'On the shortest spanning subtree and the travelling salesman problem'. In *Proceedings of the American Mathematical Society*, Pages 48-50, 1956.
- [13] J.Nesetril, E.Milkova and H.Nesetrilova. Otakar boruvka 'on minimum spanning tree problem': Translation of both the 1926 papers, comments, history. DMATH: *Discrete Mathematics*, 233, 2001.
- [14] Oleksandr Grygorash, Yan Zhou, Zach Jorgensen. 'Minimum spanning Tree Based Clustering Algorithms'. *Proceedings of the 18th IEEE International conference on tools with Artificial Intelligence (ICTAI'06)* 2006.
- [15] N.Paivinen. 'Clustering with a minimum spanning of scale-free-like structure', *Pattern Recogn. Lett.*,26(7): 921-930, 2005.
- [16] F.Preparata and M.Shamos. 'Computational Geometry: An Introduction'. *Springer-Verlag, Newyr, NY,USA*, 1985
- [17] R.Prim. 'Shortest connection networks and some generalization'. *Bell systems Technical Journal*,36:1389-1401, 1957.
- [18] Stefan Wuchty and Peter F. Stadler. 'Centers of Complex Networks'. 2006
- [19] Y.Xu, V.Olman and D.Xu. 'Minimum spanning trees for gene expression data clustering'. *Genome Informatics*,12:24-33, 2001.
- [20] C.Zahn. 'Graph-theoretical methods for detecting and describing gestalt clusters'. *IEEE Transactions on Computers*, C-20:68-86, 1971.

BIOGRAPHY OF AUTHORS



S. John Peter is working as Assistant professor in Computer Science, St.Xavier's college (Autonomous), Palayamkottai, Tirunelveli. He earned his M.Sc degree from Bharadhasan University, Trichirappali. He also earned his M.Phil from Bhradhasan University, Trichirappali. Now he is doing Ph.D in Computer Science at Manonmaniam Sundranar University, Tirunelveli. He has presented research papers on clustering algorithm in various national seminars.



Dr. S. P. Victor earned his M.C.A. degree from Bharathidasan University, Tiruchirappalli. The M. S. University, Tirunelveli, awarded him Ph.D. degree in Computer Science for his research in Parallel Algorithms. He is the Head of the department of computer science, and the Director of the computer science research centre, St. Xavier's college (Autonomous), Palayamkottai, Tirunelveli. The M.S. University, Tirunelveli and Bharathiar University, Coimbatore have recognized him as a research guide. He has published research papers in international, national journals and conference proceedings. He has organized Conferences and Seminars at national and state level.