









The framework of LSGCE methodology is illustrated in the above Fig.2. It includes three major processes which are as follows,

- 1) Generation of base clustering to form Cluster Ensemble ( $\pi$ ).
- 2) Producing Distilled Similarity Matrix (DSM) using Weighted Spectral Quality algorithm.
- 3) Extracting the ultimate data partition ( $\pi^*$ ) by exploiting the Spectral Clustering based Consensus Function.

### 4.1 Creating Cluster Ensemble

Consider the Dataset  $X = \{x_1, \dots, x_n\}$  be a set of data points and  $\pi$  denotes the cluster ensemble such that  $\pi = \{\pi_1, \dots, \pi_M\}$  are the ensemble members with base clustering. Each base clustering profits a set of clusters  $\pi_i = \{C_1^i, C_2^i, \dots, C_k^i\}$  where as  $k_i$  is number of clusters in the clustering results. The following Fig.3 illustrates the Sample Cluster Ensemble and its corresponding clusters

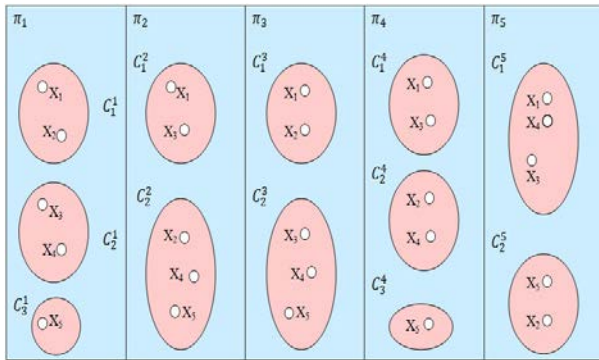


Fig.3 Sample Cluster Ensemble

Here, in this LSGCE approach Homogeneous Cluster ensemble generation method is used in which running Spectral Clustering algorithm n times obtain the base clustering results. In particular to a full space ensemble, the generations of the base clustering are extracted from the original dataset with all its instances and attributes. The systematic process of the Spectral Clustering algorithm are summarized as follows,

#### SPECTRAL CLUSTERING ALGORITHM

**Input:** A set of points  $X = \{x_1, \dots, x_n\}$  in  $R^k$ .

**Output:** The Resulting set of Base clusters.

- 1) begin
- 2) Compute Affinity Matrix  $A \in R^{n \times n}$  defined by

$$A_{ij} = \frac{\exp(-||x_i - x_j||^2)}{2\sigma^2} \text{ If } i \neq j, A_{ij} = 0$$

- 3) Define  $D$  be the Diagonal Matrix, and Build the Laplacian Matrix  $L = D^{-1/2} * AD^{-1/2}$

- 4) Determine  $e_1, e_2, \dots, e_k$  such that  $k$  be largest eigenvectors of Matrix  $L$ .
- 5) Build the Matrix  $E = [e_1, e_2, \dots, e_k] \in R^{n \times k}$
- 6) Build the Matrix  $B$  from  $E$  by stabilize each row to have a unit length  $B_{ij} = \frac{E_{ij}}{(\sum_j E^2_{ij})^{1/2}}$
- 7) Apply K-Means clustering technique over each row of  $B$  as a point in  $R^k$  to cluster them into  $k$  clusters.
- 8) end

Having obtained the set of base clusters formed from the repeated runs of the Spectral Clustering algorithm, the Cluster Ensemble is employed. From the Sample Cluster Ensemble shown in the Fig.3, Label assignment Matrix of size was created as illustrated in the Fig.4. It specifically symbolizes the cluster labels that are assigned to each data points by different base clustering.

	$\Pi_1$	$\Pi_2$	$\Pi_3$	$\Pi_4$	$\Pi_5$
$x_1$	$C_1^1$	$C_1^2$	$C_1^3$	$C_1^4$	$C_1^5$
$x_2$	$C_1^1$	$C_2^2$	$C_1^3$	$C_2^4$	$C_2^5$
$x_3$	$C_2^1$	$C_1^2$	$C_2^3$	$C_1^4$	$C_1^5$
$x_4$	$C_2^1$	$C_2^2$	$C_2^3$	$C_2^4$	$C_1^5$
$x_5$	$C_3^1$	$C_2^2$	$C_2^3$	$C_3^4$	$C_2^5$

Fig.4 Label Assignment Matrix

Moreover the Binary Cluster Association Matrix [20] illustrated in Fig.5 exposes the cluster specific nature of the original label assignment matrix. Each entry in this matrix mainly denotes the crispy association degree between the data points and the

	$C_1^1$	$C_2^1$	$C_3^1$	$C_1^2$	$C_2^2$	$C_1^3$	$C_2^3$	$C_1^4$	$C_2^4$	$C_3^4$	$C_1^5$	$C_2^5$
$x_1$	1	0	0	1	0	1	0	1	0	0	1	0
$x_2$	1	0	0	0	1	1	0	0	1	0	0	1
$x_3$	0	1	0	1	0	0	1	1	0	0	1	0
$x_4$	0	1	0	0	1	0	1	0	1	0	1	0
$x_5$	0	0	1	0	1	0	1	0	0	1	0	1

Fig.5 Binary Cluster Association Matrix

clusters formed in the ensemble. The co-association degree is based on the occurrence of the data points in the extracted clusters. It fills the matrix entry by either "1" or "0" such that if the particular data point



$$Sim_{WSQL}(\pi_i, \pi_j) = \frac{WSQL_{ij}^c}{Min[(\sum W_t(\pi_i)), (\sum W_t(\pi_j))]} * DC \quad (4)$$

where  $W_t(\pi_i)$  and  $W_t(\pi_j)$  denotes the summation of total weights associated with the clusters that forms the triple in the Linked Spectral Graph. Formally, a triple  $t = (V_t, E_t)$  is a sub graph of  $LSG$  containing two vertices  $V_t = \{v_i, v_j\} \subset V$  and three cluster nodes termed as edges  $E_t = \{e_i, e_j, e_k\} \subset E$ . Hence the Similarity measure between the two ensemble members  $\pi_i$  and  $\pi_j$  can be valued by considering minimum sum of weighted triples in the two partitions. Additionally to boost the confidence level of recognizing two non identical ensemble members being similar, a constant Decay factor  $DC \in [0,1]$  is fixed. Following the Sample Cluster Ensemble shown in the Fig.3 the similarity measures between each ensemble members in the Cluster Ensemble are estimated with the decay factor fixed to 0.9. These measures are then formulated in Fig.7 and the Distilled Similarity Matrix is illustrated in the Fig.8. Consequently, from the empirical analysis it is recognized that the estimation of similarity measures between the ensemble members rather than the clusters in the Ensemble drastically improves the similarity degrees of the clustering results over the conventional Cluster Ensemble techniques.

	$\Pi_1$	$\Pi_2$	$\Pi_3$	$\Pi_4$	$\Pi_5$
$\Pi_1$		0.41	0.36	0.41	0.17
$\Pi_2$			0.50	0.29	0.43
$\Pi_3$				1.0	0.49
$\Pi_4$					0.36
$\Pi_5$					

Fig.7 Similarity Measures between the Ensemble Members where DC = 0.9

	$c_1^1$	$c_2^1$	$c_3^1$	$c_1^2$	$c_2^2$	$c_1^3$	$c_2^3$	$c_1^4$	$c_2^4$	$c_3^4$	$c_1^5$	$c_2^5$
$X_1$	1	0.41	0.36	1	0.43	1	0.50	1	0.29	1	1	0.43
$X_2$	1	0.41	0.36	0.50	1	1	0.50	0.36	1	0.36	0.17	1
$X_3$	0.41	1	0.41	1	0.43	0.50	1	1	0.29	0.36	1	0.43
$X_4$	0.41	1	0.41	0.43	1	0.50	1	1	1	0.36	1	0.43
$X_5$	0.41	0.41	1	0.43	1	0.50	1	0.36	0.29	1	0.17	1

Fig.8 Distilled Similarity Matrix (DSM)

The Weighted Spectral Quality (WSQL) algorithm is summarized below,

**ALGORITHM:**  $WSQL(LSG, \pi_i, \pi_j)$

**Input:** A Dataset with  $x$  dimensional data objects.

**Output:** Distilled Similarity Matrix.

1)  $LSG = (V, W)$  a linked spectral graph where  $C_i, C_j \in V$  ; ;

2) **begin**

3) Compute Weight:  $W_{ij} = \frac{d_i \cap d_j}{d_i \cup d_j}$  ;

4) **init**  $WSQL_{ij}^c \rightarrow 0$  ;

5) **for each**  $C \in \pi_i$

6) **If**  $C \in \pi_j$

7)  $NP_{cut}(\pi_i, \pi_j) = \frac{MinCut(\pi_i, \pi_j)}{\sum(vol(\pi_i, \pi_j))}$  ;

8)  $WSQL_{ij}^c = \frac{1}{n} \sum_{i=1}^p NP_{cut}$  ;

9) **Return**  $WSQL_{ij}^c$  ;

10) **end**

11)  $Sim_{WSQL}(\pi_i, \pi_j) = \frac{WSQL_{ij}^c}{Min[(\sum W_t(\pi_i)), (\sum W_t(\pi_j))]} * DC$  ;

12) **end**

**4.3 Applying Spectral Clustering based Consensus Function to DSM**

Having attained the DSM, a Spectral Clustering based Consensus Function [60] is applied to extract the final clustering results. This consensus technique requires the Distilled Similarity Measures through which it applies the spectral clustering algorithm to partition the similarity measures for exploiting the ultimate clustering solutions. In the first step, it builds the affinity matrix with the entries in the obtained DSM in which it represents the similarity degrees between the two ensemble members  $\pi_i$  and  $\pi_j$ . In the second step, after obtaining the affinity matrix, spectral clustering generates the diagonal matrix through the summation of entries in the diagonal. In the third step, it normalizes the affinity matrix in order to perform efficient dimensionality reduction. In the fourth step, it produces the Eigen vectors corresponding to the first six largest Eigen values of the affinity matrix and re-normalizes each rows of the matrix. Finally in the fifth step, K-Means

algorithm is implemented to assign the samples in the newly formed data matrix to their corresponding clusters. Thus this consensus function proves to be the powerful and efficient method in obtaining absolute cluster results and also it attains the nearer optimal solutions.

## 5 Performance Evaluation

This section exposes the performance of proposed Linked Spectral Graph based Cluster Ensemble approach using few validity indices and variety of Medical datasets. The quality of each ensemble members in the total Cluster Ensemble acquired by this technique is evaluated against two different traditional clustering algorithms.

### 5.1 Examined Datasets

The experimental analysis is conducted over five medical datasets which are taken from the UCI Machine Learning Repository [61]. The details regarding the number of instances and attributes are summarized in the below Table I.

TABLE I  
DATASETS USED IN THE EXPERIMENT

Datasets	Instances	Attributes
Arrhythmia	452	279
Dermatology	366	33
Heart Disease	303	75
Hepatitis	155	19
Lung Cancer	32	56

The descriptions about the experimented Medical datasets are as follows,

1) *Arrhythmia*- This dataset mainly denotes the presence and absence of cardiac disease. Among the 279 attributes, 206 are linear valued and the remaining are nominal.

2) *Dermatology*- The data comprises of clinical features of Erythematic disease observed in the patient. It also includes the age feature and the possibility of high effect of intermediate results of the disease.

3) *Heart Disease*- It includes the details of heart disease present in the patient which was taken from the Cleveland database.

4) *Hepatitis*- It represents the data regarding the inflammation of the liver and the inflammatory cells present in tissues of the organ.

5) *Lung Cancer*- This dataset describes about the three types of pathological lung cancer cells in the patient.

### 5.2 Evaluation Criteria

The experiment set out to observe the performance of the LSGCE in contrast to few conventional clustering algorithms. In order to analyze the efficiency of the proposed work, the final clustering results of each method is evaluated with its appropriate true labels by using the following performance validity metrics.

1) *Classification Accuracy* – It is the measure [62] of number of exactly classified data objects of the clustering results compared with the known true labels divided by the total number of data points in the datasets. This Classification accuracy measures can be estimated as given below,

$$CA(\pi^*) = \frac{\sum_{i=0}^k M_i}{D} \quad (5)$$

where  $\pi^*$  denotes the final partition results,  $M_i$  illustrates the number of data objects with the majority of the cluster label points in the cluster  $i$ , then  $D$  is the total number of data objects in the dataset.

2) *Error Rate*- It is the term that describes the degree of errors or irrelevant data encountered during data clustering. This error rate is computed as given below,

$$E = 1 - CA \quad (6)$$

where  $CA$  denotes the clustering accuracy calculated from the equation (5).

3) *Rand Index*- Generally Rand index  $I$ , is the measure of the similarity between the two data clustering. In other words it is stated that a measure [63] of number of object pairs that exist in the same and different clusters. More formally it can also be stated as a proportional measure of the quantity of agreements and disagreements between the two partitions. It can be calculated as below,

$$R = \frac{(x+y)}{(x+y)+(u+v)} \quad (7)$$

where  $(x + y)$  can be denoted as the number of agreements between the two clusters  $C_i$  and  $C_j$  similarly  $(u + v)$  can be considered as the number of disagreements between the same two clusters.



4) *Compactness*- It measures [62] the average distances between the each pair of data points occurring in the same cluster. More specifically it is calculated as given below,

$$CP(\pi^*) = \frac{1}{D} \sum_{k=1}^K d_k \left( \frac{\sum_{x_i, x_j \in C_k} d(x_i, x_j)}{d_k(d_k-1)} \right) \quad (8)$$

where  $K$  denotes then number of clusters formed finally,  $d_k$  is the number of data objects corresponding to that particular cluster, and  $d(x_i, x_j)$  is the distance between the data points  $x_i$  and  $x_j$ , then  $D$  be the total number of data points in the dataset.

5) *Dunn*- Its main aspire [64] is to identify the closeness and the well separated clusters. It compares the size of the clusters with the distance between the clusters. Such that it is stated the distances between the clusters are expected to be large and the diameter of the clusters should be small. Hence it can be computed as given below,

$$Dunn(\pi^*) = \frac{\min d(C_i, C_j)}{\max \Delta(C_i)} \quad (9)$$

where  $d(C_i, C_j)$  denotes the distance computed between the two clusters  $C_i$  and  $C_j$ , and  $\Delta(C_i)$  expresses the size of the cluster .

6) *Davies Bouldin (DB)*- This DB [65] measure mainly determines average of the similarity between the two clusters  $C_i$  and  $C_j$  in which it is defined by the estimation of dispersion of a single cluster and the dissimilarity measure between the two clusters. It is evaluated as follows,

$$Sim_{ij} = \frac{s_i + s_j}{d_{ij}} \quad (10)$$

in which  $s_i$  denotes the dispersion of  $C_i$  and  $d_{ij}$  shows the dissimilarity between the two clusters, these can be calculated as given below,

$$s_i = \frac{1}{|C_i|} \sum_{v \in C_i} d(x, v_i) \quad (11)$$

$$d_{ij} = d(v_i, v_j) \quad (12)$$

where  $|C_i|$  denotes the number of data points in the cluster  $C_i$  and  $v_i$  and  $v_j$  shows the center points of the two clusters  $C_i$  and  $C_j$  respectively. Hence from the above the DB can be derived as,

$$DB(\pi^*) = \frac{1}{k} \sum_{i=1}^k Sim_i \quad (13)$$

where  $Sim_i = \max( Sim_{ij} )$  such that  $i \neq j$ .

The following Table II and Table III exemplify the results of each measure when evaluated with the LSGCE algorithm implemented in MATLAB environment to find its efficacy.

TABLE II  
AVERAGE CLUSTERING ACCURACY AND ERROR RATES OF 10 RUNS

Datasets	Clustering Accuracy	Clustering Error Rate
Heart Disease	0.785	0.215
Lung Cancer	0.468	0.532
Hepatitis	0.750	0.250
Arrhythmia	0.565	0.435
Dermatology	0.588	0.412

TABLE III  
PERFORMANCE COMPARISON AMONG EVALUATION INDICES

Datasets	Compactness	Rand Index	Davies-Bouldin	Dunn
Heart Disease	36.02	0.174	0.840	0.785
Lung Cancer	50.03	0.130	2.318	0.993
Hepatitis	66.04	0.595	0.640	1.470
Arrhythmia	62.83	0.420	2.805	0.824
Dermatology	15.20	0.212	1.818	1.574

The following Figures represent the graphical notation of the performance of LSGCE when examined with the Medical Datasets over several evaluation indices.

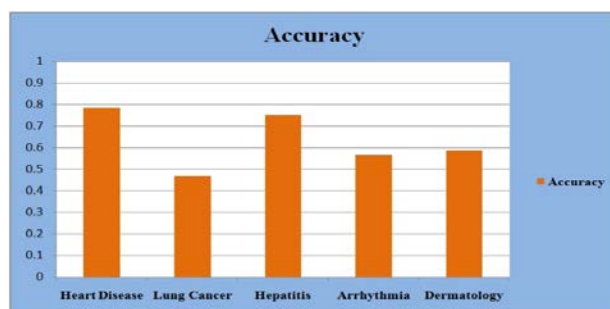


Fig.9. Average Accuracy Rates of 10 runs

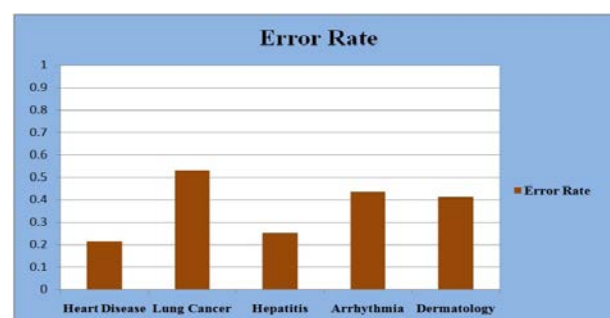


Fig.10. Average Error Rates of 10 runs

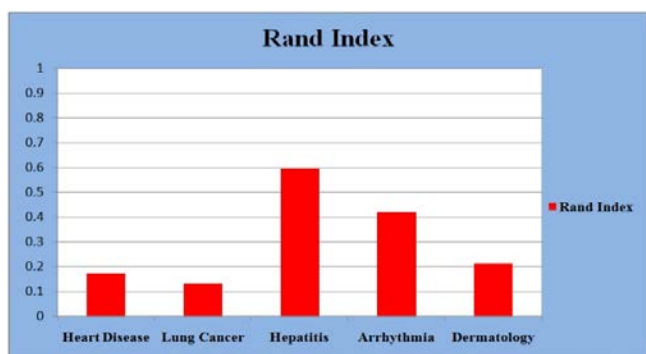


Fig.11 Rand Index Rates

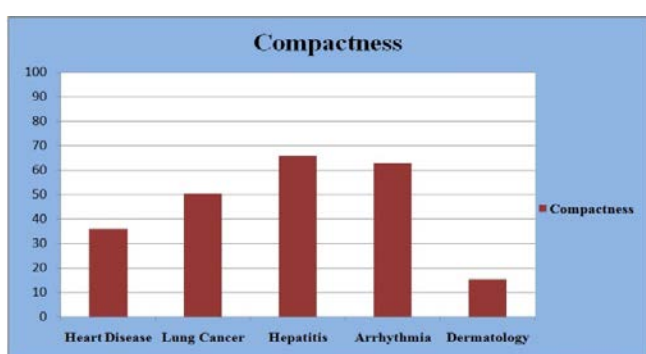


Fig.12 Compactness Rates

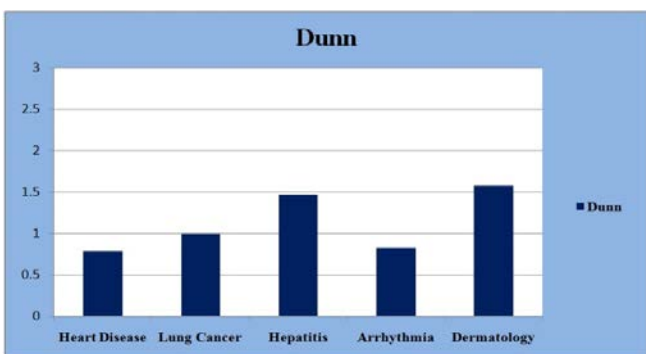


Fig.13 Dunn Rates

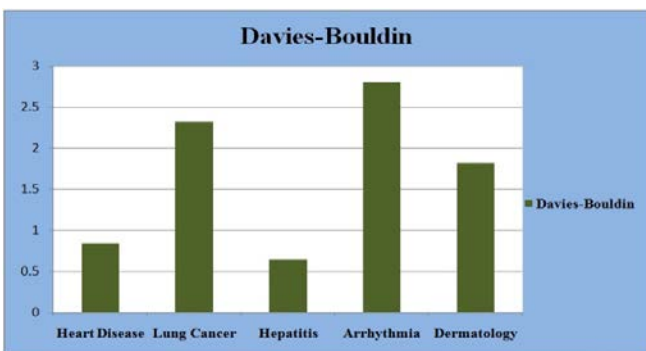


Fig.14 Davies-Bouldin Rates

### 5.3 Compared Traditional Clustering Methods

In order to estimate the efficiency of the newly proposed Linked Spectral Graph based Cluster Ensemble (LSGCE) approach, the following two traditional clustering methods have been contrasted over the Lung Cancer dataset.

#### 5.3.1 K-Means Clustering Method

K-Means is the well-known conventional clustering algorithm [66] often used to cluster the numerical data. It is an algorithm mainly framed to find the K-center point of the dataset based on the distance between the other data points and the center point. The Euclidean distance function is most preferable in nature. It initiates the K number of value as a starting point to estimate the cluster centers. In this algorithm, the main issue is to reduce the distance between the data object and the corresponding cluster center point in the dataset. Moreover this K-Means algorithm [66] faces two main challenges such as its behavior mainly depends on the initial center point and it often converges to local minima. Different initial cluster center points provide different clustering solutions. This difficulty is most widely seen when initial center points are not well separated.

#### 5.3.2 Fuzzy C-Means Clustering Method

Fuzzy C-Means clustering algorithm [67] was mainly established to smoothen the hard nature of K-Means algorithm in which a data alone can assign to the cluster. It mainly makes use of the fuzzy partitioning to let the data objects to assign to all the clusters generated with the membership grade between 0 and 1 and the sum of it is 1. By considering the highest grade the data is recorded to its appropriate cluster. Adversely the same challenging factors of K-Means happened to Fuzzy C-Means algorithm as it cannot ensure the global optima. The Clustering result is highly dependent to the randomly assigned initial membership grades.

The following Table IV compares the average clustering accuracy rates of LSGCE with traditional clustering techniques over 10 runs.

Furthermore the forthcoming Fig.15 represent the graphical illustration of the performance of newly proposed Linked Spectral Graph based Cluster Ensemble (LSGCE) approach examined with Medical datasets.

TABLE IV  
ACCURACY COMPARISON OF TRADITIONAL CLUSTERING METHODS

Datasets	LSGCE	KM	FCM
Heart Disease	0.785	0.761	0.691
Lung Cancer	0.468	0.431	0.432
Hepatitis	0.750	0.713	0.732
Arrhythmia	0.565	0.554	0.488
Dermatology	0.588	0.428	0.419

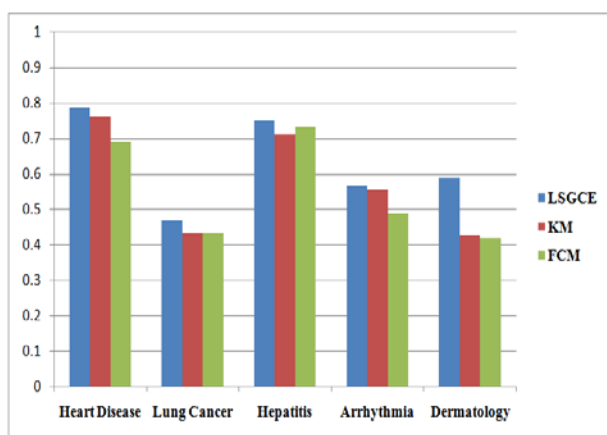


Fig.15 Accuracy Mean of Different Medical Datasets

## 6 Conclusion and Future Work

The main contribution in this paper is to exemplify the novel Linked Spectral Graph based Cluster Ensemble approach for providing efficiency in clustering Medical data and also in reducing cluster degradation problem. It greatly aims to explore and makes use of the relationship degree between the generated base clustering solutions. Additionally LSGCE performs the similarity assessment among the ensemble members of the Cluster Ensemble. This allows formation of Distilled Similarity Matrix (DSM) to be refined from the traditional Binary cluster association Matrix. The challenging issue of generating DSM is expertly resolved by Weighted Spectral Quality (WSQL) algorithm. With the

results of the extracted similarity measures, Spectral based Consensus Function is applied to finalize the ultimate cluster solutions. Hence the experimental investigation of conventional clustering algorithms tested over the Medical datasets suggests that newly proposed LSGCE approach highly overwhelms the traditional ones. Beyond these accomplishments, the future work includes the extension of LSGCE in Text data clustering. Furthermore this new approach can also be applied to other business related dataset with huge dimensions and also it further improves its efficiency in execution time.

### References:

- [1] D.S. Hochbaum and D.B. Shmoys, "A Best Possible Heuristic for the K-Center Problem," *Math. of Operational Research*, vol. 10, no. 2, pp. 180-184, 1985.
- [2] L. Kaufman and P.J. Rousseeuw, "Finding Groups in Data: An Introduction to Cluster Analysis". *Wiley Publishers*, 1990.
- [3] A.K. Jain and R.C. Dubes, *Algorithms for Clustering*. Prentice-Hall, 1998.
- [4] P. Zhang, X. Wang, and P.X. Song, "Clustering Categorical Data Based on Distance Vectors," *The J. Am. Statistical Assoc.*, vol. 101, no. 473, pp. 355-367, 2006.
- [5] A. Ng, M. Jordan, and Y. Weiss, "On Spectral Clustering: Analysis and an Algorithm," *Advances in Neural Information Processing Systems*, vol. 14, pp. 849-856, 2001.
- [6] Z. Huang, "Extensions to the K-Means Algorithm for Clustering Large Data Sets with Categorical Values," *Data Mining and Knowledge Discovery*, vol. 2, pp. 283-304, 1998.
- [7] D. Cristofor and D. Simovici, "Finding Median Partitions Using Information-Theoretical-Based Genetic Algorithms," *J. Universal Computer Science*, vol. 8, no. 2, pp. 153-172, 2002.
- [8] D.H. Fisher, "Knowledge Acquisition via Incremental Conceptual Clustering," *Machine Learning*, vol. 2, pp. 139-172, 1987.
- [9] D. Gibson, J. Kleinberg, and P. Raghavan, "Clustering Categorical Data An Approach Based on Dynamical Systems," *VLDB J.*, vol. 8, nos. 3-4, pp. 222-236, 2000.
- [10] S. Guha, R. Rastogi, and K. Shim, "ROCK: A Robust Clustering Algorithm for Categorical Attributes," *Information Systems*, vol. 25, no. 5, pp. 345-366, 2000.
- [11] M.J. Zaki and M. Peters, "Clicks: Mining Subspace Clusters in Categorical Data via Kpartite Maximal Cliques," *Proc. Int'l Conf. Data Eng. (ICDE)*, pp. 355-356, 2005.

- [12] V. Ganti, J. Gehrke, and R. Ramakrishnan, "CACTUS: Clustering Categorical Data Using Summaries," *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD)*, pp. 73-83, 1999.
- [13] D. Barbara, Y. Li, and J. Couto, "COOLCAT: An Entropy-Based Algorithm for Categorical Clustering," *Proc. Int'l Conf. Information and Knowledge Management (CIKM)*, pp. 582-589, 2002.
- [14] Y. Yang, S. Guan, and J. You, "CLOPE: A Fast and Effective Clustering Algorithm for Transactional Data," *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD)*, pp. 682-687, 2002.
- [15] Z. He, X. Xu, and S. Deng, "Squeezer: An Efficient Algorithm for Clustering Categorical Data," *J. Computer Science and Technology*, vol. 17, no. 5, pp. 611-624, 2002.
- [16] P. Andritsos and V. Tzerpos, "Information-Theoretic Software Clustering," *IEEE Trans. Software Eng.*, vol. 31, no. 2, pp. 150-165, Feb. 2005.
- [17] S. Indrajit, M. Ujjwal, & Nilanjan. "Differential Fuzzy Clustering for Categorical Data". *International Conference on Methods and Models in Computer Science*, 2009.
- [18] Sandro Vega-pons & Jose reuiz Shulcloper. "A Survey of Clustering Ensemble algorithms". *International Journal of Pattern Recognition and Artificial Intelligence*, Vol. 25, No. 3 337\_372 2011.
- [19] Domeniconi C and Al-Razgan M, "Weighted cluster ensembles: methods and analysis." *ACM Transaction on. Knowledge Discovery Data* 2(4) 1\_40. 2009.
- [20] Natthakan Iam-On, B. Tossapon, G.Simon, and Chris Price. "A Link based cluster ensemble approach for categorical data clustering." *IEEE Transactions on knowledge and data engineering*, Vol. 24, No. 3, 2012.
- [21] J. Kittler, M. Hatef, R. Duin, and J. Matas, "On Combining Classifiers," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226-239, Mar. 1998.
- [22] A. Strehl and J. Ghosh, "Cluster Ensembles: A Knowledge Reuse Framework for Combining Multiple Partitions," *J. Machine Learning Research*, vol. 3, pp. 583-617, 2002.
- [23] X. Hu and I. Yoo, "Cluster Ensemble and Its Applications in Gene Expression Analysis," *Proc. Asia-Pacific Bioinformatics Conf.*, pp. 297-302, 2004.
- [24] M. Law, A. Topchy, and A.K. Jain, "Multiobjective Data Clustering," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 2, pp. 424-430, 2004.
- [25] Huang. Z, "Extensions to the K-Means Algorithm for Clustering Large Data Sets with Categorical Values," *Data Mining and Knowledge Discovery*, vol. 2, pp. 283-304, 1998.
- [26] Fred A.L.N and Jain A.K, "Combining Multiple Clusterings Using Evidence Accumulation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 6, pp. 835-850, June 2005.
- [27] L.I. Kuncheva and D. Vetrov, "Evaluation of Stability of K-Means Cluster Ensembles with Respect to Random Initialization," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 11, pp. 1798-1808, Nov. 2006.
- [28] X.Z. Fern and C.E. Brodley, "Random Projection for High Dimensional Data Clustering: A Cluster Ensemble Approach," *Proc. Int'l Conf. Machine Learning (ICML)*, pp. 186-193, 2003.
- [29] S. Dudoit and J. Fridyand, "Bagging to Improve the Accuracy of a Clustering Procedure," *Bioinformatics*, vol. 19, no. 9, pp. 1090-1099, 2003.
- [30] B. Minaei-Bidgoli, A. Topchy, and W. Punch, "A Comparison of Resampling Methods for Clustering Ensembles," *Proc. Int'l Conf. Artificial Intelligence*, pp. 939-945, 2004.
- [31] A.P. Topchy, A.K. Jain, and W.F. Punch, "Clustering Ensembles: Models of Consensus and Weak Partitions," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1866-1881, Dec. 2005.
- [32] B. Fischer and J.M. Buhmann, "Bagging for Path-Based Clustering," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 11, pp. 1411-1415, Nov. 2003.
- [33] A. Strehl and J. Ghosh, "Cluster Ensembles: A Knowledge Reuse Framework for Combining Multiple Partitions," *J. Machine Learning Research*, vol. 3, pp. 583-617, 2002.
- [34] N. Iam-On, T. Boongoen, and S. Garrett, "Refining Pairwise Similarity Matrix for Cluster Ensemble Problem with Cluster Relations," *Proc. Int'l Conf. Discovery Science*, pp. 222-233, 2008.
- [35] L. Getoor and C.P. Diehl, "Link Mining: A Survey," *ACM SIGKDD Explorations Newsletter*, vol. 7, no. 2, pp. 3-12, 2005.
- [36] D. Cristofor and D. Simovici, "Finding Median Partitions Using Information-Theoretical-Based Genetic Algorithms," *J. Universal*

- Computer Science*, vol. 8, no. 2, pp. 153-172, 2002.
- [37] N. Nguyen and R. Caruana, "Consensus Clusterings," *Proc. IEEE Int'l Conf. Data Mining (ICDM)*, pp. 607-612, 2007.
- [38] Z. Yu, H.-S. Wong, and H. Wang, "Graph-Based Consensus Clustering for Class Discovery from Gene Expression Data," *Bioinformatics*, vol. 23, no. 21, pp. 2888-2896, 2007.
- [39] G. Karypis and V. Kumar, "Multilevel K-Way Partitioning Scheme for Irregular Graphs," *J. Parallel Distributed Computing*, vol. 48, no. 1, pp. 96-129, 1998.
- [40] A. Ng, M. Jordan, and Y. Weiss, "On Spectral Clustering: Analysis and an Algorithm," *Advances in Neural Information Processing Systems*, vol. 14, pp. 849-856, 2001.
- [41] X.Z. Fern and C.E. Brodley, "Solving Cluster Ensemble Problems by Bipartite Graph Partitioning," *Proc. Int'l Conf. Machine Learning (ICML)*, pp. 36-43, 2004.
- [42] A. Gionis, H. Mannila, and P. Tsaparas, "Clustering Aggregation," *Proc. Int'l Conf. Data Eng. (ICDE)*, pp. 341-352, 2005.
- [43] S. Dudoit and J. Fridlyand, "A Prediction-Based Resampling Method to Estimate the Number of Clusters in a Data Set," *Genome Biology*, vol. 3, no. 7, pp. 0036.1-0036.21, 2002.
- [44] S. Dudoit and J. Fridlyand, "Bagging to Improve the Accuracy of a Clustering Procedure," *Bioinformatics*, vol. 19, no. 9, pp. 1090-1099, 2003.
- [45] M. Smolkin and D. Ghosh, "Cluster Stability Scores for Microarray Data in Cancer Studies," *BMC Bioinformatics*, vol. 4, article 36, 2003.
- [46] S. Monti, P. Tamayo, J. Mesirov, and T. Golub, "Consensus Clustering: A Resampling-Based Method for Class Discovery and Visualization of Gene Expression Microarray Data," *Machine Learning*, vol. 52, pp. 9-118, 2003.
- [47] Z. Yu and H.-S. Wong, "Knowledge Based Cluster Ensemble for Cancer Discovery from Biomolecular Data," *IEEE Trans. NanoBioscience*, vol. 10, no. 2, pp. 76-85, June 2011.
- [48] Zhiwen Yu Member, IEEE, Hantao Chen Jane You Member, IEEE, Guoqiang Han Le Li "Hybrid Fuzzy Cluster Ensemble Framework for Tumor Clustering from Bio-molecular Data" *IEEE Transactions on computational biology and bioinformatics* 2013.
- [49] Zhiwen Yu, Member, IEEE, Hau-San Wong, Member, IEEE, Jane You, Member, IEEE, Qinmin Yang, Member, IEEE, and Hongying Liao "Knowledge based Cluster Ensemble for Cancer Discovery From Biomolecular Data" *IEEE Transactions on Nanobioscience*, Vol 10 No. 2, June 2011.
- [50] Yu J. & Lin Z C. "Squared error adjacency matrix clustering". *Technical report on Dept. of Computer Science*, Beijing Jiaotong University 2008.
- [51] Hongjun Wang, Hanhuai Shan & Arindam Banerjee. "Bayesian Cluster Ensembles". *Wiley Periodicals, Inc* 2011.
- [52] Gullo F, Domeniconi C, Tagarelli A "Projective clustering ensembles". In: *Proceedings of the international conference on data mining (ICDM)*, pp 794-799.
- [53] Ka Ka Ng E, Wai-Chee Fu A, Chi-Wing Wong R "Projective clustering by histograms". *IEEE Trans Knowl Data Eng (TKDE)* 17(3):369-383 2005.
- [54] Yangzihao Wang, "Spectral Clustering: A Graph Partitioning Point of View", *ECS231 Course Report*, CSE University of California, Davis.
- [55] Inderjit S.Dhillon, Yuqiang Guan, and Brian Kulis, "Kernel K-Means, Spectral Clustering and Normalized cuts", *ACM 1-58113-888-1/04/0008 KDD* 04 August 2004.
- [56] M. Al-Razgan, C. Domeniconi, and D. Barbara, "Random Subspace Ensembles for Clustering Categorical Data," *Supervised and Unsupervised Ensemble Methods and Their Applications*, pp. 31-48, Springer, 2008.
- [57] Z. He, X. Xu, and S. Deng, "A Cluster Ensemble Method for Clustering Categorical Data," *Information Fusion*, vol. 6, no. 2, pp. 143-151, 2005.
- [58] Jeh G, Widom J "SimRank: A Measure of Structural-Context Similarity." In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 538-543. ACM, New York 2002.
- [59] Weiguo Zheng, Lei Zou, Yan Song Feng, Le Chen, Dongyan, "Efficient SimRank based Similarity Join over Large graphs" *Proceedings of the VLDB endowment*, Vol 6 No.7 2011.
- [60] Zhiwen yu, Le Li, Jane You, Hau-San Wong, and Guoqiang Han, "SC3: Triple Spectral Clustering Based Consensus Clustering Framework for Class Discovery from Cancer Gene Expression Profiles", *IEEE/ACM Transaction on Computational Biology and Bioinformatics*, Vol 9, No. 6, December 2012.

- [61] A. Asuncion and D.J. Newman, "UCI Machine Learning Repository," *School of Information and Computer Science*, Univ. of California, <http://www.ics.uci.edu/~mlern/MLRepository.html>.
- [62] Nguyen N, Caruana R, "Consensus Clusterings." In *Proceedings of IEEE International Conference on Data Mining*, pp. 607-612. IEEE Computer Society, Washington, DC 2007.
- [63] Rand WM, "Objective Criteria for the Evaluation of Clustering Methods." *Journal of the American Statistical Association*, 66, 846-850 1971.
- [64] Dunn JC, "Well Separated Clusters and Optimal Fuzzy Partitions." *Cybernetics and Systems*, 4(1), 95-104 1974.
- [65] Davies DL, Bouldin DW, "A Cluster Separation Measure." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1(2), 224-227 1979.
- [66] G. Hammerly, C. Elken, "Alternatives to the K-means algorithm that find better clusterings", in: *Proceedings of the 11th International Conference on Information and Knowledge Management*, 2002, pp. 600-607.
- [67] R.L. Canon, J.Dave and J.C. Bezdek, "Efficient implementation of the fuzzy cmeans clustering algorithms". *IEEE Trans Pattern Anal Machine, Intell* 8, 248-255.
- [68] Yuzhen Zhao, Xiyu Liu, and Wenping Wang, "ROCK Clustering algorithm based on the P System with active membranes", *WSEAS Transactions on Computers*, Vol 13 2014.