

# An Efficient Link Based Cluster Approach for Data Clustering

M.Pavithra<sup>1</sup>, Dr.R.M.S.Parvathi<sup>2</sup>

Assistant Professor, Department of CSE, Jansons Institute of Technology, Coimbatore, India<sup>1</sup>

Dean- PG Studies, Sri Ramakrishna Institute of Technology, Coimbatore, India<sup>2</sup>

**ABSTRACT:** Clustering is to categorize data into groups or clusters such that the data in the same cluster are more similar to each other than to those in different clusters. The problem of clustering categorical data is to find a new partition in dataset to overcome the problem of clustering categorical data via cluster ensembles, result is observed that these techniques unluckily generate a final data partition based on incomplete information. The underlying ensemble-information matrix presents only cluster-data point relations, with many entries being left unknown. This problem degrades the quality of the clustering result. To improve clustering quality a new link-based approach the conventional matrix by discovering unknown entries through similarity between clusters in an ensemble and an efficient link-based algorithm is proposed for the underlying similarity assessment. In this paper propose C-Rank link-based algorithm improve clustering quality and ranking clusters in weighted networks. C-Rank consists of three major phases: (1) identification of candidate clusters; (2) ranking the candidates by integrated cohesion; and (3) elimination of non-maximal clusters. The finally apply this clustering result in graph partitioning technique is applied to a weighted bipartite graph that is formulated from the refined matrix.

**KEYWORDS:** Clustering, Data mining, Categorical data, Cluster Ensemble, link-based similarity, refined matrix, and C-Rank link based cluster.

## I. INTRODUCTION

DATA clustering is one of the fundamental tools we have for understanding the structure of a data set. Clustering is a collection of objects which are “similar” between them and are “dissimilar” to the objects belonging to other clusters [14]. Clustering often is a first step in data analysis. Clustering is the process of discovering homogeneous groups or clusters according to a given similarity measure. Clustering maximizes intra-connectivity among patterns in the same cluster while minimizing inter-connectivity between patterns in different clusters [7]. Clustering is an important technique in discovering meaningful groups of data points. Clustering provides speed and reliability in grouping similar objects in very large datasets [3]. Most previous clustering algorithms focus on numerical data whose inherent geometric properties can be exploited naturally to define distance functions between data points. However, much of the data existed in the databases is categorical, where attribute values can't be naturally ordered as numerical values [4].

An example of categorical attribute is *shape* whose values include *circle*, *rectangle*, *ellipse*, etc. Consensus clustering can provide benefits beyond what a single clustering algorithm can achieve. Consensus clustering algorithms often: generate better clustering's; find a combined clustering unattainable by any single clustering algorithm. The consensus clustering algorithm can be applied to the ensemble of all clustering's produced by discrete features of the data set [10]. Cluster ensemble (CE) is the method to combine several runs of different clustering algorithms to get a common partition of the original dataset, aiming for consolidation of results from a portfolio of individual clustering results [3]. A cluster ensemble consists of different partitions. Such partitions can be obtained from multiple applications of any single algorithm with different initializations, or from the application of different algorithms to the same dataset [6].

Cluster ensembles offer a solution to challenges inherent to clustering arising from its ill-posed nature: they can provide more robust and stable solutions by leveraging the consensus across multiple clustering results, while averaging out

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 12, December 2016

emergent spurious structures that arise due to the various biases to which each participating algorithm is tuned [5]. Clustering ensembles have emerged as a powerful method for improving both the robustness as well as the stability of unsupervised classification solutions. A cluster ensemble can be defined as the process of combining multiple partitions of the dataset into a single partition, with the objective of enhancing the consensus across multiple clustering results [12]. A cluster ensemble technique is characterized by two components: the mechanism to generate diverse partitions, and the consensus function to combine the input partitions into a final clustering.

## II. RELATED WORK

We briefly describe relevant work in the literature on categorical clustering and cluster ensembles. Clustering of categorical data has recently attracted the attention of many researchers [5] [6]. The *k*-modes algorithm is an extension of *k*-means for categorical features. To update the modes during the clustering process, the authors used a new distance measure based on the number of mis-matches between two points. Squeezer is a categorical clustering algorithm that processes one point at the time. ROCK (Robust Clustering using links) is a hierarchical clustering algorithm for categorical data. It uses the Jaccard coefficient to compute the distance between points [7] [8]. The COOLCAT algorithm is a scalable clustering algorithm that discovers clusters with minimal entropy in categorical data. COOLCAT uses categorical, rather than numerical attributes, enabling the mining of real-world datasets offered by fields such as psychology and statistics.

The algorithm is based on the idea that a cluster containing similar points has entropy smaller than a cluster of dissimilar points. Thus, COOLCAT uses entropy to define the criterion for grouping similar objects. LIMBO is a hierarchical clustering algorithm that uses the Information Bottleneck (IB) framework to define a distance measure for categorical tuples [9] [10]. The concepts of evolutionary computing and genetic algorithm have also been adopted by a partitioning method for categorical data, i.e., GAClust. Cobweb is a model-based method primarily exploited for categorical data sets. Different graph models have also been investigated by the STIRR, ROCK, and CLICK techniques. In addition, several density-based algorithms have also been devised for such purpose, for instance, CACTUS, COOLCAT, and CLOPE. The Cluster-based Similarity Partitioning Algorithm (CSPA) induces a graph from a co-association matrix and clusters it using the METIS algorithm [12] [13]. Hyper graph partitioning algorithm (HGPA) represents each cluster by a hyper edge in a graph where the nodes correspond to a given set of objects.



Figure 1. Overview of Cluster Ensemble Algorithm Framework.

## III. PROPOSED WORK

### MODULE DESCRIPTION

#### 1. CLUSTER ENSEMBLES OF CATEGORICAL DATA

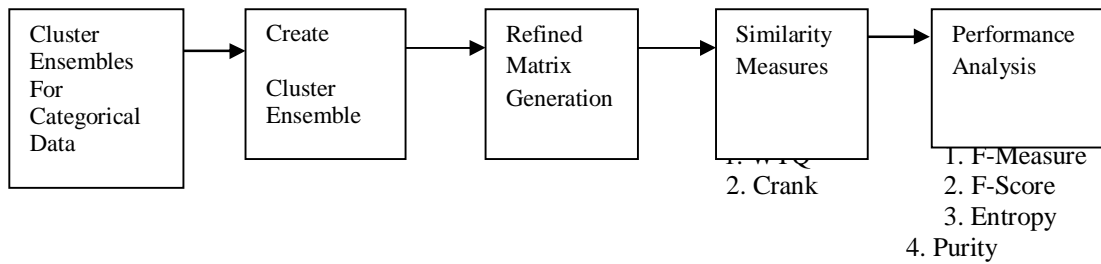
A cluster ensemble consists of different partitions. Such partitions can be obtained from multiple applications of any single algorithm with different initializations, or from the application of different algorithms to the same dataset. Cluster ensembles offer a solution to challenges inherent to clustering arising from its ill-posed nature: they can provide more robust and stable solutions by leveraging the consensus across multiple clustering results, while averaging out emergent spurious structures that arise due to the various biases to which each participating algorithm is tuned. Conventional approach, the technique developed in acquires a cluster ensemble.

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 12, December 2016

Figure 2. shows the proposed rule based similarity using link based cluster approach. A brief explanation about its phases is followed.



## 2. CREATING A CLUSTER ENSEMBLE

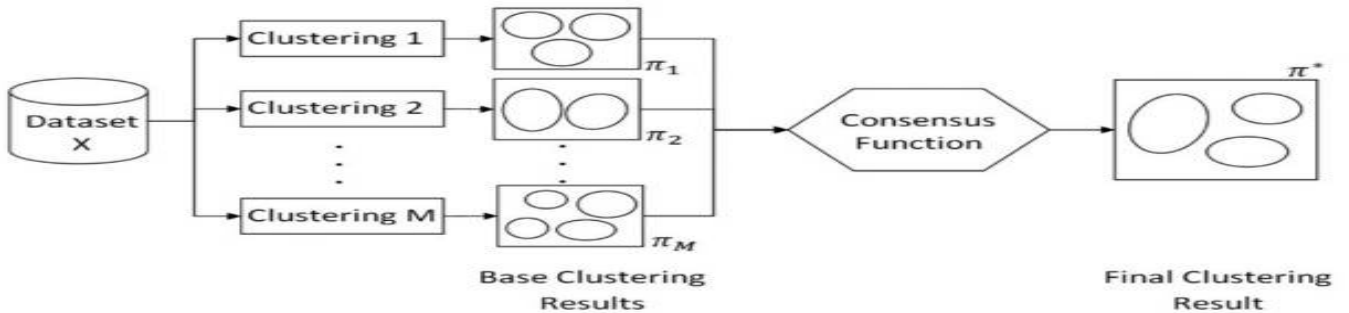


Figure 3. shows the Basic Process of Cluster Ensembles.

### TYPE I (DIRECT ENSEMBLE)

The First type of cluster ensemble transforms the problem of categorical data clustering to cluster ensembles by considering each categorical attribute value (or label) as a cluster in an ensemble. Let  $X = \{x_1 \dots x_n\}$  be a set of  $N$  data points,  $A = \{a_1 \dots a_m\}$  be a set of categorical attributes, and  $\Pi = \{\pi_1, \dots, \pi_M\}$  be a set of  $M$  partitions. Each partition  $\pi_i$  is generated for a specific categorical attribute  $a_i \in A$ . Clusters belonging to a partition  $\pi_i = \{C_1^i, \dots, C_{k_i}^i\}$  correspond to different values of the attribute  $a_i = \{a_1^i, \dots, a_{k_i}^i\}$ , where  $\bigcup_{j=1}^{k_i} C_j^i = a_i$  and  $k_i$  is the number of values of attribute  $a_i$ . With this formalism, categorical data  $X$  can be directly transformed to a cluster ensemble, without actually implementing any base clustering [14].

### TYPE II (FULL-SPACE ENSEMBLE)

In this two ensemble types are created from base clustering results, each of which is obtained by applying a clustering algorithm to the categorical data set. In particular to a full-space ensemble, base clusterings are created from the original data, i.e., with all data attributes. To introduce an artificial instability to k-modes, the following two schemes are employed to select the number of clusters in each base clusterings: 1) Fixed-k,  $k = \lceil \sqrt{N} \rceil$  (where  $N$  is the number of data points), and 2) Random-k,  $k \in \{2, \dots, \lceil \sqrt{N} \rceil\}$  [5].

### TYPE III: SUBSPACE ENSEMBLE

Another alternative to generate diversity within an ensemble is to exploit a number of different data subsets. To this extent, the cluster ensemble is established on various data subspaces, from which base clustering results are

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 12, December 2016

generated. Similar to the study in, for a given  $N * d$  data set of  $N$  data points and  $d$  attributes, an  $N * q$  data subspace (where  $q < d$ ) is generated by  $q = q_{min} + \lfloor \alpha (q_{max} - q_{min}) \rfloor$ . Where  $\alpha \in [0, 1]$  is a uniform random variable,  $q_{min}$  and  $q_{max}$  are the lower and upper bounds of the generated subspace, respectively.

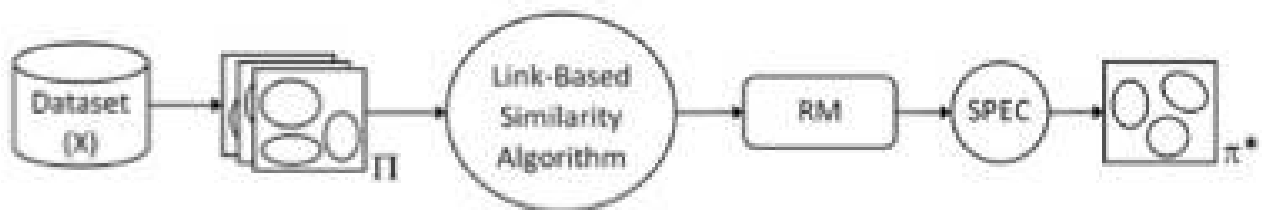
### 3. GENERATING A REFINED MATRIX

Generating a refined cluster-association matrix (RM) using a link-based similarity algorithm. Cluster ensemble methods are based on the binary cluster-association matrix. Each entry in this matrix represents a association degree between data point. Refined cluster-association matrix is put forward as the enhanced variation of the original BM. Its aim is to approximate the value of unknown associations (“0”) from known ones (“1”), whose association degrees are preserved within the RM [12] [14]. These hidden or unknown associations can be estimated from the similarity among clusters, discovered from a network of clusters.

$$RM(x_i, c_l) = \begin{cases} 1, & \text{if } c_l = C_t(x_i), \\ \text{Sim}(c_l, C_t(x_i)), & \text{otherwise} \end{cases}$$

Where  $C_t(x_i)$  is a cluster label (corresponding to a particular cluster of the clustering  $\pi$ ) to which data point  $x_i$  belongs. In addition,  $\text{sim}(C_x, C_y) \in [0, 1]$  denotes the similarity between any two clusters  $C_x, C_y$ , which can be discovered using the following link-based algorithm. Note that, for any clustering  $\pi \in \Pi$ ,  $1 \leq \sum RM(x_i, C) \leq kt$ . Unlike the measure of fuzzy membership, the typical constraint of  $\sum RM(x_i, C) = 1$  is not appropriate for rescaling associations within the RM.

### 4. WEIGHTED TRIPLE-QUALITY (WTQ): A NEW LINK-BASED SIMILARITY ALGORITHM



The Weighted Triple-Quality algorithm is efficient approximation of the similarity between clusters in a link network. WTQ aims to differentiate the significance of triples and hence their contributions toward the underlying similarity measure. A cluster ensemble of a set of data points  $X$ , a weighted graph  $G=(V,M)$  can be constructed, where  $V$  is the set of vertices each representing a cluster and  $W$  is a set of weighted edges between clusters. The weight assigned to the edge that connects clusters is estimated by the proportion of their overlapping [9] [11]. Members Shared neighbours have been widely recognized as the basic evidence to justify the similarity among vertices in a link network. The method gives high weights to rare features and low weights to features that are common to most of the pages. For WTQ can be modified to discriminate the quality of shared triples between a pair of clusters in question [14]. The quality of each cluster is determined by the rarity of links connecting to other clusters in a network.

$$W_{xy} = \frac{|L_x \cap L_y|}{|L_x \cup L_y|}$$

Following that, the similarity between clusters  $C_x$  and  $C_y$  can be estimated by,

$$\text{Sim}(C_x, C_y) = \frac{WTQ_{xy}}{WTQ_{max}} * DC$$

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 12, December 2016

## 5. CONNECTOR-BASED SIMILARITY MEASURE

We propose a Connector-based similarity measure called C-Rank. C-Rank uses both in-links and out-links at the same time. We propose a new link-based similarity measure called C-Rank, which uses both in-link and out-link by disregarding the direction of references. C-Rank, where  $L(p)$  denotes the set of undirected link neighbors of paper  $p$ . Similar to that the accuracy of Co-citation (Coupling) is improved by iterative Sim Rank (rvs-Sim Rank), Crank is defined iteratively. C-Rank unifies in-links and out-links into undirected links. C-Rank has the effect similar to increasing the weight of Co-citation (SimRank) when computing the score between old papers, increasing the weight of Coupling (rvs-SimRank) when computing the score between recent papers, and increasing the weight of a BP-based similarity measure when computing the score between old and recent papers.

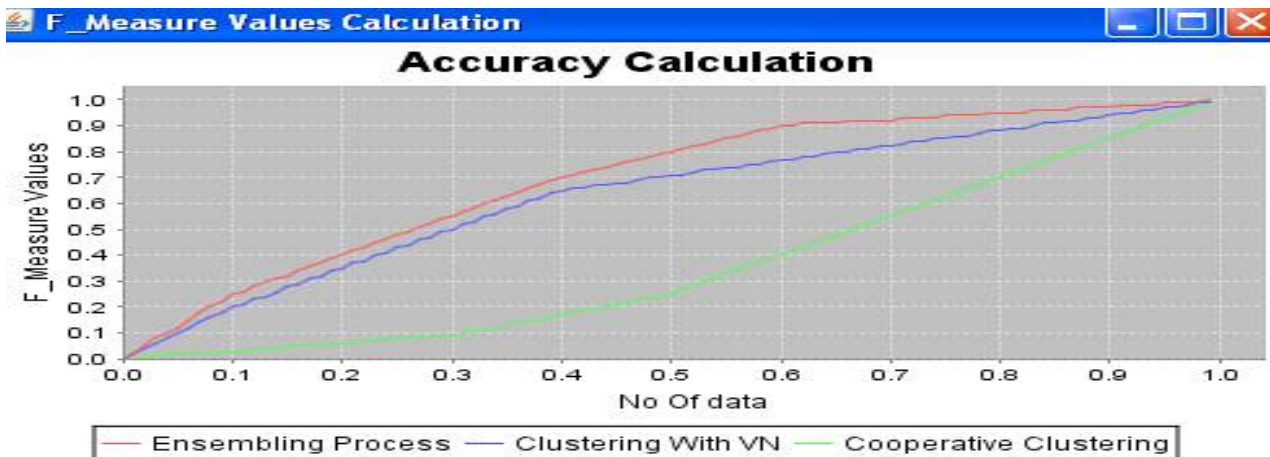
## IV. PERFORMANCE EVALUATION

This section presents the evaluation of the proposed link based method (LCE), using a variety of validity indices and real data sets. The quality of data partitions generated by this technique is assessed against those created by different categorical data clustering algorithms and cluster ensemble techniques.

### INVESTIGATED DATA SETS

The experimental evaluation is conducted over one data sets. The “UCI Machine learning repository” data set is a subset of the well known attribute values collection — UCI Machine learning repository.

## V. EXPERIMENT RESULTS

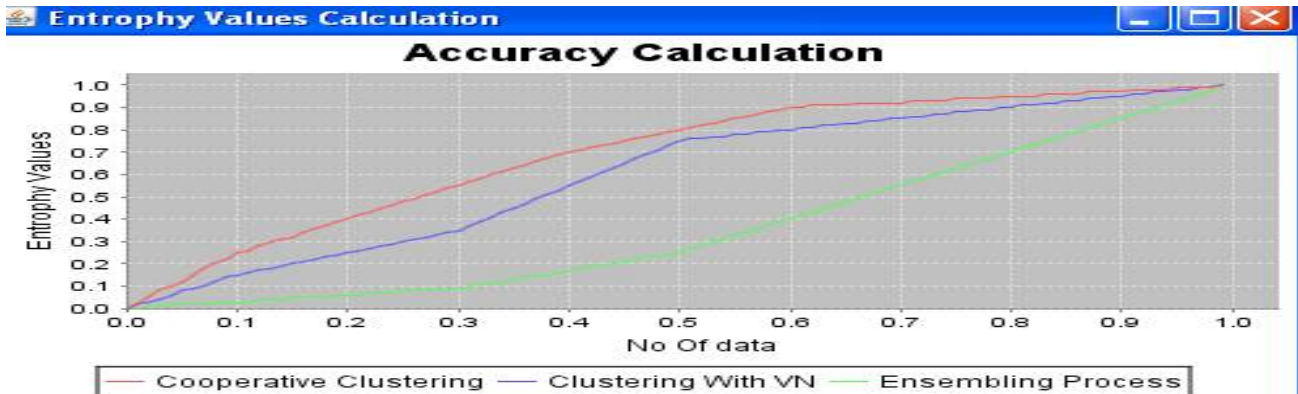


The above figure Shows that F-Measures values calculation for Ensembling process (Crank cluster similarity), clustering with VN (Weighted Triple Quality (WTQ)), Cooperative clustering (CO+SL).

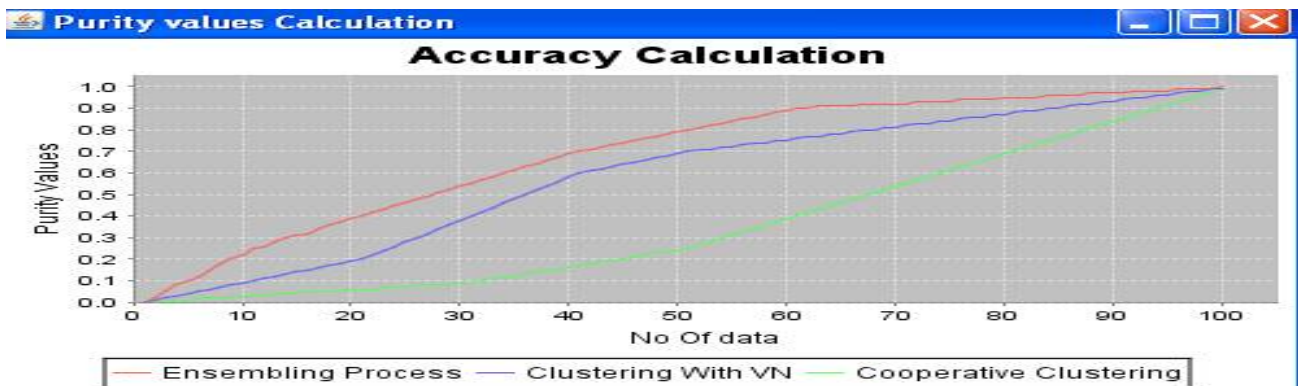
# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

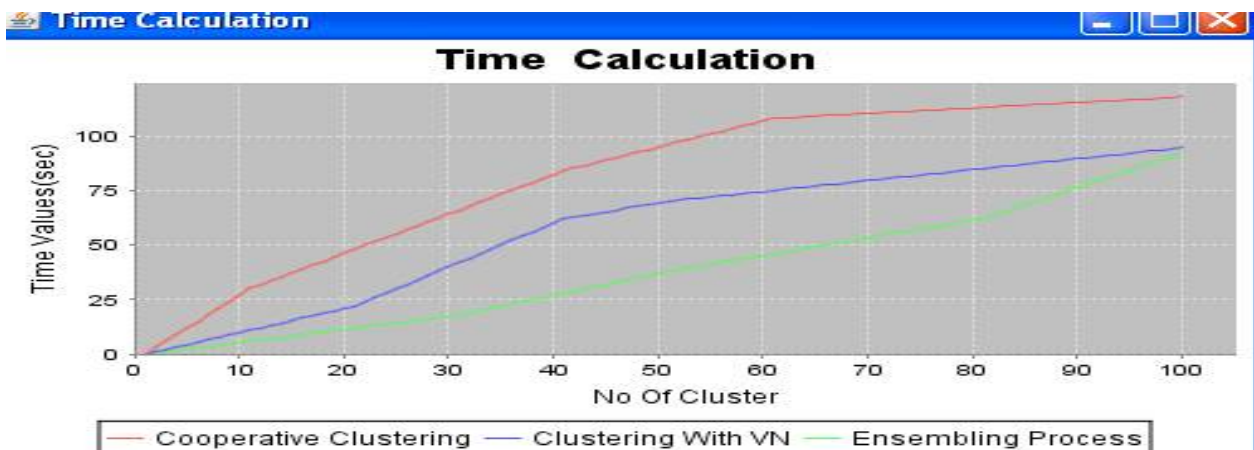
Vol. 5, Issue 12, December 2016



The above figure Shows that Entropy values calculation for Ensembling process (Crank cluster similarity), clustering with VN (Weighted Triple Quality (WTQ)), Cooperative clustering (CO+SL).



The above figure Shows that Purity values calculation for Ensembling process (Crank cluster similarity), clustering with VN (Weighted Triple Quality (WTQ)), Cooperative clustering (CO+SL).



The above figure Shows that Time calculation for Ensembling process (Crank cluster similarity), clustering with VN (Weighted Triple Quality (WTQ)), Cooperative clustering (CO+SL).

# International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 12, December 2016

## VI. CONCLUSION

This paper presents a novel, highly effective link-based cluster ensemble approach (WTQ) to categorical data clustering. It transforms the original categorical data matrix to an information-preserving numerical variation (RM), to which an effective graph partitioning technique can be directly applied. The problem of constructing the RM is efficiently resolved by the similarity among categorical labels (or clusters), using the Weighted Triple-Quality similarity algorithm. The empirical study, with different ensemble types, validity measures, and data sets, suggests that the proposed link-based method usually achieves superior clustering results compared to those of the traditional categorical data algorithms and benchmark cluster ensemble techniques. It also presents a Crank link based cluster approach for categorical data clustering.

## REFERENCES

- [1] N. Nguyen and R. Caruana, "Consensus Clusterings," Proc. IEEE Int'l Conf. Data Mining (ICDM), pp. 607-612, 2007.
- [2] 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.
- [3] C. Boulis and M. Ostendorf, "Combining Multiple Clustering Systems," Proc. European Conf. Principles and Practice of Knowledge Discovery in Databases (PKDD), pp. 63-74, 2004.
- [4] 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.
- [5] 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.
- [6] 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>, 2007.
- [7] Y. Zhang, A. Fu, C. Cai, and P. Heng, "Clustering Categorical Data," Proc. Int'l Conf. Data Eng. (ICDE), p. 305, 2000.
- [8] M. Dutta, A.K. Mahanta, and A.K. Pujari, "QROCK: A Quick Version of the ROCK Algorithm for Clustering of Categorical Data," Pattern Recognition Letters, vol. 26, pp. 2364-2373, 2005.
- [9] E. Abdu and D. Salane, "A Spectral-Based Clustering Algorithm for Categorical Data Using Data Summaries," Proc. Workshop Data Mining using Matrices and Tensors, pp. 1-8, 2009.
- [10] B. Mirkin, "Reinterpreting the Category Utility Function," Machine Learning, vol. 45, pp. 219-228, 2001.