

A Max Metric to Evaluate a Cluster

Hosein Alizadeh¹, Hamid Parvin², Sajad Parvin², Zahra Rezaei², Moslem mohamadi²

¹Islamic Azad University

Mahdi Shahr Branch

Mahdi Shahr, Iran

²Islamic Azad University

Nourabad Mamasani Branch

Mamasani Nourabad, Iran

{s.parvin, rezaei, mohamadi, halizadeh}@iust.ac.ir, hamidparvin@mamasaniiau.ac.ir



ABSTRACT: In this paper a new criterion for clusters validation is proposed. This new cluster validation criterion is used to approximate the goodness of a cluster. The clusters which satisfy a threshold of the proposed measure are selected to participate in clustering ensemble. To combine the chosen clusters, some methods are employed as aggregators. Employing this new cluster validation criterion, the obtained ensemble is evaluated on some well-known and standard datasets. The empirical studies show promising results for the ensemble obtained using the proposed criterion comparing with the ensemble obtained using the standard clusters validation criterion. Besides to reach the best results, the method gives an algorithm based on which one can find how to select the best subset of clusters from a pool of clusters.

Keywords: Clustering Ensemble, Stability Measure, Extended EAC, Co-association Matrix, Cluster Evaluation

Received: 12 October 2011, Revised 3 December 2011, Accepted 15 December 2011

© 2012 DLINE. All rights reserved

1. Introduction

Data clustering or unsupervised learning is an important and very difficult problem. The objective of clustering is to partition a set of unlabeled objects into homogeneous groups or clusters [3], [4] and [10]. There are many applications that use clustering techniques to discover latent structures of data, such as data mining [11], information retrieval [2], image segmentation [9], linkage learning [15], and machine learning. In real-world problems, clusters can appear with different shapes, sizes, data sparseness's, and degrees of separation. Clustering techniques require the definition of a similarity measure between patterns. Since there is no prior knowledge about cluster shapes, choosing a specific clustering method is not easy [16]. Studies in the last few years have tended to combinational methods. Cluster ensemble methods attempt to find better and more robust clustering solutions by fusing information from several primary data partitions [8].

Fern and Lin [8] have suggested a clustering ensemble approach which selects a subset of solutions to form a smaller but better-performing cluster ensemble than using all primary solutions. The ensemble selection method is designed based on quality and diversity, the two factors that have been shown to influence cluster ensemble performance. This method attempts to select a subset of primary partitions which simultaneously has both the highest quality and the most diversity. The Sum of Normalized Mutual Information, SNMI [5], [6] and [17], is used to measure the quality of each individual partition with respect to other partitions. Also, the Normalized Mutual Information, NMI, is employed to measure the diversity among partitions. Although the ensemble size in this method is relatively small, this method achieves significant performance improvement over full ensembles.

Law et al. proposed a multi-objective data clustering method based on the selection of individual clusters produced by several clustering algorithms through an optimization procedure [13]. This technique chooses the best set of objective functions for different parts of the feature space from the results of base clustering algorithms. Fred and Jain [7] have offered a new clustering ensemble method which learns the pairwise similarities between points in order to facilitate a proper partition of the data without the a priori knowledge of the number and the shape of the clusters. This method which is based on cluster stability evaluates the primary clustering results instead of final clustering.

We propose a new criterion for clusters validation. Then we employ this criterion to select the more robust clusters in the final ensemble. We also propose a new method named Extended Evidence Accumulation Clustering, EEAC, to construct the matrix of similarity from these selected clusters. Finally, we apply a hierarchical method over the obtained matrix to extract the final partition.

Rest of this paper is organized as follows. In section 2, we explain the proposed method. Section 3 demonstrates results of our proposed method against traditional comparatively. Finally, we conclude in section 4.

2. Proposed Method

In this section, first our proposed clustering ensemble method is briefly outlined, and then its phases are described in detail. The main idea of our proposed clustering ensemble framework is utilizing a subset of best performing primary clusters in the ensemble, rather than using all of clusters. Only the clusters that satisfy a stability criterion can participate in the combination. The cluster stability is defined according to Normalized Mutual Information, NMI. Figure 1 depicts the proposed clustering ensemble procedure.

The manner of computing stability is described in the following sections in detail. To select a subset with the most stable clusters for combination, we apply a stability-threshold to each cluster. Different sizes of the most stable clusters are explored to find the best option. After selection phase, the selected clusters are used to construct the co-association matrix. Several methods have been proposed for combination of the primary results [1] and [17]. In this work, some clusters in the primary partitions may be absent (having been eliminated by the stability criterion). Since the original EAC method [5] cannot truly identify the pairwise similarity while there is only a subset of clusters, we present a new method for constructing the co-association matrix. We call this method: Extended Evidence Accumulation Clustering method, EEAC. Finally, we use a hierarchical clustering algorithm, like single-link method, to extract the final clusters out of this matrix. For more generality, some heuristic consensus functions are also used as aggregators of selected clusters [17]. These heuristic consensus functions that are based on hypergraph partitioning and have first introduced by Strehl and Ghosh, are HperGraph Partitioning Algorithm (HGPA), Meta-Clustering Algorithm (MCLA) and Cluster-based Similarity Partitioning Algorithm (CSPA) [17].

2.1 Cluster Evaluation

Since goodness of a cluster is determined by all the data points, the goodness function $g_j(C_i, D)$ depends on both the cluster C_i and the entire dataset D , instead of C_i alone. The stability as measure of cluster goodness is used in [12]. Cluster stability reflects the variation in the clustering results under perturbation of the data by resampling.

A stable cluster is one that has a high likelihood of recurrence across multiple applications of the clustering method. Stable clusters are usually preferable, since they are robust with respect to minor changes in the dataset [13].

Now assume that we want to compute the stability of cluster C_i . In this method first a set of partitionings over resampled datasets is provided which is called the reference set. In this notation D is resampled data and $P(D)$ is a partitioning over D . Now, the problem is: "How many times is the cluster C_i repeated in the reference partitions?" Denote by $NMI(C_i, P(D))$, the Normalized Mutual Information between the cluster C_i and a reference partition $P(D)$. Most previous works only compare a *partition with another partition* [17]. However, the stability used in [13] evaluates the similarity between a *cluster and a partition* by transforming the cluster C_i to a partition and employing common partition to partition methods. To illustrate this method let $P_1 = P^a = \{C_i, D/C_i\}$ be a partition with two clusters, where D/C_i denotes the set of data points in D that are not in C_i .

Then we may compute a second partition $P_2 = P^b = \{C^*, D/C^*\}$, where C^* denotes the union of all "positive" clusters in $P(D)$ and

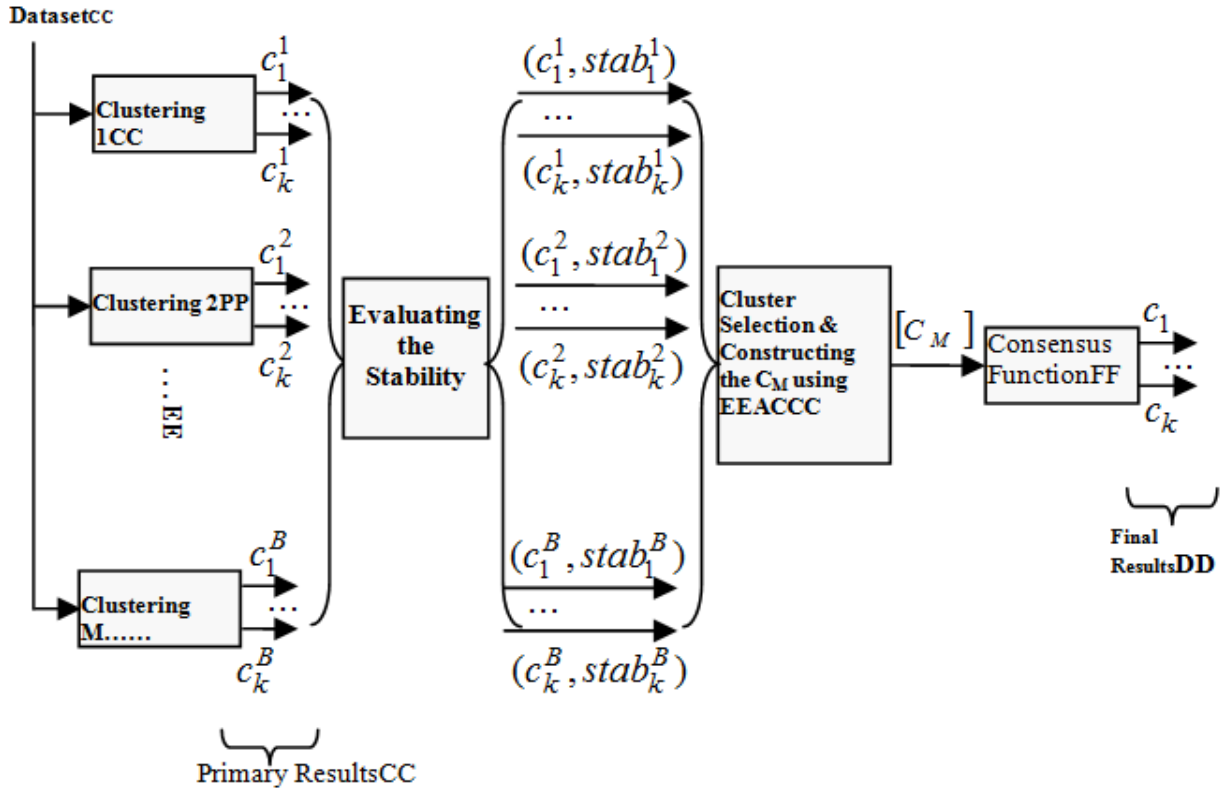


Figure 1. Training phase of the Bagging method

others are in D/C^* . A cluster C_j in $P(D)$ is positive if more than half of its data points are in C_j . Now, define $NMI(C_j, P(D))$ by $NMI(P^a, P^b)$ which is calculated as [6]:

$$NMI(P^a, P^b) = \frac{-2 \sum_{i=1}^{k_a} \sum_{j=1}^{k_b} n_{ij}^{ab} \log \left(\frac{n_{ij}^{ab} \cdot n}{n_i^a \cdot n_j^b} \right)}{\sum_{i=1}^{k_a} n_i^a \log \left(\frac{n_i^a}{n} \right) + \sum_{j=1}^{k_b} n_j^b \log \left(\frac{n_j^b}{n} \right)} \quad (1)$$

where n is the total number of samples and n_{ij}^{ab} denotes the number of shared patterns between clusters $C_i^a \in P^a$ and $C_j^b \in P^b$; n_i^a is the number of patterns in the cluster i of partition a ; also n_j^b are the number of patterns in the cluster j of partition b .

This computation is done between the cluster C_i and all partitions available in the reference set. Figure 2 shows this method.

NMI_i in Figure 2 shows the stability of cluster C_i with respect to the i -th partition in reference set. The total stability of cluster C_i is defined as:

$$Stability(C_i) = \frac{1}{M} \sum_{i=1}^M NMI_i \quad (2)$$

where M is the number of partitions available in reference set. This procedure is applied for each cluster of every primary partition.

2.2 Max Method

In this section a drawback of computing stability is introduced and an alternative approach is suggested which is named Max method. Figure 3 shows two primary partitions for which the stability of each cluster is evaluated. In this example K-means is applied as the base clustering algorithm with $K = 3$. For this example the number of all partitions in the reference set is 40. In 36 partitions the result is relatively similar to Figure 3a, but there are four partitions in which the top left cluster is divided into two clusters, as shown in Figure 3b. Figure 3a shows a true clustering. Since the well separated cluster in the top left corner is repeated several times (90% repetition) in partitions of the reference set, it has to acquire a great stability value (but not equal to 1), however it acquires the stability value of 1. Because the two clusters in right hand of Figure 3a are relatively joined and sometimes they are not recognized in the reference set as well, they have less stability value. Figure 3b shows a spurious clustering which the two right clusters are incorrectly merged. Since a fixed number of clusters are forced in the base algorithm, the top left cluster is divided into two clusters. Here the drawback of the stability measure is apparent rarely. Although it is obvious that this partition and the corresponding large cluster on the right reference set (10% repetition), the stability of this cluster is evaluated equal to 1. Since the NMI is a symmetric equation, the stability of the top left cluster in Figure 3a is exactly equal to the large right cluster in Figure 3b; however they are repeated 90% and 10%, respectively. In other words, when two clusters are complements of each other, their stabilities are always equal. This drawback is seen when the number of positive clusters in the considered partition of reference set is greater than 1. It means when the cluster C^* is obtained by merging two or more clusters, undesirable stability effects occur.

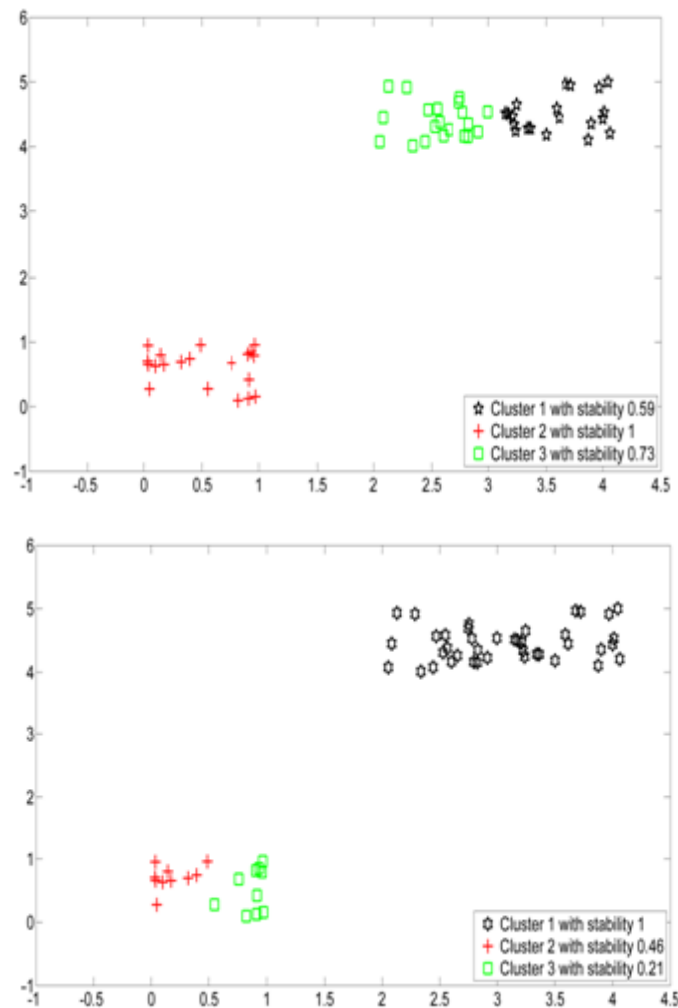


Figure 3. Two primary partitions with $k=3$. (a, Left) True clustering. (b, Right) Spurious clustering

To solve this problem we allow only one cluster in reference set to be considered as the C^* (i.e. only the most similar cluster) and all others are considered as D/C^* . In this method the problem is solved by eliminating the merged clusters.

2.3 Consensus Function

One way is to consider the selected clusters as inputs of the HGPA, MCLA and CSPA algorithms [17]. The output of the mentioned algorithms is the final partition which is also called consensus partition.

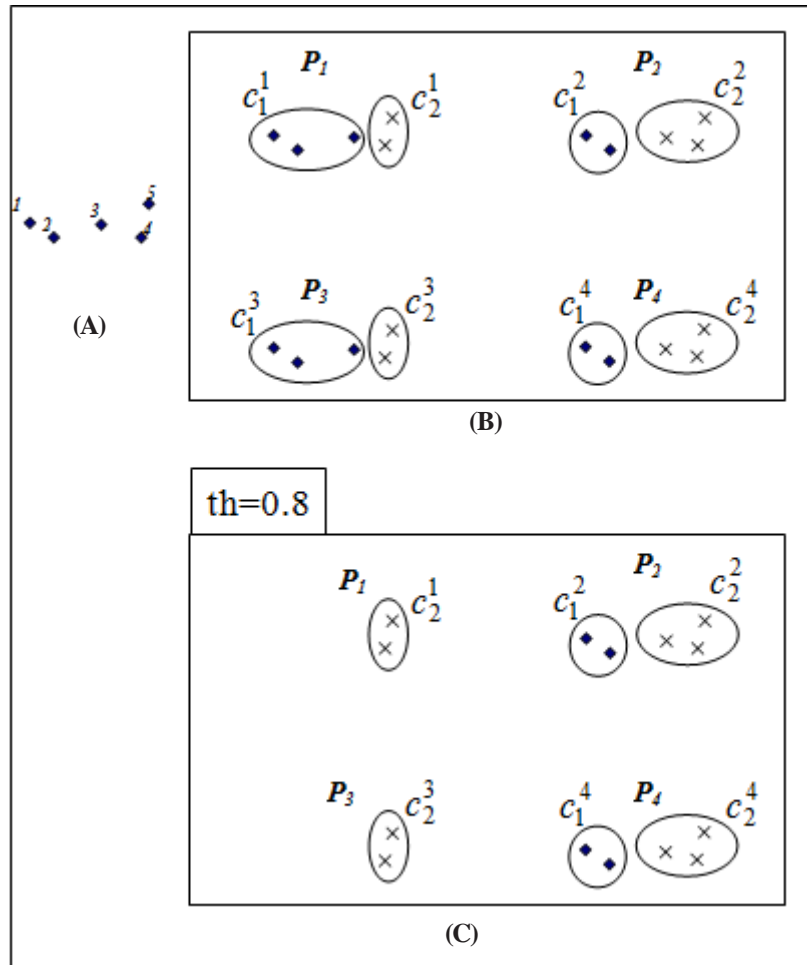


Figure 4. Computing the co-association matrix by the EEAC method. (A) Data samples. (B) 4 primary clusterings. (C) Remaining clusters after applying threshold, $th=0.8$

For the second way to extract the final partition from the selected clusters, the clusters are considered as new space for data, and a clustering algorithm, like fuzzy k-means, is employed to partition the mapped data.

Another alternative way to reach the consensus partition is to use the co-association based methods. In this method, the selected clusters are first used to construct the co-association matrix. In the EAC method the m primary results from resampled data are accumulated in an $n \times n$ co-association matrix. Each entry in this matrix is computed from equation (3).

$$C(i, j) = \frac{n_{i,j}}{m_{i,j}} \quad (3)$$

where n_{ij} counts the number of clusters shared by objects with indices i and j in the partitions over the B clusterings. Also m_{ij} is the number of partitions where this pair of objects is simultaneously present. There are only a fraction of all primary clusters available, after thresholding. So, the common EAC method cannot truly recognize the pairwise similarity for computing the co-association matrix. In our novel method (Extended Evidence Accumulation Clustering, or EEAC) each entry of the co-association matrix is computed by:

$$C(i, j) = \frac{n_{i,j}}{\max(n_i, n_j)} \quad (4)$$

where n_i and n_j are the number present in remaining (after stability thresholding) clusters for the i -th and j -th data points, respectively. Also, n_{ij} counts the number of remaining clusters which are shared by both data points indexed by i and j , respectively. To further explain, consider this example. Assume that we have five samples (Fig 4a), and that four primary clustering are applied (Figure 4b).

Also, suppose that that stability of the clusters of Figure 4b is as given bellow:

$$\begin{aligned} \text{stability}(c_2^1) &= \text{stability}(c_2^3) = 1 \\ \text{stability}(c_1^2) &= \text{stability}(c_1^4) = 1 \\ \text{stability}(c_2^2) &= \text{stability}(c_2^4) = 0.82 \\ \text{stability}(c_1^1) &= \text{stability}(c_1^3) = 0.55 \end{aligned}$$

By choosing $th = 0.8$ the first clusters from P1 and P3 are deleted (Figure 4c). According to equation (4), each entry of the co-association matrix is:

$$\begin{aligned} C(1,2) &= \frac{2}{\max(2,2)} = \frac{2}{2} = 1 & C(1,3) = C(2,3) &= \frac{0}{\max(2,2)} = \frac{0}{2} = 0 \\ C(3,4) = C(3,5) &= \frac{2}{\max(2,4)} = \frac{2}{4} = 0.5 & C(4,5) &= \frac{4}{\max(4,4)} = \frac{4}{4} = 1 \end{aligned}$$

In Figure 4a-c, the data points may be “tracked” by their geometrical arrangement. Example: in computing $C(3,4)$, note that points 3 and 4 both are in cluster 2 of partitions P2 and P4, so that numerator $n_{34} = 2$; also note that $n_3 = 2$, since point 3 is only in cluster 2 of P2 and P4, but $n_4 = 4$ since point 4 is not only in these clusters, but also in cluster 2 of P1 and P3. Before and after applying threshold, the co-association matrix is given by equation (5) and (6), respectively:

$$C_{before} = \begin{bmatrix} 1 & 1 & 0.5 & 0 & 0 \\ 1 & 1 & 0.5 & 0 & 0 \\ 0.5 & 0.5 & 1 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 1 & 1 \\ 0 & 0 & 0.5 & 1 & 1 \end{bmatrix} \quad (5)$$

In this matrix the 3rd object can be considered as both clusters with an equal probability of 50%. The stability measure adds some information to this matrix by applying the threshold.

$$C_{after} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 1 & 1 \\ 0 & 0 & 0.5 & 1 & 1 \end{bmatrix} \quad (6)$$

By comparing these two matrices and also considering the stability values, it can be seen that deletion of unstable clusters improves the co-association matrix. By eliminating the unstable cluster with samples $\{1, 2, 3\}$ which is spuriously created by primary clusterings.

After computing the co-association matrix by the EEAC method, a consensus function is employed to extract the final clusters from the matrix. Here, the single-link method is used for this task.

3. Experimental Study

Evaluation metric based on which a consensus partition is evaluated is discussed in the first subsection of this section. The details of the used datasets are given in the subsequent section. Then the settings of experimentations are given. Finally the experimental results are presented.

| | <i>Dataset Name</i> | <i># of Class</i> | <i># of Features</i> | <i># of Samples</i> |
|----|---------------------|-------------------|----------------------|---------------------|
| 1 | Breast-Cancer* | 2 | 9 | 683 |
| 2 | Iris* | 3 | 4 | 150 |
| 3 | Bupa* | 2 | 6 | 345 |
| 4 | SAHeart* | 2 | 9 | 462 |
| 5 | Ionosphere | 2 | 34 | 351 |
| 6 | Glass* | 6 | 9 | 214 |
| 7 | Half rings | 2 | 2 | 400 |
| 8 | Galaxy* | 7 | 4 | 323 |
| 9 | Yeast* | 10 | 8 | 1484 |
| 10 | Wine | 3 | 13 | 178 |

Table 1. Brief information about the used datasets

3.1 Evaluation Metric

After producing the consensus partition, the most important question is “how good a partition is?”. The evaluation of a partition is very important as it is mentioned. Here the NMI between the consensus partition and real labels of the dataset is considered as an evaluation metric of the consensus partition. Also accuracy between the consensus partition and real labels of the dataset is considered as another metric.

3.2 Datasets

The proposed method is examined over 9 different standard datasets and one artificial dataset. It is tried for datasets to be diverse in their number of true classes, features and samples. A large variety in used datasets can more validate the obtained results. Brief information about the used datasets is available in Table 1. More information is available in [14].

Note that some of datasets which are marked with star (*) in Table 1 are normalized. All experiments are done over the normalized features in the starred dataset. It means each feature is normalized with mean of 0 and variance of 1, $N(0, 1)$. The artificial Half Ring dataset is depicted in the Figure 5.

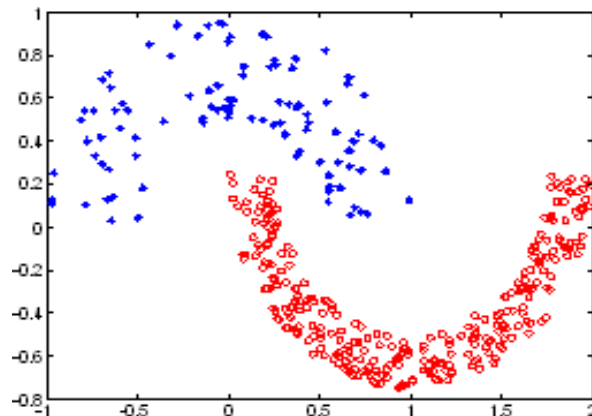


Figure 5. Half Ring dataset

3.3 Experimental Settings

To be more general and fair, all experiments are averaged over 10 independent runs. In all experimentations there are 120 independent partitions obtained by 120 independent runs of k-means clustering algorithm with different initialized seed points and different k parameter, ranging from k to 2*k.

After selecting a subset of clusters, to extract the final partition from them, the real number of clusters, i.e. the column three of the Table 1, is served by the consensus functions.

As it is known in fuzzy k-means clustering algorithm, each data point belongs to all clusters with different membership values. To extract the final partition from output of fuzzy k-means algorithm as consensus function, each data point is assigned to the most membership value.

3.4 Experimental Results

As it is inferred from the Figure 6, the best ratio of selection of the stable clusters is 60% and the best option for consensus function is CSPA for Iris dataset.

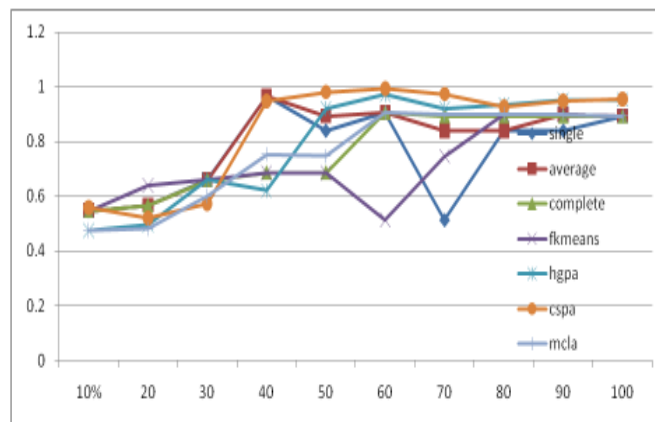


Figure 6. The horizontal axis stands for the rate of stable clusters that are selected. The vertical axis stands for the NMI values between the labels of Iris dataset and the consensus partitions obtained by different consensus functions over the selected clusters

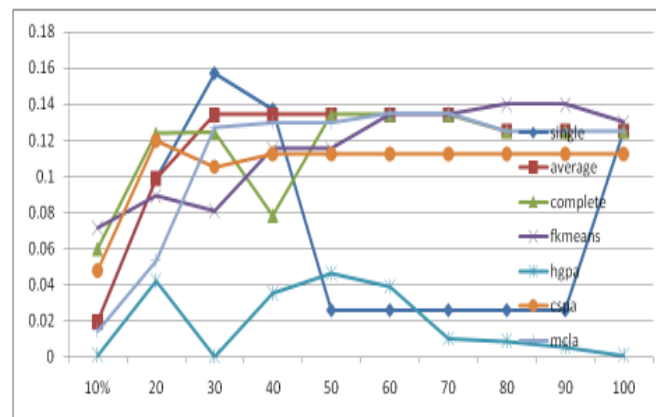


Figure 7. The horizontal axis stands for the rate of stable clusters that are selected. The vertical axis stands for the NMI values between the labels of Ionosphere dataset and the consensus partitions obtained by different consensus functions over the selected clusters

Figure 7 makes it clear that the best ratio of selection of the stable clusters is 30% and the best option for consensus function is Single-Linkage for Ionosphere dataset.

To see whether the use of a subset of the most stable clusters can affect the quality of the final cluster or not, consider Figure 8. To make a general decisive conclusion, the results for all ten datasets of Table 1 are averaged and the final results are illustrated in the Figure 8. The Averaged-Linkage consensus function over 50% of the most stable clusters generally reaches the maximum for all dataset.

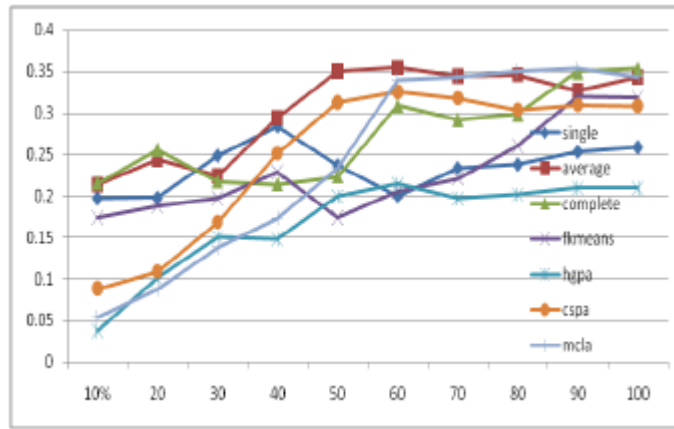


Figure 8. The horizontal axis stands for the rate of stable clusters that are selected. The vertical axis stands for the averaged NMI values for all ten datasets of Table 1

Table 2 shows the performance of the proposed method comparing with most common base and ensemble methods.

| Dataset | Simple Methods (%) | | | | Ensemble Methods (%) | | | |
|---------|--------------------|-----------------|------------------|--------|----------------------|---------------|---------------------------------|---------------------------------|
| | Single Linkage | Average Linkage | Complete Linkage | Kmeans | Kmeans Ensemble | Full Ensemble | Cluster Selection by NMI Method | Cluster Selection by Max Method |
| Wine | 37.64 | 38.76 | 83.71 | 96.63 | 96.63 | 97.08 | 97.75 | 97.44 |
| BreastC | 65.15 | 70.13 | 94.73 | 95.37 | 95.46 | 95.10 | 95.75 | 96.49 |
| Yeast | 34.38 | 35.11 | 38.91 | 40.20 | 45.46 | 47.17 | 47.17 | 51.27 |
| Glass | 36.45 | 37.85 | 40.65 | 45.28 | 47.01 | 47.83 | 48.13 | 47.35 |
| Bupa | 57.68 | 57.10 | 55.94 | 54.64 | 54.49 | 55.83 | 58.09 | 58.40 |

Table 2. Experimental results

4. Conclusion and Future Works

In this paper a new clustering ensemble framework is proposed which is based on a subset of total primary spurious clusters. Also a new alternative method for common NMI is suggested. Since the quality of the primary clusters are not equal and presence of some of them can even yield to lower performance, here a method to select a subset of more effective clusters is proposed. A common cluster validity criterion which is needed to derive this subset is based on normalized mutual information. In this paper some drawbacks of this criterion is discussed and an method is suggested which is called max mehod. The experiments show that the proposed framework commonly outperforms in comparison with the full ensemble; however it uses just 50% of primary clusters. Another innovation of this chapter is a method for constructing the co-association matrix where some of clusters and respectively some of samples do not exist in partitions. This new method is called Extended Evidence Accumulation Clustering, EEAC.

References

[1] Ayad, H., Kamel, M. S. (2008). Cumulative Voting Consensus Method for Partitions with a Variable Number of Clusters, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 30 (1) 160-173.

[2] Bhatia, S. K., Deogun, J. S. (1998). Conceptual Clustering in Information Retrieval, *IEEE Trans. Systems, Man, and Cybernetics*, 28 (3) 427-536.

- [3] Dudoit, S., Fridly, J. (2003). Bagging to improve the accuracy of a clustering procedure, *Bioinformatics*, 19(9), 1090-1099.
- [4] Faceli K., Marcilio C. P. Soutod. (2006). Multi-objective Clustering Ensemble, *In: Proceedings of the Sixth International Conference on Hybrid Intelligent Systems (HIS'06)*.
- [5] Fred, A., Jain, A. K. (2002). Data Clustering Using Evidence Accumulation, *In: Proc. of the 16th Intl. Conf. on Pattern Recognition, ICPR02, Quebec City*, p. 276 – 280.
- [6] Fred, A., Jain A. K. (2005). Combining Multiple Clusterings Using Evidence Accumulation, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 27 (6) 835–850.
- [7] Fred, A., Jain, A. K. (2006). Learning Pairwise Similarity for Data Clustering, *In: Proc. of the 18th Int. Conf. on Pattern Recognition (ICPR'06)*.
- [8] Fred, A., Lourenco, A. (2008), Cluster Ensemble Methods: from Single Clusterings to Combined Solutions, *Studies in Computational Intelligence (SCI)*, 126, 3–30.
- [9] Frigui, H., Krishnapuram, R. (1999). A Robust Competitive Clustering Algorithm with Applications in Computer Vision, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 21 (5) 450-466.
- [10] Jain, A. K., Murty, M. N., Flynn, P. (1999). Data clustering: A review. *ACM Computing Surveys*, 31(3), 264–323.
- [11] Judd, D., Mckinley, P., Jain, A.K. (1997). Large-Scale Parallel Data Clustering, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19 (2) 153-158.
- [12] Lange, T., Braun M. L., Roth V., Buhmann J. M. (2003). Stability-based model selection. In *Advances in Neural Information Processing Systems 15*. MIT Press.
- [13] Law, M. H. C., Topchy, A. P., Jain A. K. (2004). Multiobjective data clustering. *In: Proc. of IEEE Conference on Computer Vision and Pattern Recognition, V. 2*, p. 424–430, Washington D.C.
- [14] Newman, C. B. D. J., Hettich S., Merz, C. (1998). UCI repository of machine learning databases, <http://www.ics.uci.edu/~mlern/MLSummary.html>.
- [15] Parvin, H., Minaei-Bidgoli, B., Alinejad, H. (2011). Linkage Learning Based on Differences in Local Optimums of Building Blocks with One Optima. *International Journal of the Physical Sciences, IJPS*, p. 3419 – 3425.
- [16] Roth, V., Lange, T., Braun, M., Buhmann, J. (2002). A Resampling Approach to Cluster Validation, *Intl. Conf. on Computational Statistics, COMPSTAT*.
- [17] Strehl, A., Ghosh, J. (2002). Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3(Dec), 583–617