

A NEW CLUSTERING ENSEMBLE FRAMEWORK

Hamid Parvin, Hosein Alizadeh, Sajad Parvin

Islamic Azad University, Nourabad Mamasani Branch, Nourabad Mamasani, Iran

{ parvin, halizadeh, s.parvin }@mamasaniiu.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 this measure are selected to participate in clustering ensemble. For combining the chosen clusters, a co-association based consensus function is applied. Since the EAC method cannot derive the co-association matrix from a subset of clusters, a new EAC based method which is called Extended EAC, EEAC, is applied for constructing the co-association matrix from the subset of clusters. Employing this new cluster validation criterion, the obtained ensemble is evaluated on some well-known and standard data sets. 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.

Keywords: *Clustering Ensemble, Stability Measure, Cluster Evaluation.*

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]. There are many applications which use clustering techniques for discovering structure in data, such as data mining [10], information retrieval [2], image segmentation [9], and machine learning. In real-world problems, clusters can appear with different shapes, sizes, data sparseness, 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 [14]. 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 partitionings [8].

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.

Fern and Lin [4] 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 diversity. The Sum of Normalized Mutual Information, SNMI [15], [5] and [6], is used to measure the quality of an individual partition with respect to other partitions. Also, the Normalized Mutual Information, NMI, is employed for measuring 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 [12]. 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 similarity between points in order to facilitate a proper partitioning of the data without the a priori knowledge of the number of clusters and of the shape of

these clusters. This method which is based on cluster stability evaluates the primary clustering results instead of final clustering.

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.

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 which 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. After, a subset of the most stable clusters is selected for combination. This is simply done by applying a stability-threshold to each cluster. In the next step, the selected clusters are used to construct the co-association matrix. Several methods have been proposed for combination of the primary results [1] and [15]. In our 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 the hierarchical single-link clustering to extract the final clusters from this matrix.

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 [11]. Cluster stability reflects the variation in the clustering results under perturbation of the data by resampling.

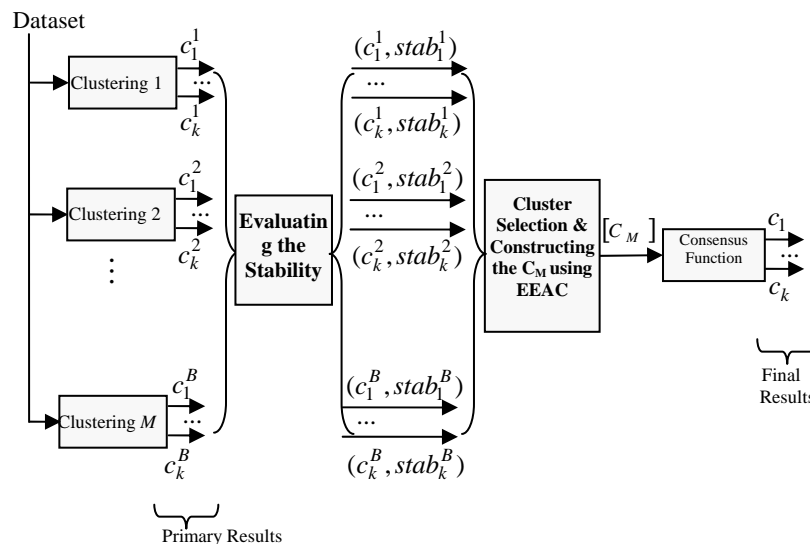


Fig.1 Training phase of the Bagging method

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 [12].

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* [15]. However, the stability used in [12] 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 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_i . Now, define $NMI(C_i, P(D))$ 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. Fig. 2 shows this method.

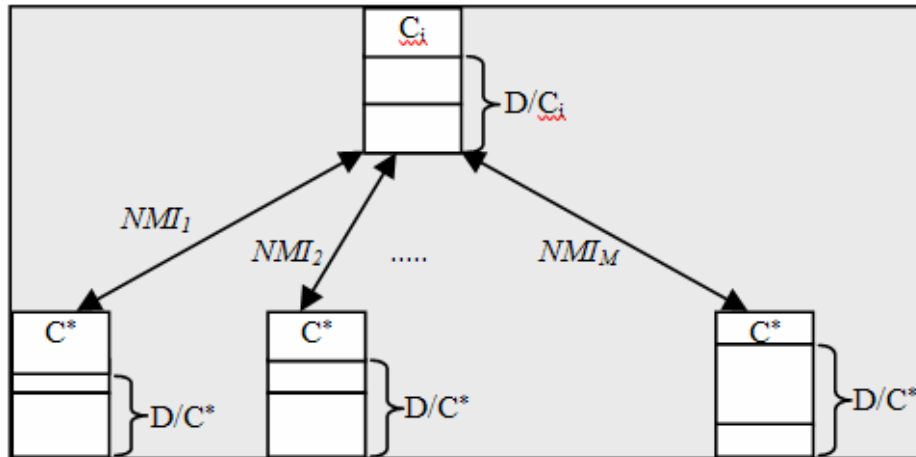


Fig.2 Computing the Stability of Cluster C_i .

NMI_i in Fig. 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. Fig. 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 Fig 3a, but there are four partitions in which the top left cluster is divided into two clusters, as shown in Fig 3b. Fig 3a shows a true clustering. Since the well separated cluster in the top left corner is repeated several times (90% repetition) in partitionings 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 Fig 3a are relatively joined and sometimes they are not recognized in the reference set as well, they have less stability value. Fig. 3.b shows a spurious clustering which the two right clusters are incorrectly merged. Since a fixed number of clusters is 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 fig 3.a is exactly equal to the large right cluster in fig 3.b; 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.

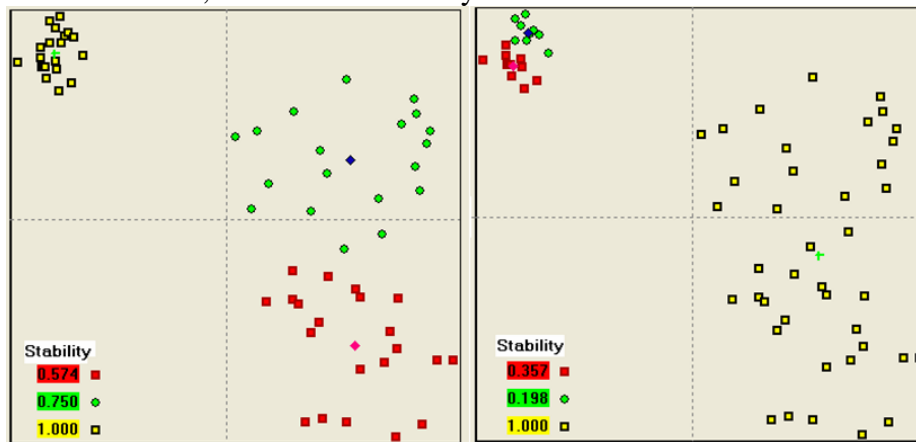


Fig.3 Two primary partitions with $k=3$. (a) True clustering. (b) Spurious clustering

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

In this step, the selected clusters are 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 this equation:

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

where n_{ij} counts the number of clusters shared by objects with indices i and j in the partitions over the primary 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.

Experimental Results

This section reports and discusses the empirical studies. The proposed method is examined over 5 different standard datasets. 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 [13].

Table 1. Brief information about the used datasets.

	Class	Features	Samples
Glass	6	9	214
Breast-C	2	9	683
Wine	3	13	178
Bupa	2	6	345
Yeast	10	8	1484

All experiments are done over the normalized features. It means each feature is normalized with mean of 0 and variance of 1, $N(0, 1)$. All of them are reported over means of 10 independent runs of algorithm. The final performance of the clustering algorithms is evaluated by re-labeling between obtained clusters and the ground truth labels and then counting the percentage of the true classified samples. Table 2 shows the performance of the proposed method comparing with most common base and ensemble methods.

Table 2. Experimental results.

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
Breast-C	65.15	70.13	94.73	95.37	95.46	95.10	95.75	96.49
Wine	37.64	38.76	83.71	96.63	96.63	97.08	97.75	97.44
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

The four first columns of Table 2 are the results of some base clustering algorithms. The results show that although each of these algorithms can obtain a good result over a specific dataset, it does not perform well over other datasets. For example, according to Table 2 the K-means algorithm has a good clustering result over Wine dataset in comparison with linkage methods. But, it has lower performance in comparison to linkage methods in the case of Bupa dataset. Also, the complete linkage has a good performance in Breast-Cancer dataset in comparison with others; however it is not in the case of all datasets. The four last columns show the performance of some ensemble methods in comparison with the proposed one. Taking a glance at the last four columns in comparison with the first four columns shows that the ensemble methods do better than the simple based algorithms in the case of performance and robustness along with different datasets. The first column of the ensemble methods is the results of an ensemble of 100 K-means which is fused by EAC method. The 90% sampling from dataset is used for creating diversity in primary results. The sub-sampling (without replacement) is used as the sampling method. Also the random initialization of the seed points of K-means algorithm helps them to be more diverse. The single linkage

algorithm is applied as consensus function for deriving the final clusters from co-association matrix. The second column from ensemble methods is the full ensemble which uses several clustering algorithms for generating the primary results. Here, 70 K-means with the above mentioned parameters in addition to 30 linkage methods provide the primary results. The third column of the ensemble methods is consensus partitioning using EEAC algorithm of top 33% stable clusters, employing NMI method as measure of stability. The fourth column of the ensemble methods is Also consensus partitioning using EEAC algorithm of top 33% stable clusters, employing max method as measure of stability.

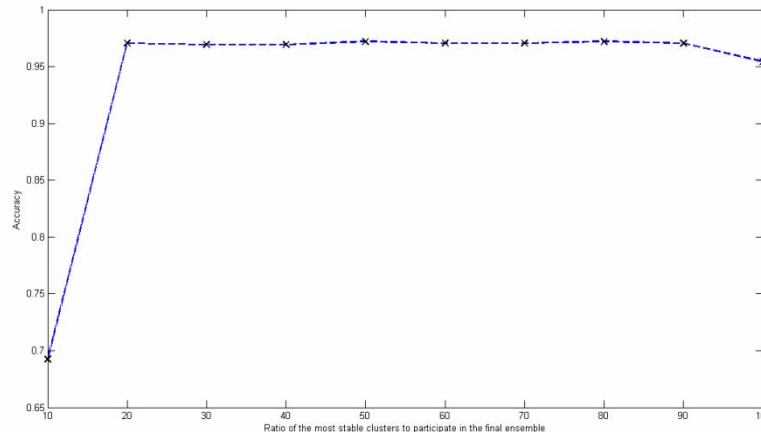


Fig.4 Two primary partitions with $k=3$. (a) True clustering. (b) Spurious clustering

To better understand the effect of proposed clustering ensemble framework, consider Fig. 4 which is different accuracies of the consensus partitions obtained out of different ratios of the most stable clusters in Breast-C dataset. In Fig. 4, the different size of the most stable clusters in terms of max metric are selected to participate in final ensemble. The accuracy of consensus partition extracted out of the selected clusters is presented in vertical axis. As it is obvious participating 20~30% of total clusters in the final ensemble is a very promising option. Also participation all clusters is not a good option.

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 a method is suggested which is called max method. The experiments show that the proposed framework commonly outperforms in comparison with the full ensemble; also participation all clusters in the final ensemble is not a good option; however it uses just 33% of primary clusters. Also the proposed max criterion does slightly better than NMI criterion generally. Because of the symmetry which is concealed in NMI criterion and also in NMI based stability, it yields to lower performance whenever symmetry is also appeared in the dataset. 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.: Cumulative Voting Consensus Method for Partitions with a Variable Number of Clusters. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, VOL. 30, NO. 1, 160-173, (2008)
2. Bhatia S.K., Deogun J.S.: Conceptual Clustering in Information Retrieval. *IEEE Trans. Systems, Man, and Cybernetics*, vol. 28, no. 3, 427-536, (1998)
3. Faceli K., Marcilio C.P., Souto D.: Multi-objective Clustering Ensemble. *Proceedings of the Sixth International Conference on Hybrid Intelligent Systems (HIS'06)* (2006)
4. Fern X.Z., Lin W.: Cluster Ensemble Selection. *SIAM International Conference on Data Mining (SDM08)* (2008)
5. Fred A., Jain A. K.: Data Clustering Using Evidence Accumulation. *Proc. of the 16th Intl. Conf. on Pattern Recognition, ICPR02, Quebec City*, 276 – 280, (2002)
6. Fred A., Jain A.K.: Combining Multiple Clusterings Using Evidence Accumulation. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 27(6):835–850, (2005)
7. Fred A., Jain A.K.: Learning Pairwise Similarity for Data Clustering. In *Proc. of the 18th Int. Conf. on Pattern Recognition (ICPR'06)*, (2006)
8. Fred A., Lourenco A.: Cluster Ensemble Methods: from Single Clusterings to Combined Solutions. *Studies in Computational Intelligence (SCI)*, 126, 3–30, (2008)
9. Frigui H., Krishnapuram R.: A Robust Competitive Clustering Algorithm with Applications in Computer Vision. *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 21, no. 5, 450-466, (1999)
10. Judd D., Mckinley P., Jain A.K.: Large-Scale Parallel Data Clustering. *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 2, 153-158, (1997)
11. Lange T., Braun M.L., Roth V., Buhmann J.M.: Stability-based model selection. In *Advances in Neural Information Processing Systems 15*. MIT Press, (2003)
12. Law M.H.C., Topchy A.P., Jain A.K.: Multiobjective data clustering. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, 424–430, (2004)
13. Newman C.B.D.J., Hettich S., Merz C.: UCI repository of machine learning databases. <http://www.ics.uci.edu/~mlearn/MLSummary.html>, (1998)
14. Roth V., Lange T., Braun M., Buhmann J.: A Resampling Approach to Cluster Validation. *Intl. Conf. on Computational Statistics, COMPSTAT*, (2002)
15. Strehl A., Ghosh J.: Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3(Dec):583–617, (2002)

Article received: 2011-08-20