



# Combining multiple clusterings using similarity graph

Selim Mimaroglu <sup>\*</sup>, Ertunc Erdil

Department of Computer Engineering, Bahcesehir University, Ciragan Caddesi, 34353 Besiktas, Istanbul, Turkey

## ARTICLE INFO

### Article history:

Received 13 May 2010

Received in revised form

25 August 2010

Accepted 8 September 2010

### Keywords:

Clustering

Combining clustering partitions

Cluster ensemble

Evidence accumulation

Robust clustering

Mutual information

## ABSTRACT

Multiple clusterings are produced for various needs and reasons in both distributed and local environments. Combining multiple clusterings into a final clustering which has better overall quality has gained importance recently. It is also expected that the final clustering is novel, robust, and scalable. In order to solve this challenging problem we introduce a new graph-based method. Our method uses the evidence accumulated in the previously obtained clusterings, and produces a very good quality final clustering. The number of clusters in the final clustering is obtained automatically; this is another important advantage of our technique. Experimental test results on real and synthetically generated data sets demonstrate the effectiveness of our new method.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

Clustering is the process of organizing objects into groups which have similar members; a distance metric is used for evaluating the similarity. Clustering, which has broad applications, is also known as unsupervised classification in the literature. In a cluster, hopefully, objects are similar to each other and they are dissimilar to other objects in other clusters.

Clustering has a long and rich history in a variety of scientific fields [1]: Taxonomists, social scientists, psychologists, biologists, statisticians, mathematicians, engineers, computer scientists, medical researchers, and others who collect and process real data have all contributed to clustering methodology. Some of the clustering algorithms partition a set of objects into groups, which are known as partitional clustering algorithms. *k*-means [2] is a well-known partitional clustering algorithm that divides a set of objects into *k* clusters, where *k* is a user specified parameter. Hierarchical clustering algorithms produce a nested sequence of clusters, with a single all-inclusive cluster at the top and single point clusters at the bottom. Agglomerative hierarchical algorithms [3] start with all the data points as a separate cluster. Each step of the algorithm involves merging two clusters that are the most similar. Yet another type of clustering is density based clustering; DBSCAN is a popular algorithm that can correctly cluster arbitrary shape data sets, when right parameters are

provided. An expert, by using his/her experience, can also produce a clustering. Therefore, for the same data set there may be multiple clusterings.

Finding natural groupings of a data set is a hard task as attested by hundreds of clustering algorithms [1] in the literature. Each clustering technique makes some assumptions about the underlying data set. If the assumption holds, good clusterings can be expected. But, more than often, assumptions about the data set do not hold, which in turn means bad clusterings are generated. It is very hard to select an appropriate clustering method, because the natural grouping is unknown. Many clustering algorithms effect the result by taking input parameters, e.g. *k*-means, *k*-medoids, DBSCAN, agnes, etc. It is beneficial to apply different clustering methods on the same data set, or the same method with varying input parameters or both. Combining multiple clusterings, which can be visualized in Fig. 1, has been studied in machine learning, pattern recognition, and data mining. Our main contribution is a new graph-based method for combining clusterings efficiently and effectively which produces very good quality clusters.

Some recent work on combining multiple clusterings can be found in [4–6]. CSPA, which is introduced in [7], is based on a co-association matrix, and METIS, which is a software package for partitioning unstructured graphs and hyper-graphs [8,9].

HGPA is introduced in [7] as well: Multiple clusterings construct a hyper-graph where each object is a vertex, and each cluster is an hyper-edge. Main idea is to have *k* unconnected components of the hyper-graph by using HMETIS [10,11]. Combining multiple clusterings problem is formulated as partitioning the hyper-graph by cutting a minimal number of hyper-edges. A set of hyper-edges are removed and *k* unconnected components are obtained, which provides the final clustering.

<sup>\*</sup> Corresponding author. Tel.: +90 212 381 05 55; fax: +90 212 381 05 50.  
E-mail addresses: [selim.mimaroglu@bahcesehir.edu.tr](mailto:selim.mimaroglu@bahcesehir.edu.tr) (S. Mimaroglu),  
[ertunc.erdil@stu.bahcesehir.edu.tr](mailto:ertunc.erdil@stu.bahcesehir.edu.tr) (E. Erdil).

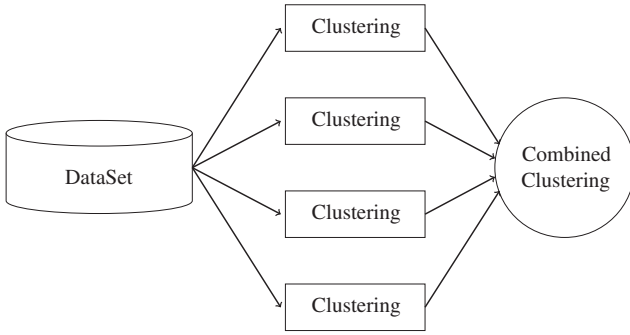


Fig. 1. Combining multiple clusterings.

In meta-clustering algorithm (MCLA) [7] each cluster is represented by a hyper-edge, like HGPA. MCLA is composed of the following steps: (1) constructing the meta-graph, (2) partitioning the meta-graph, and (3) computing cluster members.

Evidence accumulation (EAC) [12] accumulates the evidence in each cluster to form a co-association matrix,  $SM$ . Each entry in this matrix,  $SM_{ij}$ , is the number of times that objects  $i$  and  $j$  are assigned to the same clusters. The similarity matrix is provided as input to an agglomerative clustering algorithm.

The data set may be distributed at different sites, in this case a distributed clustering solution with a final merging of clusters is needed. Ref. [13] proposes two methods (BM and MM) for combining distributed clusters. Using cluster centers (prototypes) instead of clusters reduces computation and memory requirements. BM works on several clusterings each having  $n$  clusters. It groups the centroids according to their similarity and merges them to have a final clustering with  $n$  clusters. MM uses METIS, and it is more flexible: clusterings can have different number of clusters. Good results of both BM and MM are reported in [13].

In the literature, methods for partitioning unstructured graphs and hyper-graphs are also used for combining multiple clusterings. Graph partitioning methods compute a final clustering by cutting some of the edges (or hyper-edges), thus  $k$  unconnected components represent  $k$  final clusters. In our experimental evaluations, we used PMETIS, KMETIS, and HMETIS.

All the methods mentioned in this section require the number of final clusters in advance. BM and MM are prototype based clustering methods, therefore they apply to only globular shape data sets. For this reason, we will not consider BM and MM in the experimental evaluations. EAC requires a lot of computations, therefore it is slow on all the data sets. MCLA is slow on small data sets. Neither EAC, nor MCLA scales well. PMETIS, KMETIS, and HMETIS are very fast, but generally they provide low quality final clusterings. PMETIS, KMETIS, HMETIS, MCLA and EAC have trouble on some of the challenging data sets presented in Figs. 11 and 12.

The paper is structured as follows: In Section 2, the problem of combining multiple clusterings is introduced formally. Section 3 introduces our novel, similarity graph-based algorithm for combining multiple clusterings. Discussion of our main algorithm is presented in Section 4. Section 5 provides experimental evaluations. In the final section, we present our conclusions and future work.

## 2. Combining multiple clusterings

In this section we provide a formal definition to the combining multiple clusterings problem, which is also known as *cluster ensemble* in the literature, and refer to combining multiple clusterings into a new, final clustering. Combining multiple clusterings requires reusing pre-existing knowledge, employing

distributed data mining methods, and producing a final clustering having better overall quality. Some solutions for combining multiple clusterings are based on genetic algorithms such as [14,15] or simulated annealing [16]. A greedy optimization, information-theoretical method is presented in [7]. And, [17] presents a method based on ant colony optimization. All these methods operate on a consensus function.

Let  $D$  be a data set. A clustering (partition) of  $D$ ,  $\pi(D)$ , can be stated as follows:

$$\pi(D) = \{C_1, C_2, \dots, C_{|\pi(D)|}\}$$

where  $C_i$  is a cluster (block) of  $\pi(D)$ ,  $1 \leq i \leq |\pi(D)|$ , and

$$D = \bigcup_{i=1}^{|\pi(D)|} C_i$$

Note that we have a partial clustering (i.e. not complete) when  $\bigcup_{i=1}^{|\pi(D)|} C_i \subset D$ . Given a set of clusterings  $\Pi(D) = \{\pi_1(D), \pi_2(D), \dots, \pi_m(D)\}$ , the problem of combining multiple clusterings is defined as finding a new clustering  $\pi^*(D) = \{C_1^*, C_2^*, \dots, C_{|\pi^*(D)|}^*\}$  by using the information provided by  $\Pi(D)$ . This is achieved by using a consensus function  $cns(\Pi(D)) = \pi^*(D)$  such that

$$\forall i (\phi(\pi^*(D)) \geq \phi(\pi_i(D))), \quad 1 \leq i \leq |\Pi(D)| \quad (1)$$

where function  $\phi$  is a cluster validity measure. Exhaustively searching all the possible clusterings for finding the best clustering is not an option, since there are too many possible clusterings (see [7]).

Many consensus functions have been proposed in the literature. In [16,14], consensus functions based on median partition approach have been proposed. This approach searches differences between the clusterings by working on a coarser level. At another direction, in [7], consensus functions based on hyper-graphs have been proposed. In this technique, a hyper-edge represents a cluster, and a hyper-graph represents a clustering. A co-association measure, which is shown in (2), is used in [18].

$$coassoc(i, j) = votes_{ij} \quad (2)$$

where  $votes_{ij}$  is the number of times that objects  $i$  and  $j$  are assigned to the same clusters. This information is stored in a  $|D| \times |D|$  co-association matrix.

We represent each cluster by its characteristic bit vector which is as long as the size of the data set,  $|D|$ . Three clusterings on a data set are presented in Fig. 2. For example,  $C_{11}$  cluster has  $d_1$ ,  $d_3$ , and  $d_6$  objects as shown below.

$C_{11}$							
1	0	1	0	0	1	0	0

Fig. 3 shows the co-association matrix of Fig. 2. Note that, Fig. 3 represents the evidence accumulated by the pre-existing multiple clusterings [12,19,18,7].

Clustering	Cluster	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$d_7$	$d_8$
$\pi_1(D)$	$C_{11}$	1	0	1	0	0	1	0	0
	$C_{12}$	0	0	0	1	1	0	0	0
	$C_{13}$	0	1	0	0	0	0	1	1
$\pi_2(D)$	$C_{21}$	1	1	0	1	0	0	0	0
	$C_{22}$	0	0	0	0	0	0	0	1
	$C_{23}$	0	0	0	0	1	0	1	0
	$C_{24}$	0	0	1	0	0	1	0	0
$\pi_3(D)$	$C_{31}$	0	0	1	0	0	1	0	0
	$C_{32}$	1	1	0	1	0	0	0	1
	$C_{33}$	0	0	0	0	1	0	1	0

Fig. 2. Binary representation of multiple clusterings,  $\Pi$ .

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$d_7$	$d_8$
$d_1$	3	2	1	2	0	1	0	1
$d_2$	2	3	0	2	0	0	1	2
$d_3$	1	0	3	0	0	3	0	0
$d_4$	2	2	0	3	1	0	0	1
$d_5$	0	0	0	1	3	0	2	0
$d_6$	1	0	3	0	0	3	0	0
$d_7$	0	1	0	0	2	0	3	1
$d_8$	1	2	0	1	0	0	1	3

Fig. 3. Co-association matrix,  $SM$ , of Fig. 2.

### 3. Combining multiple clusterings via similarity graph

In this section we present our novel algorithm, which works on a *similarity graph*. The similarity graph is an undirected and weighted graph that represents a co-association (similarity) matrix. In a similarity graph,  $SG=(D,E)$ , each edge  $(d_i, d_j)$  has a label which corresponds to the entry  $SM_{ij}$  in the co-association matrix. To simplify a similarity graph we can omit edges with 0 weight, and self loops (i.e. all the edges  $(d_i, d_i)$ ), because this information is not necessary for constructing a final clustering. Also, edge labels of value 1, i.e.  $SM_{ij}=1$ , are not written on the similarity graph for simplicity.

**Definition 3.1.** The *degree of freedom* of a vertex  $d_i$  is

$$df(d_i) = |\{d_j | (d_j, d_i) \in E\}|$$

For a data set  $D$ , and a family of clusterings  $\Pi(D)$ , let  $SM$  be the corresponding co-association matrix. Edges are labeled by the function weight defined by

$$\text{weight}(d_i, d_j) = SM_{ij}$$

where  $SM_{ij}$  is the entry at row  $i$  and column  $j$  of  $SM$ .

The *sum of weights* of edges incident to a vertex  $d_i$  is the

$$sw(d_i) = \sum_{j=1, j \neq i}^{|D|} \text{weight}(d_j, d_i)$$

**Lemma 3.2.** For any vertex  $d_i$  of a similarity graph we have

$$sw(d_i) \geq df(d_i)$$

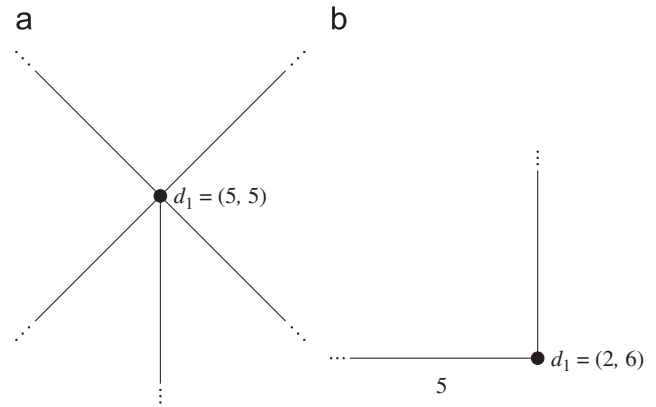
**Proof.** The inequality is immediate from the definitions of  $df$  and  $sw$ .  $\square$

In a similarity graph each vertex  $d_i$ , is labeled by a pair  $(df(d_i), sw(d_i))$ , and each edge  $(d_i, d_j)$  is labeled by  $\text{weight}(d_i, d_j)$ . Two examples demonstrating the degree of freedom and sum of weights are shown in Fig. 4. High values of  $df(d_i)$  mean that  $d_i$  is connected to many other vertices. Similarly, large values of sum of weights (i.e. large values of  $sw(d_i)$ ) indicate high similarity with several other vertices.

**Definition 3.3.** The *attachment* of a vertex is given by

$$\text{attachment}(d_i) = \frac{sw(d_i)}{df(d_i)}$$

We use attachment in our algorithm to initiate new clusters; an object having the highest attachment is selected as initial object

Fig. 4.  $df$  and  $sw$ : (a)  $df(d_1) = sw(d_1) = 5$  and (b)  $df(d_1) = 2$ ,  $sw(d_1) = 6$ .

(pivot) for the reasons that will be explained in Section 4.1. There may be isolated nodes having, 0 degree of freedom and 0 sum of weights. By convention, attachment value of such a vertex is considered as 0.

Combining Multiple Clusterings via Similarity Graph (COMUSA) is introduced in Algorithm 1.

**Algorithm 1.** Combining Multiple Clusterings via Similarity Graph COMUSA

**Input:**  $D$ : Data Set,  $\Pi(D)$ : Multiple Clusterings

**Output:**  $\pi^*(D)$ : Final Clustering

```

1 Initialize an empty queue  $Q$ ;
2  $clusterId = 1$ ;
3 Construct similarity graph  $SG = (D, E)$  using  $\Pi(D)$ , and  $D$ ;
4 Sort  $D$  in decreasing order with respect to attachment;
5 while there are unmarked objects do
6   Add unmarked object,  $d_i$ , with highest attachment( $d_i$ )
   to  $Q$ ;
7   while  $Q$  is not empty do
8     // pivot object
9      $v =$  remove first element from  $Q$ ;
10    Add  $v$  to cluster  $clusterId$ ;
11    Mark  $v$ ;
12    foreach  $(w, v) \in E$  do
13      if  $w$  is marked then
14        continue;
15      else
16         $strWeight = \text{weight}(w, v)$ ;
17         $isMax = \text{true}$ ;
18        foreach  $(z, w) \in E$  do
19          // maximum constraint
20          if  $strWeight \neq \text{weight}(z, w)$  then
21             $isMax = \text{false}$ ;
22            break;
23          if  $isMax$  then
24            Add  $w$  to  $Q$ ;
25     $clusterId++$ ;

```

For a data set and a family of clusterings of this data set (in Fig. 2) the corresponding co-association matrix and the similarity graph are shown in Figs. 3 and 5, respectively. COMUSA is explained on the similarity graph shown in Fig. 5. In Table 1, attachment values of vertices are shown.  $d_3$  and  $d_6$  have the highest attachment values; we randomly pick  $d_3$  as a pivot to start a new cluster.

Next step is to discover the objects that will be in the same cluster as  $d_3$ . Starting from nearest neighbors, we have to evaluate

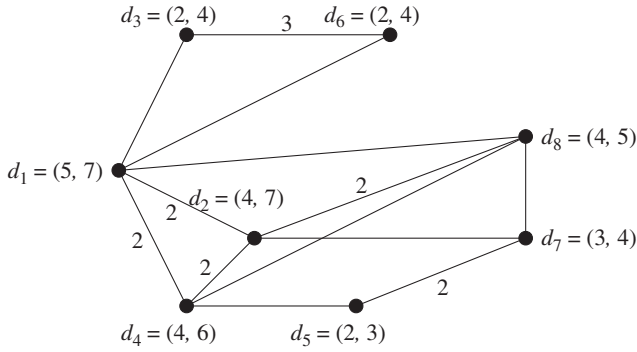


Fig. 5. Similarity graph of Fig. 3.

Table 1

Attachment values in decreasing order.

Vertex	Attachment (vertex)
$d_3$	2.00
$d_6$	2.00
$d_2$	1.75
$d_4$	1.50
$d_5$	1.50
$d_1$	1.40
$d_7$	1.33
$d_8$	1.25

the edge weights and include objects into first cluster if they are strongly linked to  $d_3$ . Immediate neighbors of  $d_3$  are  $d_6$  and  $d_1$ .  $d_6$  is included in the cluster since  $\text{weight}(d_3, d_6) \geq \text{weight}(d_6, d_1)$ .  $d_1$  is not included in the same cluster, because  $\text{weight}(d_3, d_1) \neq \text{weight}(d_1, d_4)$  (similarly  $\text{weight}(d_3, d_1) \neq \text{weight}(d_1, d_2)$ ). Our first cluster has two objects  $C_{1.1} = \{d_3, d_6\}$ .

Since there are some objects that are not included in any cluster (not marked), the algorithm keeps running. Next,  $d_2$  is selected as a pivot, because among all the unmarked objects it has the highest attachment value. The algorithm considers including  $d_1$ ,  $d_4$ ,  $d_7$  and  $d_8$ . The object  $d_1$  is included in the cluster since among all the edges passing through  $d_1$ ,  $\text{weight}(d_2, d_1)$  has the maximum value. Similarly,  $d_4$  is added to the same cluster since  $\text{weight}(d_2, d_4) \geq \text{weight}(d_4, d_8)$ ,  $\text{weight}(d_2, d_4) \geq \text{weight}(d_4, d_5)$ , and  $\text{weight}(d_2, d_4) \geq \text{weight}(d_4, d_1)$ . For the same reason  $d_8$  is also added to the second cluster, but not  $d_7$ . Next, the algorithm considers the neighbors of newly added objects:  $d_1$  does not have any unmarked neighbors,  $d_4$  becomes an acting pivot, and the algorithm considers adding  $d_5$ . This is not possible because  $\text{weight}(d_4, d_5) \neq \text{weight}(d_5, d_7)$ . The object  $d_4$  does not have any other neighbors, so  $d_8$  becomes acting pivot. For  $d_8$  the algorithm considers adding  $d_7$  to the current cluster, but it is not possible since  $\text{weight}(d_8, d_7) \neq \text{weight}(d_7, d_5)$ . There are no acting pivots, therefore no further extension is possible: we have  $C_{2.2} = \{d_2, d_1, d_4, d_8\}$ .

There are only two unmarked objects left which are  $d_5$  and  $d_7$ . The object  $d_5$  has the highest attachment value so it is chosen as a pivot object. Its only unmarked neighbor is  $d_7$ .  $\text{weight}(d_5, d_7) \geq \text{weight}(d_7, d_2)$  and  $\text{weight}(d_5, d_7) \geq \text{weight}(d_7, d_8)$ , therefore  $d_7$  and  $d_5$  are clustered together, so  $C_{3.3} = \{d_5, d_7\}$ . All the objects are marked, COMUSA halts. Final clustering having three clusters are shown in Fig. 6.

COMUSA starts a new cluster with an object having the highest attachment value, then extends the cluster at hand as much as possible. In a similarity graph neighbors of a pivot are checked with respect to their similarity to the pivot. Then, each neighbor is considered as an acting pivot. In a final clustering the number of clusters depends on the data set: COMUSA detects this number automatically, which is a big advantage.

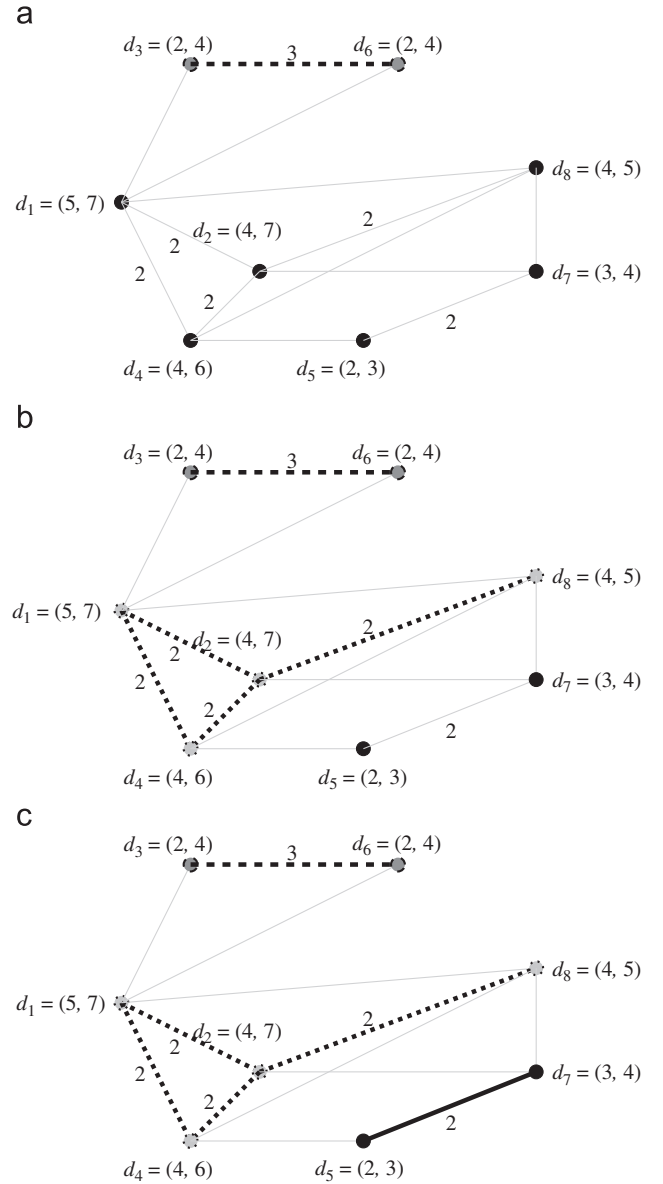


Fig. 6. Generating final clustering using COMUSA on Fig. 5: (a) First cluster in dashed pattern; (b) Second cluster in dotted pattern and (c) Third cluster in straight bold line.

#### 4. Discussion of COMUSA

In this section, we explain the important features of COMUSA.

##### 4.1. Selection of pivot objects

Selecting a pivot (seed) object for initiating a new cluster is essential in COMUSA. High attachment values indicate high sum of weights and low degree of freedom. Therefore, an object having a high attachment value is *strongly connected* somewhere. An unmarked object having the highest attachment value is selected as pivot for starting a new cluster.

##### 4.2. Expansion of a cluster

Initially each cluster is a singleton. Pivot object expands the cluster by considering all the immediate neighbors. A neighbor is included in a pivot's cluster if it is most similar to the pivot. Once

a neighbor is included, it is marked and then it acts like a pivot by considering its immediate neighbors for further expansion.

#### 4.3. Termination condition

Each cluster is expanded by its neighbors as explained previously. Extension of a cluster comes to an end if pivots of a cluster cannot add any other objects into the cluster. In a data set, if there are unmarked objects, COMUSA starts a new cluster by choosing a new pivot. COMUSA halts when all the objects are marked.

#### 4.4. COMUSA works well

Arbitrary shape clusters can be found by our algorithm, we do not make any assumptions about the input data set. COMUSA works very well because pivot objects are good starting points, and an object is included into a cluster if the object is most similar to a pivot in that cluster. Experimental results show that in a short amount of time COMUSA creates very good quality clusters on real and synthetic data sets, even on very challenging ones (see Figs. 10a, 11, 12), therefore remedies the weaknesses of the related work.

#### 4.5. Relaxation

COMUSA finds the natural number of clusters efficiently and automatically. This feature, presently missing in many clustering algorithms, is very important, since determining the correct number of clusters is very hard. Experimental evaluations demonstrate that number of clusters found by COMUSA are either identical to the natural number of clusters in a data set or very close to this number.

In some cases it is desirable to have larger clusters. Larger clusters can be obtained with COMUSA by relaxing the *maximum* constraint. Maximum edge weight constraint can be observed in lines 15–20 of Algorithm 1. Instead of requiring the maximum edge weight, we can relax this constraint by a user input *relaxation* percentage.

COMUSA with relaxation is explained on a (partial) similarity graph that is shown in Fig. 7. Highest attachment value node,  $d_1$ , is selected as pivot. The algorithm tries to extend the cluster with  $d_2$ , but this is not possible since  $\text{weight}(d_2, d_3) > \text{weight}(d_1, d_2)$ . With a relaxation of 25%,  $d_2$  can be added to the cluster since  $\text{weight}(d_1, d_2) + \text{weight}(d_1, d_2) \cdot 25\% \geq \text{weight}(d_2, d_3)$ . Then,  $d_3$  is added to the same cluster. Notice that COMUSA with a positive relaxation value produces larger, and fewer clusters.

The next two examples show that COMUSA is robust, and intuitive with respect to the similarity graph.

**Example 4.1.** Let us consider the similarity graph shown in Fig. 8a. Notice that all the edge labels are 2, and attachment values for all the vertices are constant. Each vertex is qualified to be a pivot, and no matter what vertex is selected as the pivot

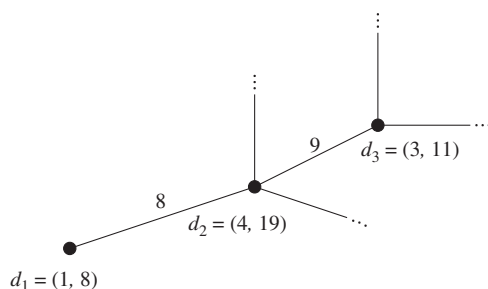


Fig. 7. A partial similarity graph.

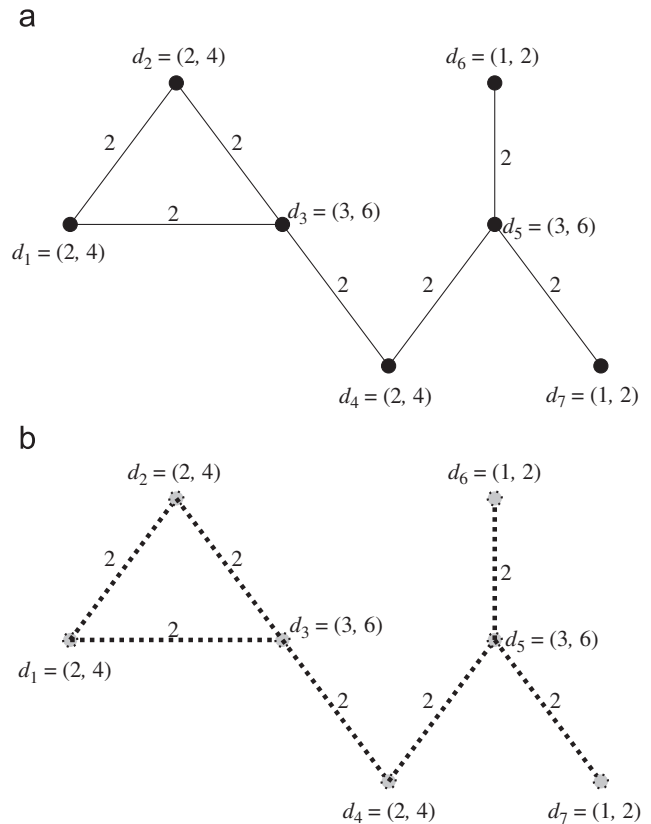


Fig. 8. COMUSA on a data set: (a) Similarity graph of a data set and (b) Final clustering of (a) in dashed pattern.

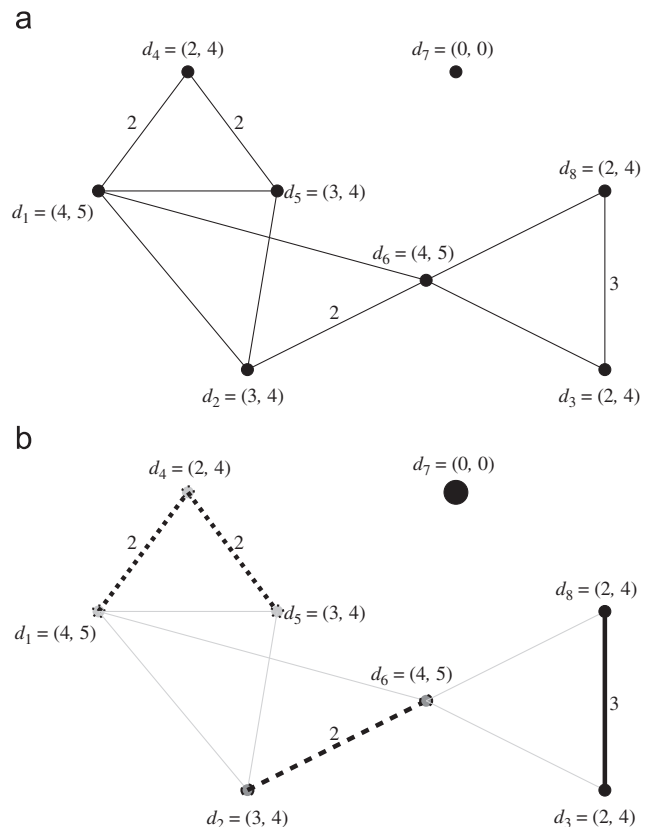


Fig. 9. COMUSA on another data set: (a) Similarity graph of another data set and (b) Final clustering of (a) having 4 clusters. Three clusters are shown with distinct patterns,  $d_7$  is a singleton cluster.



we end up with one big cluster having all the vertices:  $C_{*1} = \{d_1, d_2, d_3, d_4, d_5, d_6, d_7, d_8\}$  as shown in Fig. 8b.

**Example 4.2.** Running COMUSA on the similarity graph shown in Fig. 9a generates 4 clusters:  $C_{*1} = \{d_3, d_8\}$ ,  $C_{*2} = \{d_1, d_4, d_5\}$ ,  $C_{*3} = \{d_2, d_6\}$ , and  $C_{*4} = \{d_7\}$ . This result is very intuitive too, objects having high values of similarity are grouped in the same clusters. Also, note that isolated object  $d_7$  is left by itself in a cluster. Final clustering produced by COMUSA is shown in Fig. 9b.

## 5. Experimental evaluations

In this section we present some objective cluster quality measures, experimental data sets, and test results.

### 5.1. Evaluating validity of the final clustering

Several objective measures have been proposed to evaluate goodness of the final clustering such as inter-cluster similarity, intra-cluster similarity [20], rand index [21], adjusted rand index (ARI) [22], normalized mutual information [7], jaccard index [23], and silhouette coefficients [24]. Cluster validity functions that are used for evaluating COMUSA are explained below.

**Table 2**  
Contingency table.

Class/cluster	$v_1$	$v_2$	...	$v_p$	Sums
$u_1$	$n_{11}$	$n_{12}$	...	$n_{1p}$	$n_{1.}$
$u_2$	$n_{21}$	$n_{22}$	...	$n_{2p}$	$n_{2.}$
$\vdots$	$\vdots$	$\vdots$		$\vdots$	
$u_r$	$n_{r1}$	$n_{r2}$	...	$n_{rp}$	$n_{r.}$
Sums	$n_{.1}$	$n_{.2}$		$n_{.p}$	$n_{..} = n$

**Table 3**  
Cluster ensembles of test data sets.

INPUT	Number of clusterings	Number of clusters	Ensemble generation methods
1-Spiral, hand clustered	2	3	Manually
1-Spiral, k-means clustered	4	[2, 3]	k-means
2-Spiral	2	[2, 3]	Manually
2-Half rings	3	[2, 5]	k-means
2-Curve	2	[2, 8]	k-means, randomly, manually
2D2K	3	[2, 3]	k-means
8D5K	3	[3, 5]	k-means
Iris	3	[2, 3]	k-means
Glass	4	[6, 8]	k-means, randomly, manually
Breast Cancer	5	[2, 5]	k-means, randomly, manually
Image Segmentation	10	7	Randomly
Syn5K	10	5	Randomly

#### 5.1.1. Intra and inter cluster similarities

Intra-cluster similarity measures the inner similarity of a cluster, where large values are preferred. For a clustering,  $\pi(D) = \{C_1, C_2, \dots, C_{|\pi(D)|}\}$ , intra-cluster similarity is defined as follows:

$$ICS(\pi(D)) = \sum_{i=1}^{|\pi(D)|} \frac{1}{|C_i|^2} \sum_{d, d' \in C_i} sim(d, d') \quad (3)$$

In (3)  $sim(d, d')$  is the similarity of the objects  $d$  and  $d'$ . When working on multiple clusterings, evidence accumulated in  $\Pi(D)$  is used as a similarity measure as mentioned earlier. Intra-cluster similarity of a final clustering  $\pi^*(D)$  with respect to multiple clusterings  $\Pi(D)$  is shown in the following equation [20]:

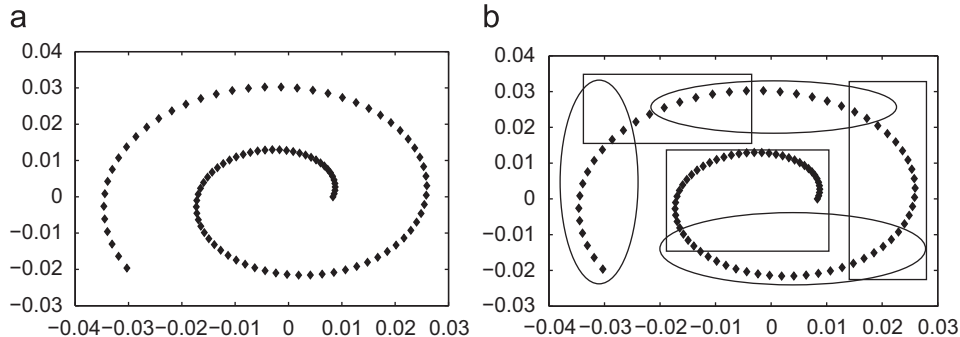
$$ICS_{\Pi}(\pi^*(D)) = \sum_{k=1}^{|\pi^*(D)|} \frac{1}{|C_{*k}|^2} \sum_{i=1}^{|\Pi|} \sum_{j=1}^{|\pi_i|} \binom{|C_{*k} \cap C_{ij}|}{2} \quad (4)$$

Inter-cluster similarity of a cluster  $\pi(D)$  is defined as follows:

$$ECS(\pi(D)) = \sum_{i=1}^{|\pi(D)|} \sum_{j=i+1}^{|\pi(D)|} \frac{1}{|C_i||C_j|} \sum_{d \in C_i, d' \in C_j} sim(d, d') \quad (5)$$

**Table 4**  
COMUSA on 1-spiral data set.

INPUT	Relaxation (%)	ARI
1-Spiral, hand clustered	0	1.0
1-Spiral, k-means clustered	34	1.0



**Fig. 10.** 1-Spiral data set and a clustering.

Low values of (5) indicate that  $\pi(D)$  has isolated clusters, which is preferred. Inter-cluster similarity of a final clustering  $\pi^*(D)$  with respect to multiple clusterings  $\Pi(D)$  is shown in Formula (6) [20].

$$ECS_{\Pi}(\pi^*(D)) = \sum_{k=1}^{|\pi^*(D)|} \sum_{l=k+1}^{|\pi^*(D)|} \frac{1}{|C_{*,k}| |C_{*,l}|} \sum_{i=1}^{|\Pi|} \sum_{j=1}^{|\pi_i|} \left( \frac{|(C_{*,k} \vee C_{*,l}) \wedge C_{ij}|}{2} - \binom{|C_{*,k} \wedge C_{ij}|}{2} - \binom{|C_{*,l} \wedge C_{ij}|}{2} \right) \quad (6)$$

Inter-cluster similarity and intra-cluster similarity can be combined to form a clustering validity function as shown in (7).

$$\phi(\pi^*(D)) = k_1 \cdot ICS(\pi^*(D)) + k_2 \cdot ECS(\pi^*(D)) \quad (7)$$

where  $k_1 > 0$  and  $k_2 < 0$ . In our tests, we use (7) as an unsupervised cluster evaluation technique and we refer to it as ICS+ECS.

### 5.1.2. Adjusted rand index (ARI)

We use adjusted rand index (ARI) in order to measure the extent to which the clustering structure discovered by COMUSA matches some external criteria, i.e. class labels. Given a data set  $D = \{d_1, \dots, d_n\}$ , suppose  $U = \{u_1, \dots, u_r\}$  represents classes, and  $V = \{v_1, \dots, v_p\}$  represents a clusterings of the  $D$ .

$$\bigcup_{i=1}^r u_i = D = \bigcup_{j=1}^p v_j$$

and  $u_i \cap u_j = \emptyset$  for  $1 \leq i, j \leq r$  and  $i \neq j$ . Also,  $v_i \cap v_j = \emptyset$  for  $1 \leq i, j \leq p$  and  $i \neq j$ .

In Table 2,  $n_{ij} = |u_i \cap v_j|$ ,  $n_i = \sum_{j=1}^p n_{ij}$ , and  $n_j = \sum_{i=1}^r n_{ij}$ . ARI can be formulated as follows:

$$\frac{\sum_{i,j} \binom{n_{ij}}{2} - \left( \sum_i \binom{n_i}{2} \sum_j \binom{n_j}{2} \right) / \binom{n}{2}}{\frac{1}{2} \left( \sum_i \binom{n_i}{2} + \sum_j \binom{n_j}{2} \right) - \left( \sum_i \binom{n_i}{2} \sum_j \binom{n_j}{2} \right) / \binom{n}{2}}$$

ARI takes maximum value at 1, which indicates perfect match to the external criteria.

### 5.2. Generating cluster ensembles

In our experiments we use three different approaches for generating cluster ensembles: manually constructing clusters, randomly constructing clusters or randomly injecting error into the original clusters, and using  $k$ -means algorithm with varying  $k$ -values. Note that the diversity and quality of a cluster ensemble impacts the quality of the final clustering. We choose different approaches for generating cluster ensembles because we expect

them to produce a diverse set of clusterings with different properties and qualities. For all the experimental data sets, Table 3 provides the number clusterings (partitions), number of clusters at each clustering, and the ensemble generation methods. For example, Breast Cancer data set has 5 clusterings, each clustering with 2–5 clusters, and each clustering is generated by  $k$ -means, manually or at random. Experimental data sets, and the actual cluster ensembles used for testing along with the Java implementation of COMUSA is available at <http://www.cs.umb.edu/~smimarog/implementations/comusa>.

### 5.3. Test results

We have conducted experiments on a computer having 2.8GHz processor with 4GB of main memory, running on Linux kernel 2.6. Our choice of implementation language is Java, which provides built-in support for bit vectors, and operations on bit vectors. COMUSA, MCLA, and EAC are all implemented in Java and are tested with Java Development Kit 1.6.0.16. We obtained PMETIS, KMETIS, and HMETIS from the corresponding authors. PMETIS and KMETIS belong to the METIS package and are implemented in C language. HMETIS is also implemented in C.

The 1-spiral data set is synthetically generated and it contains 100 objects. The 2-spiral, 2-half rings, 2-curve data sets are also synthetically generated and contain 200, 118, and 192 objects respectively. Although these are small data sets having 2 dimensions only, finding the natural clusters of these data sets is challenging for many clustering algorithms and combining multiple clusterings methods.

The 2D2K and the 8D5K data sets are obtained from [7]: 2D2K is synthetically generated and contains 500 points each of two 2-dimensional Gaussian clusters with means (0.227, 0.077) and

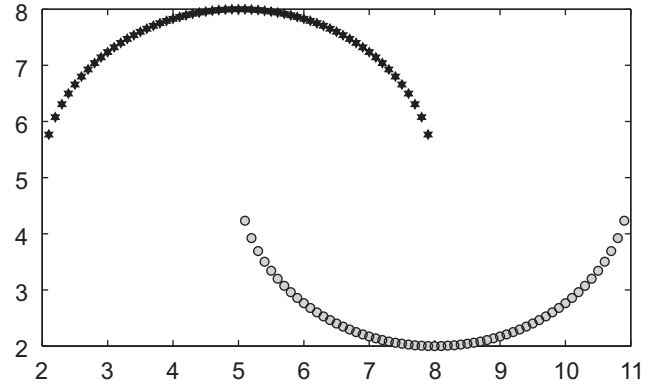


Fig. 12. 2-Half rings data set.

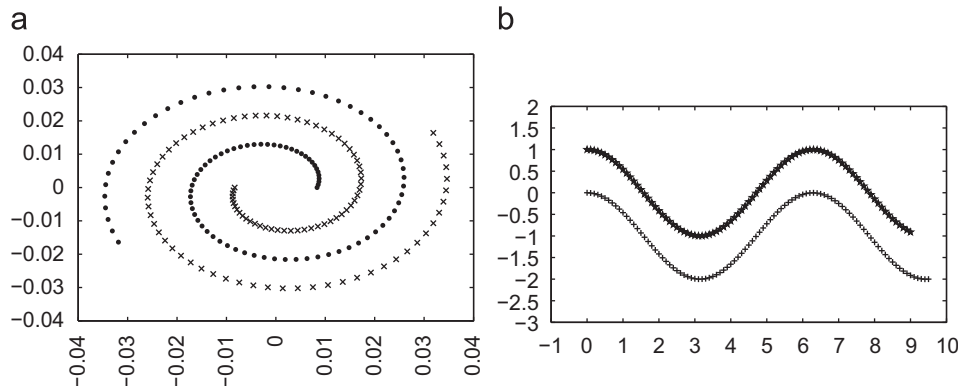


Fig. 11. 2-Spiral and 2-curve data sets.

(0.095, 0.323) and diagonal covariance matrices with 0.1 for all diagonal elements. The 8D5K contains 1000 points from five multivariate Gaussian distributions (200 points each) in 8-dimensional space. The clusters all have the same variance (0.1), but different means. Means were drawn from a uniform distribution within the unit hypercube. Syn5K data set is an artificially generated data set containing 5000 objects and 5 classes.

Iris, Glass, Breast Cancer, and Image Segmentation are real data sets that are obtained from University of California Irvine Machine Learning Repository [25]. Iris is a multivariate, real data set having 4 dimensions, 150 objects and 3 classes. Glass Identification data set is a multivariate, real data set having 10 dimensions, 214 objects, 6 classes. Breast Cancer data set is a multivariate, real data set, having 9 attributes, 286 objects and 2 classes. Image Segmentation is also a multivariate data set, with 19 real type attributes, 2310 objects and 7 classes.

On the data set shown in Fig. 10a, we created two partial clusterings. These clusters are shown in Fig. 10b, where rectangular objects form a clustering, and elliptical objects form another clustering. Notice that although both of the clusterings are partial, COMUSA successfully identifies the 1-spiral data set. On the same data set,  $k$ -means with  $k=2$ , and  $k=3$  are performed twice, for each value, to obtain four different clusterings. These multiple clusterings are provided as input to COMUSA with 34% relaxation, and again COMUSA successfully discovered the natural clusters as shown in Table 4.

The 2-spiral, 2-curve, and 2-half rings data sets are shown in Figs. 11 and 12. Hand created clusters, as well as  $k$ -means clusters are obtained on these data sets. Partitions generated by  $k$ -means on the 2-half rings data set are shown in Fig. 13. COMUSA, PMETIS, KMETIS, HMETIS, MCLA, and EAC results are compared for cluster validity. As shown in Table 5, COMUSA produces perfect output on all the data sets. For these data sets PMETIS, KMETIS, HMETIS, MCLA and EAC are requested to produce 2 clusters for fairness.

On 2D2K and 8D5K data sets COMUSA results are compared to PMETIS, KMETIS, HMETIS, MCLA, and EAC results with respect to ECS+ICS validity measure as shown in Table 6. Clearly, COMUSA provides final clusterings having superior quality. Number of clusters is provided as input to PMETIS, KMETIS, HMETIS, MCLA, and EAC, which can be seen in the same table.

**Table 5**

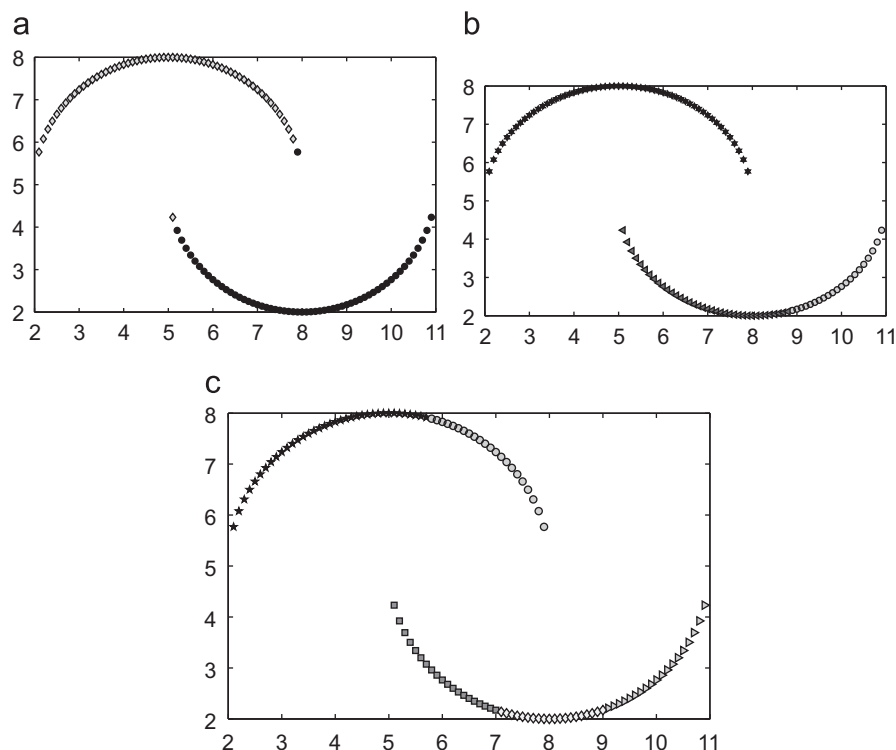
Cluster validity results on 2-spiral, 2-half ring and 2-curves data sets.

INPUT	COMUSA		PMETIS	KMETIS	HMETIS	MCLA	EAC
	Relaxation	ARI	ARI	ARI	ARI	ARI	ARI
2-Spiral	%34	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	−0.005	<b>1.0</b>	0.0
2-Half rings	%34	<b>1.0</b>	0.966	0.966	−0.008	<b>1.0</b>	<b>1.0</b>
2-Curve	%50	<b>1.0</b>	0.057	0.057	−0.004	0.086	<b>1.0</b>

**Table 6**

Cluster validity results on 2D2K and 8D5K.

INPUT	COMUSA					
	Relaxation (%)					ECS+ICS
2D2K	0					<b>54.27</b>
8D5K	0					<b>238.68</b>
INPUT	Number of clusters	PMETIS ECS+ICS	KMETIS ECS+ICS	HMETIS ECS+ICS	MCLA ECS+ICS	EAC ECS+ICS
2D2K	2	26.80	26.93	12.22	26.93	26.93
	3	33.06	32.39	36.40	41.07	39.38
8D5K	5	214.18	214.18	57.98	219.87	192.31
	6	211.64	213.42	157.59	219.87	<b>238.68</b>



**Fig. 13.** Partitions of 2-half rings data set generated with  $k$ -means: (a)  $k=2$ ; (b)  $k=3$  and (c)  $k=5$ .



We also conducted experimental evaluations on real data sets: Iris, Glass, Breast Cancer, Image Segmentation. We compare the generated clusters with the real class labels. Table 7 shows that COMUSA produces very good quality clusterings.

Table 8 shows the average execution time results of COMUSA, PMETIS, KMETIS, HMETIS, MCLA, and EAC for the data sets and inputs in Tables 5–7. COMUSA is much faster than both MCLA and EAC on all the data sets except Image Segmentation and Syn5K. For most of the data sets COMUSA is faster than both PMETIS and

**Table 7**  
Cluster validity results on Iris, Glass, Breast Cancer, Image Segmentation, and Syn5K data sets.

INPUT	COMUSA	
	Relaxation (%)	ARI
Iris	0	0.676
	50	0.698
Glass	0	0.308
	25	0.403
	34	<b>0.964</b>
Breast Cancer	0	0.156
	17	0.301
	20	0.604
	25	<b>1.0</b>
Image Segmentation	0	0.350
	12	0.932
	13	0.831
Syn5K	0	0.301
	15	<b>1.0</b>
	25	0.999

INPUT	Number of clusters	PMETIS ARI	KMETIS ARI	HMETIS ARI	MCLA ARI	EAC ARI
Iris	3	0.691	0.688	0.096	<b>0.711</b>	<b>0.711</b>
	4	0.412	0.368	0.036	0.662	0.690
Glass	6	0.4740	0.4698	0.1568	0.9633	0.8041
	8	0.3835	0.4151	0.0063	0.6439	0.7862
	22	0.2026	0.1971	0.0654	0.2898	0.6404
Breast Cancer	24	0.1892	0.1705	0.0615	0.2847	0.4292
	2	0.3597	0.3942	0.0024	0.5778	0.8322
	6	0.2174	0.2039	0.1311	0.2967	0.8981
	11	0.1189	0.1122	0.0684	0.3268	0.5463
Image Segmentation	16	0.0836	0.0828	0.0188	0.4547	0.5316
	7	0.987	<b>0.988</b>	0.535	0.985	0.840
Syn5K	9	0.633	0.639	0.352	0.938	0.837
	11	0.574	0.539	0.354	0.969	0.838
	3	0.488	0.480	0.317	0.611	0.482
	5	<b>1.0</b>	<b>1.0</b>	0.426	<b>1.0</b>	<b>1.0</b>
	7	0.578	0.556	0.607	0.953	0.999

**Table 8**  
Execution time results (ms).

INPUT	COMUSA	PMETIS	KMETIS	HMETIS	MCLA	EAC
2-Spiral	1.1	4.0	4.0	1.0	201.0	148.0
2-Half rings	2.3	2.0	2.0	1.0	196.0	114.0
2-Curve	1.1	2.0	3.0	3.0	203.0	150.0
2D2K	25.1	48.5	46.0	2.0	197.5	1203.5
8D5K	21.7	43.0	35.5	2.5	195.5	1207
Iris	3.6	3.0	2.5	2.5	199.5	117.5
Glass	1.3	5.5	7.3	34.3	198.5	149.5
Breast Cancer	2.1	9.5	12	12	198.8	601
Image Segmentation	7083.6	251.0	206.0	2	922.3	15466.3
Syn5K	71769.0	1371.0	1159.0	6	372.3	30774.6

**Table 9**  
Number of clusters.

INPUT	Relaxation (%)	Number of clusters	
		COMUSA	Original
1-Spiral, hand clustered	0	1	1
1-Spiral, k-means clustered	34	1	1
2-Spiral	20	2	2
2-Half rings	34	2	2
2-Curve	20	2	2
2D2K	50	2	2
8D5K	0	6	5
Iris	0	4	3
Glass	50	3	
	0	24	6
	25	22	
	34	6	
Breast Cancer	0	16	2
	17	11	
	20	6	
	25	2	
Image Segmentation	0	256	7
	12	32	
	13	16	
	20	9	
Syn5K	0	276	5
	15	5	
	25	5	

KMETIS, and for a few data sets COMUSA is comparable to PMETIS and KMETIS. COMUSA is faster than HMETIS on three data sets. Note that COMUSA is implemented in Java, which is known to be slower than C language implementations. Changing relaxation values of COMUSA does not effect the execution time considerably: COMUSA iterates over all the edges of the similarity graph regardless of the relaxation input.

COMUSA is good at finding correct number of clusters, depending on the relaxation rate. Number of clusters generated by COMUSA are compared with the original number of clusters for real and synthetic data sets in Table 9. For all the data sets COMUSA is able to find correct number of clusters or come close to this number.

COMUSA's feature of finding the number of clusters automatically can be explained as follows. A vertex  $v_i$  is added to the cluster of  $v_p$  if it is more closely attached to  $v_p$  than to any other vertex. COMUSA applies this principle repeatedly and generates connected components. Similarity graph is the backbone structure for discovering connected components, because each connected components resides on the similarity graph. When COMUSA comes to a stop, each connected component forms a cluster. Since the similarity graph is computed from the input clustering ensembles, ensemble generation step impacts both the quality of the final clustering and the number of clusters in the final clustering (see Examples 4.1 and 4.2). By applying several different cluster ensemble generation methods as explained in Section 5.2, we demonstrated that COMUSA produces good quality clusters, with correct number of clusters on a diverse collection of cluster ensembles and data sets.

All the test data sets and COMUSA implementation are available at <http://www.cs.umb.edu/~smimarog/implementations/comusa>.

## 6. Conclusions and future work

In this paper, we presented a novel algorithm for combining multiple clusterings. COMUSA takes multiple clusterings as input and creates a similarity graph by using the evidence accumulated

from the clusterings. By using the similarity graph, COMUSA computes pivot objects for starting a new cluster, and expands the clusters as much as possible. Number of clusters is automatically found by our algorithm with respect to relaxation rate. COMUSA is partitional, exclusive, and complete. Extensive experimental evaluations on many real and artificial data sets demonstrate that COMUSA: (1) finds arbitrary shape clusters, (2) is not affected by the cluster size, (3) is not affected by noise and outliers, (4) is not affected by the sparseness of the data set, (5) is order independent, (6) is deterministic, and (7) scales well.

As future work, we are planning to construct the similarity graph at cluster level, and improve COMUSA accordingly. As a result of this modification, the final clustering validity may diminish, and improving the validity may be challenging.

## References

- [1] A.K. Jain, Data clustering: 50 years beyond k-means, in: W. Daelemans, B. Goethals, K. Morik (Eds.), *ECML/PKDD (1)*, Lecture Notes in Computer Science, vol. 5211, Springer2008, pp. 3–4.
- [2] J.B. MacQueen, Some methods for classification and analysis of multivariate observations, in: *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.
- [3] A.K. Jain, R.C. Dubes, *Algorithms for Clustering Data*, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1988.
- [4] H.G. Ayad, M.S. Kamel, On voting-based consensus of cluster ensembles, *Pattern Recognition* 43 (5) (2010) 1943–1953.
- [5] X. Wang, C. Yang, J. Zhou, Clustering aggregation by probability accumulation, *Pattern Recognition* 42 (5) (2009) 668–675.
- [6] S. Vega-Pons, J. Correa-Morris, J. Ruiz-Shulcloper, Weighted partition consensus via kernels, *Pattern Recognition* 43 (8) (2010) 2712–2724.
- [7] A. Strehl, J. Ghosh, Cluster ensembles—a knowledge reuse framework for combining multiple partitions, *Journal of Machine Learning Research* 3 (2003) 583–617.
- [8] G. Karypis, V. Kumar, Multilevel algorithms for multi-constraint graph partitioning, in: *Supercomputing '98: Proceedings of the 1998 ACM/IEEE Conference on Supercomputing (CDROM)*, IEEE Computer Society, Washington, DC, USA1998, pp. 1–13.
- [9] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, *SIAM Journal on Scientific Computing* 20 (1) (1998) 359–392.
- [10] G. Karypis, R. Aggarwal, V. Kumar, S. Shekhar, Multilevel hypergraph partitioning: application in VLSI domain, in: *DAC '97: Proceedings of the 34th Annual Design Automation Conference*, ACM, New York, NY, USA1997, pp. 526–529.
- [11] G. Karypis, R. Aggarwal, V. Kumar, S. Shekhar, Multilevel hypergraph partitioning: applications in VLSI domain, *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 7 (1) (1999) 69–79.
- [12] A.L.N. Fred, A.K. Jain, Combining multiple clusterings using evidence accumulation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27 (6) (2005) 835–850.
- [13] P. Hore, L.O. Hall, D.B. Goldgof, A scalable framework for cluster ensembles, *Pattern Recognition* 42 (5) (2009) 676–688.
- [14] D. Cristofor, D. Simovici, Finding median partitions using information-theoretical-based genetic algorithms, *Journal of Universal Computer Science* 8 (2) (2002) 153–172.
- [15] M. Mohammadi, A. Nikanjam, A. Rahmani, An evolutionary approach to clustering ensemble, in: *ICNC '08: Proceedings of the 2008 Fourth International Conference on Natural Computation*, IEEE Computer Society, Washington, DC, USA2008, pp. 77–82.
- [16] V. Filkov, S. Skiena, Heterogeneous data integration with the consensus clustering formalism, in: *Proceedings of Data Integration in the Life Sciences*, Springer2004, pp. 110–123.
- [17] J. Azimi, P. Cull, X. Fern, Clustering Ensembles Using Ants Algorithm, *Methods and Models in Artificial and Natural Computation. A Homage to Professor Miras Scientific Legacy*, 2009, pp. 295–304.
- [18] A.K. Jain, A.L.N. Fred, Data clustering using evidence accumulation, *ICPR '02: Proceedings of the 16th International Conference on Pattern Recognition (ICPR'02)*, vol. 4, IEEE Computer Society, Washington, DC, USA2002, p. 40276.
- [19] A. Topchy, A. Jain, W. Punch, Combining multiple weak clusterings, in: *Proceedings of the Third IEEE International Conference on Data Mining*, IEEE Computer Society2003, p. 331.
- [20] S. Mimaroglu, A.M. Yagci, A binary method for fast computation of inter and intra cluster similarities for combining multiple clusterings, in: *ICIS '09: Proceedings of the 2nd International Conference on Interaction Sciences*, ACM, New York, NY, USA2009, pp. 452–456.
- [21] W.M. Rand, Objective criteria for the evaluation of clustering methods, *Journal of the American Statistical Association* 66 (336) (1971) 846–850  
ArticleType: primary\_article/Full publication date: December 1971/Copyright 1971 American Statistical Association.
- [22] L. Hubert, P. Arabie, Comparing partitions, *Journal of Classification* 2 (1) (1985) 193–218.
- [23] L. Denud, A. Gunoche, Comparison of distance indices between partitions, in: *Data Science and Classification*, Springer2006, pp. 21–28.
- [24] L. Kaufman, P. Rousseeuw, *Finding Groups in Data. An Introduction to Cluster Analysis*, Wiley Interscience, New York, 1990.
- [25] A. Asuncion, D. Newman, UCI machine learning repository, 2007. <[www.ics.uci.edu/mllearn/MLRepository.html](http://www.ics.uci.edu/mllearn/MLRepository.html)>.

**Selim Mimaroglu** is an Assistant Professor in the Department of Computer Engineering at Bahcesehir University. He received Ph.D. in Computer Science from University of Massachusetts Boston, USA. He received M.S. in Computer Science from the same university. His research interests include binary methods in data mining, frequent item set detection, pattern recognition, clustering, text document clustering, and classification.

**Ertunc Erdil** is a M.S. student in the Department of Computer Engineering at Bahcesehir University. He is mainly interested in data mining. His M.S. Thesis is being supervised by Dr. Selim Mimaroglu.