



VICTORIA UNIVERSITY
MELBOURNE AUSTRALIA

Elastic differential evolution for automatic data clustering

This is the Published version of the following publication

Chen, Jun-Xian, Gong, YJ, Chen, WN, Li, Mengting and Zhang, Jun (2019)
Elastic differential evolution for automatic data clustering. *IEEE Transactions on Cybernetics*, 51 (8). pp. 4134-4147. ISSN 2168-2267

The publisher's official version can be found at
<https://ieeexplore.ieee.org/document/8864092>
Note that access to this version may require subscription.

Downloaded from VU Research Repository <https://vuir.vu.edu.au/45263/>

Elastic Differential Evolution for Automatic Data Clustering

Jun-Xian Chen, Yue-Jiao Gong^{ID}, *Senior Member, IEEE*, Wei-Neng Chen, *Senior Member, IEEE*,
Mengting Li, and Jun Zhang^{ID}, *Fellow, IEEE*

Abstract—In many practical applications, it is crucial to perform automatic data clustering without knowing the number of clusters in advance. The evolutionary computation paradigm is good at dealing with this task, but the existing algorithms encounter several deficiencies, such as the encoding redundancy and the cross-dimension learning error. In this article, we propose a novel elastic differential evolution algorithm to solve automatic data clustering. Unlike traditional methods, the proposed algorithm considers each clustering layout as a whole and adapts the cluster number and cluster centroids inherently through the variable-length encoding and the evolution operators. The encoding scheme contains no redundancy. To enable the individuals of different lengths to exchange information properly, we develop a subspace crossover and a two-phase mutation operator. The operators employ the basic method of differential evolution and, in addition, they consider the spatial information of cluster layouts to generate offspring solutions. Particularly, each dimension of the parameter vector interacts with its correlated dimensions, which not only adapts the cluster number but also avoids the cross-dimension learning error. The experimental results show that our algorithm outperforms the state-of-the-art algorithms that it is able to identify the correct number of clusters and obtain a good cluster validation value.

Index Terms—Clustering, differential evolution, elastic encoding, subspace.

I. INTRODUCTION

DATA clustering plays a crucial role in various fields, such as industrial informatics [1], [2]; bioinformatics [3]–[6]; and pattern recognition [7]–[9]. The task is to partition a

dataset into different data subsets named clusters. Commonly, data clustering is performed in an unsupervised way using only the data distribution information. The principle is that we need to maximize the intracluster similarity while minimizing the intercluster similarity. Therefore, data clustering can be regarded as an optimization problem.

Different types of clustering algorithms were developed in the literature, such as the partitioning methods, the agglomerative or divisive hierarchical clustering, the prototype-based methods, and the density-based methods. In this article, we focus on partitioning the clustering which organizes the data objects into a number of exclusive clusters. The most well-known technique in this category is the K -means clustering algorithm [10]. Although it was proposed as early as 1979, the algorithm is still widely used in various fields in recent years [11], [12]. However, there are several shortcomings of the traditional partition-based clustering algorithms such as the K -means. First, they require *a priori* knowledge, namely, the predefined number of clusters. Given a dataset, one may need to use an analytical technique to run the algorithm repeatedly to identify the most suitable K value to improve the performance. Second, these algorithms are typically the local search algorithms which provide only locally optimal solutions. The results would be highly sensitive to the initial partitions of the data objects.

From the perspective of optimization, clustering is shown to be a nondeterministic polynomial NP-hard problem [13] and, therefore, the approximation algorithms, such as the evolutionary computation (EC) algorithms become the potential solvers. Many attempts have been devoted to this area. For example, Krishna and Murty [14] proposed a genetic K -means algorithm (GKA) with a specified cluster number for data clustering and replaced the crossover operator in the genetic algorithm with one step of K -means. In [15], a faster GKA (FGKA) was put forward by Lu *et al.* Merz and Zell [16] adopted some memetic algorithms for clustering gene expression data. In these algorithms, the cluster number is often regarded as an input, due to the fixed-length encoding scheme of EC algorithms. However, when dealing with high-dimensional and big-data problems, it becomes a tedious or even impossible task to predefine the optimal number of clusters. Therefore, we need an automatic way to adjust the number of clusters during the optimization process.

Recently, several techniques have been developed to address this issue. Das *et al.* [17] adopted a fixed-length structure to encode the coordinates of the cluster centroids, together

Manuscript received May 15, 2019; revised August 12, 2019 and September 6, 2019; accepted September 11, 2019. Date of publication October 11, 2019; date of current version August 4, 2021. This work was supported in part by the National Natural Science Foundation of China under Grant 61873095, Grant 61873097, and Grant U1701267, in part by the Guangzhou Science and Technology Planning Project under Grant 201904010211, in part by the Guangdong Natural Science Foundation Research Team Project under Grant 2018B030312003, and in part by the Guangdong-Hong Kong Joint Innovation Platform Project under Grant 2018B050502006. This article was recommended by Associate Editor P. N. Suganthan. (*Corresponding author: Yue-Jiao Gong.*)

J.-X. Chen, Y.-J. Gong, and W.-N. Chen are with the School of Computer Science and Engineering, South China University of Technology, Guangzhou 510006, China (e-mail: gongyuejiao@gmail.com).

M. Li is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou 510006, China.

J. Zhang is with the Victoria University, Melbourne, VIC 8001, Australia (e-mail: junzhang@ieee.org).

This article has supplementary material provided by the authors and color versions of one or more figures available at <https://doi.org/10.1109/TCYB.2019.2941707>.

Digital Object Identifier 10.1109/TCYB.2019.2941707

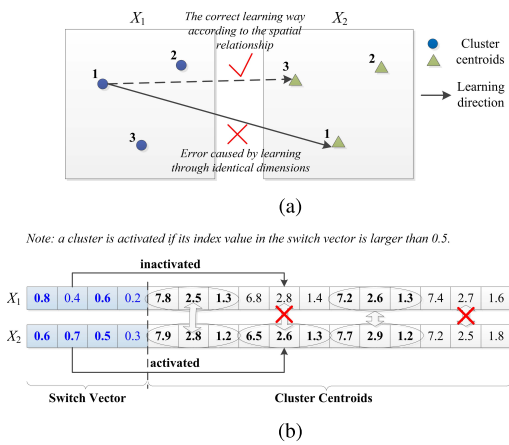


Fig. 1. Illustrations of the cross-dimension learning error. (a) Individuals exchange information in the wrong way (the individual X_1 wrongly takes the index-based dimensions as learning direction). (b) Invalid learning due to the switch mechanism (the activated centroid learns invalid information from the inactivated one).

with a zero-one switch vector to represent the state of each cluster centroid. The cluster centroid is activated when its corresponding switch value is decoded into one. Since the switch vector is updated iteratively, the number of used clusters is adjusted during the optimization. The work of Liu *et al.* [18] follows this method, but it makes improvements in representing the activate state of the cluster centroids. The three adaptive encoding schemes reduce the number of bits in the chromosome to label the ON/OFF state of the cluster centroids. In [19], a genetic algorithm was developed to perform the clustering task, and the algorithm found the right number of clusters according to the average Silhouette width criterion. But the encoding scheme remained the constant-length parameter vector. Hruschka *et al.* [20] improved a genetic algorithm to estimate the number of clusters, and the algorithm performed well in gene expression data. Although it is able to adjust the cluster numbers, this kind of algorithm endures two deficiencies. First, the algorithms encode and optimize the maximum number of cluster centroids, regardless of whether they are activated or not. The search space is large and full of redundancy. The increase of problem dimensionality could lower the performance of the EC algorithms. Second, to the best of our knowledge, all existing EC-based clustering algorithms encounter the problem of “cross-dimension learning error.” Cross-dimension learning error occurs when the individuals in the population interact with each other (e.g., the crossover of the GA and DE algorithms, or when a particle learns information from the others in the PSO). It has two situations, which are illustrated in Fig. 1 and described as follows.

First, the individuals in the traditional EC algorithms exchange information according to the dimension indices. When performing data clustering, this mechanism neglects the spatial relationship between the cluster centroids that a cluster centroid can exchange/learn the information from a totally unrelated cluster centroid. Fig. 1(a) shows a normal situation that two parameter vectors exchange the incorrect information: according to the dimension indices (cluster labels), the centroid 1 of X_1 interacts with the centroid 1 of X_2 . However,

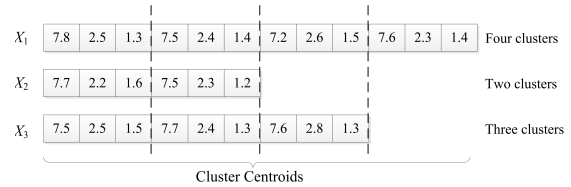


Fig. 2. Elastic encoding scheme with variable lengths.

according to the spatial layout, the correct interaction should be between the centroid 1 of X_1 and the centroid 3 of X_2 because they correlate with each other.

Second, for automatic data clustering, many EC algorithms with a fixed-length structure encode the coordinates of all possible cluster centroids, while using a switch vector to represent the state (activated or not) of each cluster centroid. In these cases, the activated centroid may learn invalid information from the inactivated one, as shown in Fig. 1(b). However, the coordinates of inactivated centroids could be obsolete, which would provide wrong guidance. Because of the above cross-dimension learning error, the EC-based clustering algorithms become inefficient or vulnerable to premature convergence.

To summarize, the state-of-the-art clustering algorithms have two remaining problems to solve. First, the existing differential evolution algorithms use a fixed-length encoding scheme for automatic data clustering, but apply some auxiliary space to represent the activate state of the cluster centroids. The search space is redundant, which reduces clustering efficiency. Second, the existing evolutionary clustering algorithms commonly suffer cross-dimension learning error when performing crossover or individual learning. The existence of the cross-dimension learning error disorders the interactions between individuals and, hence, reduces the performance of a population-based optimization approach. To address the problems, we propose a novel elastic differential evolutionary (E-DE) algorithm in this article. The novelties and contributions of the E-DE algorithm are summarized as follows.

- 1) E-DE adopts an elastic encoding scheme where the population consists of variable-length parameter vectors, each denotes a different number of clusters. As illustrated in Fig. 2, the encoding scheme contains no redundancy: each vector possesses only the coordinates of the cluster centroids that take effect in the clustering of data. E-DE considers each clustering layout (number and positions of clusters) as a whole to perform optimization, and it allows the parameter vector to change its cluster number flexibly during the search process.
- 2) Based on our new elastic encoding scheme, the mutation and crossover operators of DE are sophisticatedly redefined to fit the requirements arising from the variation of parameter vector length. The new mutation adapts the cluster number and centroids according to the differential information of individuals in the current population. Meanwhile, we develop a subspace crossover that takes the spatial layout of clusters into consideration, letting each dimension of one vector learn from the correlated dimensions of the other vectors. To

the best of our knowledge, this article makes the first attempt to eliminate the cross-dimension learning error in evolutionary clustering algorithms.

- 3) Unlike some other peer algorithms that introduce many complex auxiliary techniques, we keep our algorithm conceptually simple. The framework of E-DE is identical to the classical DE algorithm, which iteratively performs the crossover, mutation, and selection steps to refine the population. The adjustment and optimization of cluster number and coordinates is realized through this process in an inherent way. Our algorithm is lightweight and computationally efficient.
- 4) We compare the performance of our algorithm with several state-of-the-art clustering algorithms on synthetic and real-life datasets of various characteristics. The experimental results indicate that the E-DE algorithm outperforms the others.

The remainder of this article is organized as follows. Section II discusses the traditional and state-of-the-art clustering algorithms and points out some drawbacks of the remaining methods. In Section III, we describe the elastic differential evolution algorithm in detail. Section IV shows the experiments and results. The conclusions are demonstrated in Section V.

II. PRELIMINARIES

A. Related Work

The data clustering problem, which aims at partitioning the dataset into different clusters, is an NP-hard problem [13]. The evolutionary algorithms (EAs) are commonly known as the powerful optimizers for the NP-hard problems, so they have received wide attention for data clustering. There are some related survey papers on this topic [21]–[23]. In this section, we review the evolutionary clustering algorithms according to two categories: 1) algorithms with a fixed cluster number and 2) algorithms with a variable cluster number.

1) *Clustering With Fixed Number of Clusters*: In this category, the algorithms require a specified cluster number as *a priori* knowledge, in order to fit the fixed-length encoding scheme of EAs. Many studies made improvements on the evolution operators to enhance the performance of clustering.

- 1) *Single-Objective Optimization Methods*: Krishna and Murty [14] proposed a GKA for data clustering, where the elementary step of K -means was taken to replace the crossover step in the evolution process. Lu *et al.* [15] improved the GKA in terms of the selection and mutation operators and made it run faster. Merz and Zell [16] adopted the memetic algorithms with two mutation and recombination operators for clustering gene expression data. These methods enable the algorithm to quickly converge to a near-optimal solution. Xiang *et al.* [24] employed the DE/best/1 mutation strategy to develop a dynamic shuffled differential EA (SDE) to improve the convergence performance of data clustering. The SDE algorithm divides the initial population into two sub-populations to evolve and combine the best results to form the best optimum at the end of optimization. More recently, a density-based NK hybrid genetic algorithm

was developed in [25]. To determine the number of clusters, the algorithm first generates a large number of solutions and selects the solution with the best validation index. In 2019, Guan *et al.* [26] proposed a particle-swarm-optimized density-based clustering and classification (PODCC) algorithm, which successfully offset the drawbacks of the well-known DBSCAN algorithm for density clustering.

- 2) *Multiobjective Optimization Methods*: In recent years, the evolutionary multiobjective optimization (EMO) methods have also received attention for data clustering. Some studies focus on fuzzy clustering, such as the multiobjective genetic algorithm by Mukhopadhyay *et al.* [27] and the multiobjective differential EA by Saha *et al.* [28]. The algorithms optimize more than one measure for the clustering goodness and generate a set of Pareto optimal solutions for the subsequent operations.

2) *Clustering With a Variable Number of Clusters*: Because the number of clusters can hardly be known in advance, some efforts are made to estimate the number of clusters before performing clustering, such as the last leap and the last major leap methods proposed in [29]. The term “automatic clustering” commonly refers to the methods that are able to adapt the number of clusters automatically without external specification. To fulfill this goal, different kinds of heuristic or local search strategies are incorporated into EAs to enable the algorithms to change the cluster number during the clustering process.

- 1) *Single-Objective Optimization Methods*: Das *et al.* [17] put forward an improved differential EA called ACDE that adjusted the cluster number through a new individual encoding scheme. The ACDE algorithm embeds a zero-one switch vector in the individual representation for activating or inactivating the cluster centroids. An improved ACDE version was proposed by Tam *et al.* [30] using a new changing schema of threshold values. Further, Liu *et al.* [18] developed three adaptive encoding schemes that represent the activate state of the cluster centroids more efficiently. Arellano-Verdejo *et al.* [31] proposed a new hybrid differential EA (DELA), and Tvrdik and Křivý [32] combined differential evolution and K -means. They paid more attention to the center rearrangement. Maulik and Saha [33] proposed a modified differential evolution for automatic fuzzy clustering, which adopted the same individual encoding scheme like the ACDE. In addition, the global best and local best information are utilized to accelerate the searching process. Besides, Chang and Zhang [34] put forward a dynamic niching genetic algorithm for data clustering. Since the niching method is able to preserve the diversity of the population, the algorithm is more likely to find the correct cluster number and well-formed data structure. Sheng *et al.* [35] proposed a multilocal search and adaptive niching-based memetic algorithm (MANM) for data clustering. Two local search strategies merge and split clusters in the evolution process. Then, at the

end of each iteration, an adaptive niching strategy is applied to update the population information in the offspring. Furthermore, Sheng *et al.* [36] improved the performance of the MANM and designed a new memetic algorithm with adaptive multisubpopulation competition and multiniche crowding methods. In [37], the two-stage genetic clustering algorithm (TGCA) adopts a variable-length encoding like the proposed E-DE algorithm. Then, as the name indicated, the TGCA treats the number of clusters and the positions of clusters separately by a two-stage method during the evolution. The experiments showed good performance of TGCA. Besides, in [38], a weighted Gaussian means algorithm was developed, which adopts a feature weighting learning strategy for automatic data clustering.

- 2) *Multiobjective Optimization Methods*: The multiobjective optimization methods are also well suited in this category. Saha *et al.* [39] proposed a framework of multiobjective optimization to solve the semisupervised clustering. Four objective functions are applied: a) the total compactness of the partitioning; b) the total symmetry present in the clusters; c) the cluster connectedness; and d) the adjust rand index. Mukhopadhyay *et al.* [40] put forward a fuzzy C-medoid method for categorical data clustering using the multiobjective genetic algorithm. In the evolution process, the fuzzy cluster variance and cluster separation are taken as two objectives to be optimized simultaneously. Wan *et al.* [41] put forward an adaptive multiobjective memetic algorithm for spectral-spatial fuzzy clustering. The method contains two automatic layers and accomplishes the clustering task of remote sensing images very well. Garza-Fabre *et al.* [42] proposed an improved and more scalable evolutionary approach to multiobjective clustering. It utilizes two encoding schemes for better population representation, that is, the delta locus and the delta binary encoding schemes.

The existing clustering algorithms have the following limitations. First, the algorithms with a fixed number of clusters require prior knowledge that can hardly be known in practical applications. The EC-based automatic data clustering algorithms relieve this problem. However, as discussed in Section I, these algorithms endure the drawbacks of encoding redundancy and cross-dimension learning error. Besides, many algorithms rely on auxiliary techniques, such as niching [35] and population competition [36] to adjust the number of clusters, which increase computational complexity. To tackle the aforementioned limitations, we proposed an elastic differential EA for automatic data clustering. The algorithm is achieved by a new elastic encoding scheme and the corresponding new mutation and crossover operator under the concept of subspace swapping.

B. Validation Indices for Clustering

For the clustering problem, the ultimate purpose of the algorithm is to find a well-formed clustering structure for

the dataset with maximum intercluster distance and minimum intracluster distance. Numerical measures are applied to judge the performance of data clustering, which include external indices and internal indices [43]. An external index measures the clustering results with a reference model, like the partition result provided by a domain expert. In contrast, an internal index directly measures the results based on the intrinsic characteristics of clusters, which does not refer to the external information. There exist a large number of internal indices for clustering validation, such as the Dunn index, CH index, Gamma index, C-index, SIL index, DB index, CS index, I index, and so on. Since there does not exist a “perfect” index so far, many works were devoted to make comparisons between different indices, such as the review work [44].

In this article, we focus on proposing an efficient automatic clustering algorithm. The developed E-DE is generic, which is suitable to adopt any of the above-mentioned indices in the fitness evaluation step. In the major body of this article, we use the DB [45] and I [46] indices that are commonly adopted in the literature. The two indices are simple to understand, lightweight to implement, and are validated to be efficient for various clustering datasets. Besides, we also implement and test the algorithms by adopting the other indices, which are reported in the supplementary file of this article.

1) *DB Index*: In the DB index, each cluster centroid corresponds to an index value R_i

$$R_i = \max_{j,j \neq i} \left(\frac{e_i + e_j}{D_{ij}} \right) \quad (1)$$

where $e_i = (1/N_i) \sum_{x_i \in C_i} \|x_i - m_i\|^2$ is the average error within cluster i ; $D_{ij} = \|m_i - m_j\|^2$ is the distance between the cluster centroids; and N_i denotes the number of data points in cluster C_i . Then, the DB index is formulated as

$$DB(K) = \frac{1}{K} \sum_{i=1}^K R_i. \quad (2)$$

Since we aim at maximizing the intercluster distance and minimizing the intracluster distance, the smaller the DB index value is, the better the clustering result will be.

2) *I Index*: The I index includes three factors that compete with each other, which can be defined as

$$I(K) = \left(\frac{1}{K} \times \frac{P}{E_k} \times D_{\max} \right)^p \quad (3)$$

where K is the number of clusters; P is a constant related to the given dataset, which is set to the number of data samples (N); and E_k and D_{\max} are calculated as

$$E_k = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - m_k\| \quad (4)$$

$$D_{\max} = \max_{i,j=1,2,\dots,K} \|m_i - m_j\| \quad (5)$$

where m_k denotes the k th cluster centroid, x_i is the i th data sample, and C_k means the k th cluster set. Besides, the power $p > 1$ controls the contrast of the I index for different cluster configurations. The three factors in (3), namely, $1/K$, P/E_k , and D_{\max} , have a positive or negative correlation effect on the I index value. Generally, the larger the I index value is, the

better the clustering result will be. The objective of the E-DE algorithm is to maximize the I index value in the evolution process.

3) *Others*: The average intracluster and intercluster distances are commonly adopted to give the details of the clustering results, which can be defined as

$$\text{Intra_dist} = \sum_{k=1}^K \frac{\sum_{\mathbf{x}_i, \mathbf{x}_j \in C_k, j>i} \|\mathbf{x}_i - \mathbf{x}_j\|^2}{\frac{1}{2} \times N_{C_k} \cdot (N_{C_k} - 1)} \quad (6)$$

$$\text{Inter_dist} = \frac{\sum_{i,j=1,2,\dots,K, j>i} \|\mathbf{m}_i - \mathbf{m}_j\|^2}{\frac{1}{2} \times K \cdot (K - 1)} \quad (7)$$

where K is the number of clusters, C_k denotes the k th cluster, and N_{C_k} is the number of data assigned to the cluster C_k .

Besides, the adjusted rand index (ARI) [47] is a commonly used external index to check clustering accuracy. Given two partitions of a dataset of N samples, namely, $X = \{X_1, X_2, \dots, X_r\}$ and $Y = \{Y_1, Y_2, \dots, Y_s\}$, the overlap between X and Y can be summarized in a contingency table $[n_{ij}]$, where each entry n_{ij} denotes the number of objects in common between X_i and Y_j : $n_{ij} = |X_i \cap Y_j|$. Then, the ARI is calculated as

$$\text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}} \quad (8)$$

where n_{ij} , a_i , and b_j are values from the contingency table. a_i represents the row sum of the contingency table while b_j denotes the column sum of the table.

III. PROPOSED ALGORITHM

In this article, we propose an E-DE algorithm to deal with the automatic data clustering problem. Generally, the DE framework is simple and flexible, which suits our proposed elastic encoding and reproduction method. In addition, the literature studies widely validate that DE exhibits powerful performance on continuous numerical optimization as well as real-world applications [48]–[51].

The clustering problem can be defined as follows. Suppose that a dataset $A = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ contains N unlabeled samples, while each sample is represented by a D -dimensional feature vector, namely, $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,D})$. The clustering algorithm aims at partitioning the dataset A into K clusters that are disjoint with each other, which means that each sample in the dataset can only be assigned to one cluster (hard partitions). The partition result is denoted as $\{C_l \mid l = 1, 2, \dots, K\}$, where the clusters have no intersection and the union of them is equal to the dataset A . We use $\lambda_j \in \{1, 2, \dots, K\}$ to denote the cluster label for each sample \mathbf{x}_j , which means that \mathbf{x}_j belongs to the cluster C_{λ_j} . Thus, the clustering result can be represented by a cluster label vector, that is, $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)$.

The proposed E-DE algorithm follows the framework of the basic differential EA, which evolves the candidate solutions (cluster configurations) under some criteria like the clustering validation indices. In this article, we use the DB index and the I index in the fitness evaluation of individuals. To automatically adapt the number of clusters, an elastic individual encoding

Algorithm 1 E-DE for Automatic Data Clustering

- 1: Randomly generate an initial population with np individuals, each has a variable cluster number. (See Section III-B)
- 2: Evaluate the fitness of the individuals in the initial population.
- 3: Repeat the following steps for each individual $i = 1, 2, \dots, np$ until the maximum evaluation times is met:
 - 3.1 Mutation (See Section III-C)
 - 3.1.1 Randomly select three individuals X_{r1} , X_{r2} and X_{r3} ;
 - 3.1.2 Determine the number of clusters K_v in the mutant vector according to Eq. (8);
 - 3.1.3 Combine the three individuals to produce the initial mutant vector V_i ;
 - 3.1.4 Fine-tune the cluster centroids in V_i by the Gaussian location mutation in Eq. (9) to obtain the ultimate mutant vector V_i .
 - 3.2 Crossover (See Section III-D)
 - 3.2.1 Determine the length of crossover by the crossover rate CR according to Eq. (10);
 - 3.2.2 Retrieve the sequence of cluster centroids in V_i to determine the subspace of crossover;
 - 3.2.3 Exchange partial cluster layouts between V_i and the individual X_i by Eq. (12) to generate the trial vector U_i .
 - 3.3 Selection (See Section III-E)
 - 3.3.1 Assign the data objects to the nearest centroids in the trial vector U_i ;
 - 3.3.2 Evaluate the quality of clustering by the DB or I index;
 - 3.3.3 Preserve U_i or X_i to the next generation according to Eq. (13).
- 4: Output the partition represented by the best individual.

scheme is designed and employed. To meet the variable-length condition in our encoding scheme, we put forward the new crossover and mutation operators which are defined on the spatial layout represented by the individuals. These new reproduction operators also avoid the cross-dimension learning error. After a number of iterations, the algorithm terminates and outputs the global best individual. The overall process of the E-DE algorithm is shown in Algorithm 1. For public use, we provide the source code of E-DE online, which can be downloaded at [52]. The following sections describe different components of E-DE in detail.

A. Elastic Encoding Scheme

The number of clusters influences the clustering performance significantly in considering various types of cluster validation measures. However, how to identify the optimal cluster number remains a challenging issue currently. To fulfill this goal, in the EC-based algorithms, it would be promising to let different individuals possess different cluster numbers. Then, with the evolution of population, because of the survival of fitness effect, the individual with the optimal cluster number would dominate the others and survive to the last. In this way, the optimal cluster number is found. The traditional fixed-length encoding scheme is rigid, while the method of embedding a cluster switch vector works but possesses deficiencies. We adopt an elastic encoding scheme: the population consists of variable-length parameter vectors, each denoting a cluster configuration with a different number of clusters. Illustrated in Fig. 2, unlike the switch mechanism that contains some redundant cluster centroids in the parameter vector [17], [18], our individual possesses only the coordinates of the cluster centroids that take effect in the clustering of data. During the evolution process, the cluster number in the individuals X_1 , X_2 , and X_3 can be automatically changed, along with the length of the parameter vectors.

B. Population Initialization

Initialization of the population is to generate np individuals, each with a specified number of clusters and the corresponding cluster centroids. For $i = 1, 2, \dots, np$, we first use a random number generator to determine the cluster number K_i , where K_i should comply with the range $[2, 20]$. Afterward, K_i cluster centroids are randomly selected from the origin dataset, which constitutes the parameter vector of individual i . The length of each parameter vector is $D \times K_i$, where D represents the number of features in the dataset. For example, it can be observed from Fig. 2 that for the dataset of three features, the first individual has four clusters ($K_1 = 4$) and, hence, its parameter vector contains 12 dimensions. After randomly selecting the cluster centroids, the first individual becomes (7.8, 2.5, 1.3, 7.5, 2.4, 1.4, 7.2, 2.6, 1.5, 7.6, 2.3, 1.4).

C. Mutation

The traditional DE mutation randomly selects three individuals and then performs some arithmetic combination operation on them, in order to generate a mutant vector. The premise of the method is that the individuals, or parameter vectors, in the population possess the same number of dimensions. In the proposed E-DE, since we develop the variable-length encoding for different individuals, the traditional DE mutation is no longer suitable. Therefore, we design a new mutation operator, which consists of two phases described as follows.

1) *Phase One*: In the first phase, three individuals X_{r1} , X_{r2} , and X_{r3} are randomly selected from the current population. The number of clusters in the mutant vector V_i is calculated as

$$K_v = K_{r1} + F \times (K_{r2} - K_{r3}) \quad (9)$$

where K_i represents the cluster number of the i th individual, F is the scaling factor, and K_v should be a round number. Note that (8) is similar to the original DE/rand/1 mutation: the random selection of the base vector enhances the global exploration of the algorithm, while the differential length-based perturbation realizes the contour match effect when searching the optimum. This way, the cluster number would be optimized and converged eventually. Besides, we restrict the cluster number in the range of $[2, 20]$. If K_v is smaller than 2, we set $K_v = 2$; otherwise, if K_v is larger than 20, we set $K_v = 20$. After deciding the number of clusters in the mutant vector, the cluster centroids are determined according to the following three situations.

- 1) $K_v > K_{r1}$, which means that the mutant vector V_i has more clusters than the individual X_{r1} : the cluster centroids of individual X_{r1} are used as the fundamental structure to construct the mutant vector V_i . Then, we need to add $(K_v - K_{r1})$ more cluster centroids to V_i . They are randomly selected from the individuals X_{r2} and X_{r3} . It should be noted that the selected cluster centroids are distinct from each other.
- 2) $K_v < K_{r1}$, which means that the mutant vector V_i has fewer clusters than the individual X_{r1} : we also start with the cluster centroids of X_{r1} and then randomly deleted $(K_{r1} - K_v)$ cluster centroids from it. The resulting mutant vector V_i contains exactly K_v clusters.

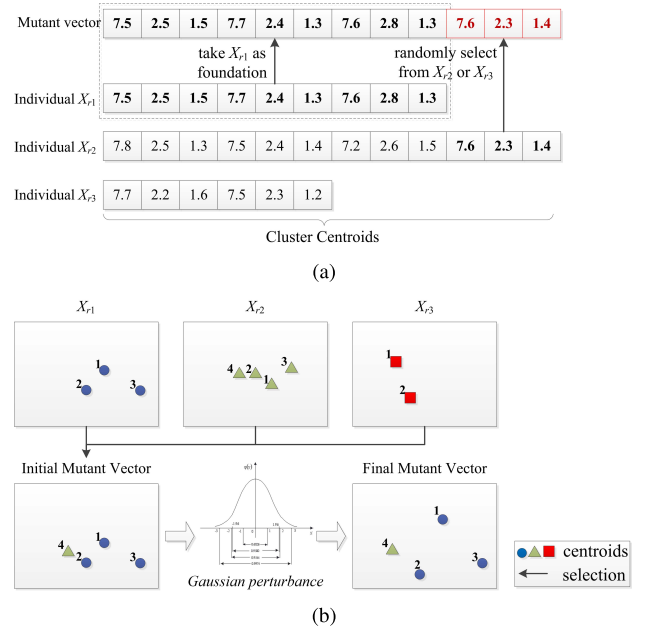


Fig. 3. Mutation process of the E-DE algorithm. (a) Parameter vectors (X_{r1} , X_{r2} , and X_{r3} are individuals randomly selected from the current population). (b) Spatial layouts.

- 3) $K_v = K_{r1}$, which means that the mutant vector V_i and the individual X_{r1} possess the same number of clusters: in this case, the cluster centroids of individual X_{r1} are directly copied to the mutant vector V_i .
- 2) *Phase Two*: By considering the above procedures, the increment or decrement of the cluster centroids comes from the pool of cluster centroids held by the previous population. To alter the centroids, in the second phase of our mutation, the initial mutant vectors further undergo a Gaussian mutation. The operation changes the geographic locations of centroids, which improves the exploration ability of the algorithm. For each cluster centroid that is contained in the mutant vector V_i , a random number $r \in [0, 1]$ is generated. If r is smaller than the mutation probability M_u , the cluster centroid in the mutant vector V_i is adjusted to

$$\mathbf{m}_i^k = \mathbf{m}_i^k + G(0, 0.1) \times (\mathbf{ub} - \mathbf{lb}) \quad (10)$$

where \mathbf{m}_i^k represents the k th cluster centroid represented by the mutant vector, $G(0, 0.1)$ is a Gaussian random number generator with 0 mean and 0.1 standard deviation, and \mathbf{ub} and \mathbf{lb} represent the upper and lower bounds of the dataset. The utilization of the Gaussian perturbation considers the following factors: the Gaussian mutation is able to enrich the population diversity without destroying the formulated solution structures. In comparison, if the uniformly random mutation is applied instead, the locations of centroids would change to a large extent, which alters the cluster layout severely. So the Gaussian mutation helps to keep a more well-formed clustering structure and enables the DE algorithm to search the problem in a smooth way.

Fig. 3(a) and (b) shows a mutation example. Suppose that $K_{r1} = 3$, $K_{r2} = 4$, $K_{r3} = 2$, and $F = 0.5$, then the calculated K_v is 4. Since K_v is larger than K_{r1} , one cluster centroid

should be selected from the individual X_{r2} or X_{r3} to fill the mutant vector. After that, the location perturbation is further performed on the mutant vector. The new operator not only adapts the cluster number in generating the mutant vector, but also incorporates the centroid distribution information from different individuals in the population. Next, the ultimate mutant vector will intersect with the target vector of individual i by the crossover operator.

D. Crossover

The existing DE crossover methods will lead to the cross-dimension learning error in data clustering. This is because, as illustrated in Fig. 1 and Section I, they exchange the information between the parameter vectors based on the dimension index abruptly, which ignores the spatial relationships of different cluster centroids. Besides, they are unavailable to be directly applied to variable-length parameter vectors. Therefore, we put forward a new crossover operator to solve the above two problems. The crossover can effectively exchange the information between the variable-length individuals without incurring the cross-dimension learning error.

Based on the exponential crossover of DE, we define and incorporate a subspace crossover in E-DE, which is described as follows and illustrated in Fig. 4.

Step a): First, determine the crossover segment for each individual: a random location n is chosen as the starting point and the following L genes will be used for crossover. The crossover length L is determined according to the crossover rate CR as

$$L = L + 1, \text{ while } \text{rand}(0, 1) < CR \quad (11)$$

where the initial value of L is zero. The L value determines the length of genes that are going to be swapped in the crossover process. This incremental way has been demonstrated in the original DE algorithm, which is called the exponential crossover [53].

Step b): Retrieve the cluster centroids that are contained in the mutant vector V_i from the selected gene segment in step a). Namely, the first cluster centroid comes from the randomly selected gene location n , while the last centroid is from the gene location $n+L-1$. Based on the retrieved sequence of cluster centroids, the subspace to exchange information in the crossover is defined as follows. First, the center of the subspace, denoted as cs_i , is calculated as the average of all the retrieved centroids. Then, suppose that the first and last centroids in the sequence are m_i^{first} and m_i^{last} , respectively; and the range of the subspace on each dimension is defined as the distance between m_i^{first} and m_i^{last} on each dimension

$$d_{i,j} = \left| m_{i,j}^{\text{first}} - m_{i,j}^{\text{last}} \right| \quad (12)$$

where $m_{i,j}^{\text{first}}$ and $m_{i,j}^{\text{last}}$ denote the j th attribute value of the first and last centroids. Particularly, assume that $c_{i,j}$ is the j th coordinate of the center cs_i and that $d_{i,j}$ is the distance between m_i^{first} and m_i^{last} on dimension j , and the j th dimension of the subspace ranges from $(c_{i,j} - d_{i,j}/2)$ to $(c_{i,j} + d_{i,j}/2)$.

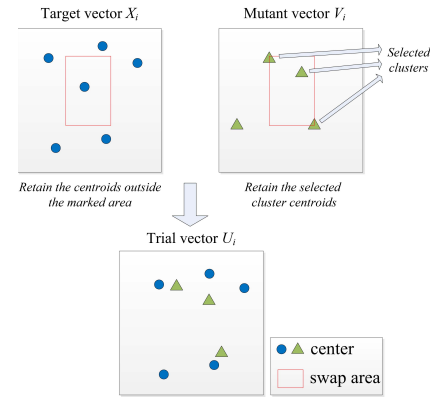


Fig. 4. Illustration of the subspace crossover process of the E-DE algorithm.

We use the distance between m_i^{first} and m_i^{last} to let the range of the subspace be more sensitive to the retrieved sequence of cluster centroids. It helps to increase the population diversity. In contrast, if we take the distance between the maximum and minimum value of all retrieved centroids for calculation, different selected gene sequences could correspond to an identical range of subspace. Besides, in the current stage, our E-DE concentrates on continuous data clustering. If the features are discrete, some discrete distance measures (Hamming distance, permutation distance, etc.) should be used according to the practical requirements.

Step c): For the other parameter vector that joins the crossover, that is, the target vector X_i , we use the identical area that is defined in step b) to retrieve the subset of cluster centroids in X_i . Then, the selected cluster centroids of V_i and X_i are exchanged to generate the trial vector U_i

$$U_i = V_i \otimes X_i. \quad (13)$$

The operator \otimes corresponds to the information exchange process in step c), which is illustrated in Fig. 4. Note that the cluster numbers in the crossover area in X_i and V_i could be different. This is so that the cluster number in the generated trial vector could be distinct from the current individual X_i . Therefore, the proposed elastic length encoding scheme, together with the mutation and crossover, is well suited for the automatic data clustering problem that requires adjusting the number of clusters. Besides, if the cluster number in the trial vector U_i exceeds the maximum cluster number $K_{\max} = 20$ or lower than the minimum cluster number $K_{\min} = 2$, it will adopt the cluster number of target vector X_i , but each cluster centroid in the individual X_i will be replaced by the nearest cluster centroid in the mutant vector V_i .

In addition to being capable of adjusting the cluster numbers, the proposed subspace crossover inherently avoids the cross-dimension learning error since it considers the spatial relationship between clusters. Note that the target and mutant vectors represent different spatial layouts of cluster centroids, as depicted in Fig. 4. In the crossover process, first, a subspace is selected as the swap area and, then, the target and mutant vectors exchange their subsets of cluster centroids that are located within the subspace. This mechanism ensures that the exchanged clusters have a spatial relationship, which

TABLE I
PROPERTIES OF 19 DATASETS USED IN THE CLUSTERING EXPERIMENTS

Datasets	# Samples	# Features	# Clusters (Ground truth)
Iris	50	3	3
Cancer	683	9	2
Glass	214	9	6
Vowel	871	3	6
Zoo	101	7	7
Vote	435	16	2
Ecoli	336	7	8
Seed	210	7	3
Compound	339	2	6
R15	600	2	15
Dim2	1351	2	9
G2	2048	2	2
Dim32	1024	32	16
Customer	440	6	6
Unbalance	6500	2	8
Spiral	312	2	3
Mockdata1	1000	2	4
Mockdata2	1000	2	2
Mockdata3	600	2	4

are different from the traditional DE algorithms that perform crossover through the cluster indices. Thus, the generated trial vector inherits different cluster centroids from the target or mutant vectors from a spatial point of view. To summarize, the subspace crossover not only fits our elastic encoding scheme but also avoids the cross-dimension learning error of the traditional DE algorithms.

E. Selection

To perform selection first, the fitness of the trial vector U_i is calculated. A data assignment step is required when the cluster centroids are determined. For each data point in the dataset, calculate its distance to all cluster centroids in U_i . Then, each data object is assigned to its nearest cluster centroid.

Note that sometimes there may exist some clusters which contain less than 2 data points, which are considered empty clusters. The empty clusters will affect the well-formed cluster structure. However, if they are simply removed, the population would have a hidden property to reduce the total number of clusters during optimization, which goes against our intention. Therefore, we use a random reassignment strategy [17] instead. Particularly, for the individual that contains an empty cluster, it is rearranged according to the following procedure. Select N/K data points from the dataset, use the mean of the data points as the new cluster centroid, and redo the data assignment. Repeat this step until no empty cluster comes up. Then, update the individual accordingly.

After the data assignment step, the dataset is partitioned into clusters. The validation indices can be applied to evaluate the quality of the partition. Since how to design a good validation index for clustering is beyond the scope of this article, here, we simply adopt the well-known DB index and I index described in Section II-B. The selection of E-DE is identical to the classical DE selection, namely, the algorithm compares the fitness of the current individual X_i and the trial vector U_i , and it selects the better one to the next generation

$$X_{i,G+1} = \begin{cases} X_{i,G}, & \text{if } f(X_{i,G}) \text{ is better than } f(U_{i,G}) \\ U_{i,G}, & \text{otherwise.} \end{cases} \quad (14)$$

Note that if the DB index is applied, the smaller the f value is, the better the clustering results are. Differently, for the I index, the higher f value is preferred.

IV. EXPERIMENTS

A. Experimental Setup

In this section, we conduct several experiments to validate the performance of the proposed E-DE algorithm. The experiments are based on 19 datasets (summarized in Table I) that are widely adopted to examine the state-of-the-art clustering algorithms [29], [54]–[61]. These datasets possess various characteristics. The Iris, Cancer, Glass, Zoo, Vote, Ecoli, Seed, and Customer datasets are standard datasets coming from the UCI machinery repository [54]. In addition, the Compound [55] is a graphical-shaped dataset with six clusters and 339 data objects in two dimensions. The R15 [56] contains Gaussian clusters in a ringed region, with 600 data objects in total. The Dim2 dataset [57] is a synthetic clump including nine clusters, with 1351 data objects in two dimensions. The G2 dataset [58] is a typical 2-class collection, with 2048 data in two dimensions. The Dim32 [59] is a high-dimensional dataset that contains 32 dimensions, with 1024 data objects in total. The Unbalance dataset [60] has eight clusters of quite different sizes and densities, with 6500 data objects in total. The Spiral [61] is a nonconvex dataset composed of three different trajectories of the spiral, with 312 data objects in two dimensions. Further, we test three mock datasets, that is, Mockdata1, Mockdata2, and Mockdata3. They are produced according to three functions that are written by Jeroen Kools to generate different pathological clustering datasets. Commonly, the centroid-based clustering algorithms can hardly handle nonconvex datasets like the spiral. The nonconvex datasets are tested in this article for the purpose of investigations.

We compare the E-DE algorithm with six representative clustering algorithms, namely, the K -means [10], BF- K -means [44], S-DE [24], ACDE [17], MANM [35], and TGCA [37]. The K -means algorithm requires a predefined cluster number as input to partition the dataset into K clusters. The brute-force K -means (BF- K -means) is to run K -means several times varying the number of clusters from the K_{\min} to K_{\max} (i.e., 2 to 20), and select the best clustering solution. The S-DE algorithm is a global optimization algorithm, while it also requires the cluster number as the prior condition. The ACDE algorithm adopts a relatively automatic parameter vector encoding scheme, adjusting the cluster number through the zero-one switch. The MANM algorithm takes the memetic algorithm as a framework, and it utilizes the multilocal search and adaptive niching technique for automatic data clustering. The TGCA algorithm adopts the two-stage selection and mutation operations to adjust the cluster number and centroids. The individuals that have the same cluster number will be partitioned into the same subpopulation to perform crossover. In the experiments, we adopt the parameter settings that are suggested in their original references. Besides, all EC algorithms are assigned a maximum number of fitness evaluation times, 10^6 . Both the K -means and BF- K -means algorithms are terminated when the cluster centroids are not updated.

For the E-DE algorithm, the size of population np is set to $10 \times D$, where D is the data dimensions. The minimum and maximum cluster numbers are set to 2 and 20. The scalar factor F and the crossover rate CR are set to 0.5 and 0.4, while

TABLE III
COMPARING THE CLUSTERING RESULTS (MEAN AND STANDARD DEVIATION) OF THE E-DE ALGORITHM WITH THE OTHER SIX CLUSTERING ALGORITHMS USING THE I INDEX. NOTE: THE CORRECT CLUSTER NUMBER IS PROVIDED UNDER THE NAME OF EACH DATASET

Problem	#Counts	K-means	S-DE	BF-K-means	ACDE	MANM	TGCA	E-DE
Iris (3)	Avg. Fitness	3.238±0.164	5.143±0.032	6.640±0.277	7.239±1.012	4.449±0.330	5.310±0.533	5.329±0.433
	Intra Distance	0.734±0.029	0.633±0.031	0.933±0.053	0.842±0.067	1.032±0.019	0.854±0.066	0.833±0.054
	Inter Distance	2.323±0.325	3.537±0.032	3.587±0.187	4.642±0.547	2.039±0.420	4.033±0.027	3.777±0.529
	Avg. ARI	0.347	0.372	0.906	0.922	0.914	0.929	0.953
Cancer (2)	Avg. Fitness	0.417±0.215	0.724±0.143	1.187±0.050	2.053±1.449	1.277±0.287	1.783±0.550	2.338±0.640
	Intra Distance	8.220±0.045	7.239±0.053	8.785±0.042	9.323±0.405	10.662±0.153	9.326±0.523	9.796±0.380
	Inter Distance	10.723±1.169	13.436±0.564	12.032±0.270	18.325±3.733	15.829±1.435	20.067±2.325	21.721±3.659
	Avg. ARI	0.249	0.313	0.827	0.924	0.892	0.875	0.929
Glass (6)	Avg. Fitness	2.321±0.442	4.248±1.321	2.132±0.029	6.223±1.420	4.126±0.420	4.023(6)	6.282(3)
	Intra Distance	2.122±0.432	2.324±0.560	1.933±0.012	2.121±0.223	1.546±0.223	6.517±0.220	3.325±0.367
	Inter Distance	5.320±0.442	6.620±0.320	5.432±0.325	9.323±0.630	8.329±0.524	2.333±0.340	2.313±0.420
	Avg. ARI	0.848	0.906	0.550	0.215	0.238	0.360	0.860
Vowel (6)	Avg. Fitness	1.703±0.132	3.450±0.227	2.929±0.103	4.032±0.323	3.222±0.324	3.332±0.145	3.814±0.547
	Intra Distance	232.325±3.120	246.425±12.430	387.453±59.450	291.545±52.391	312.324±24.030	259.314±22.830	253.526±12.556
	Inter Distance	738.524±32.670	1149.625±42.535	942.374±70.265	1193.472±128.566	1489.125±32.240	1207.643±41.250	1232.332±66.430
	Avg. ARI	0.887	0.938	0.270	0.563	0.697	0.263	0.385
Zoo (7)	Avg. Fitness	0.846±0.067	1.643±0.140	1.162±0.069	1.514±0.261	1.222±0.082	1.862±0.009	1.761±0.062
	Intra Distance	1.756±0.162	1.988±0.106	2.513±0.380	2.414±0.494	3.654±0.166	2.278±0.159	2.360±0.293
	Inter Distance	4.367±0.339	5.768±0.293	5.430±0.417	6.646±0.699	5.700±0.259	6.574±0.724	6.728±0.733
	Avg. ARI	0.2400	0.234	0.313	0.287	0.526	0.369	0.246
Vote (2)	Avg. Fitness	0.188±0.078	0.124±0.142	0.718±0.030	0.584±0.236	0.473±0.040	0.601±0.019	1.328±0.025
	Intra Distance	1.921±0.014	1.552±0.012	2.017±0.042	1.584±0.042	2.048±0.023	1.270±0.088	2.025±0.006
	Inter Distance	1.732±0.106	2.433±0.120	2.394±0.044	2.929±0.390	2.108±0.093	2.489±0.220	3.953±0.048
	Avg. ARI	0.129	0.163	0.891	0.906	0.906	0.891	0.992
Ecoli (8)	Avg. Fitness	0.374±0.157	1.443±0.064	0.976±0.032	2.317±0.527	2.045±0.035	1.547±0.067	1.354±0.230
	Intra Distance	0.302±0.025	0.323±0.020	0.364±0.058	0.365±0.094	0.374±0.018	0.368±0.054	0.324±0.022
	Inter Distance	0.523±0.024	0.742±0.024	0.555±0.032	1.202±0.328	0.942±0.029	0.933±0.051	0.923±0.087
	Avg. ARI	0.408	0.594	0.112	0.112	0.231	0.231	0.262
Seed (3)	Avg. Fitness	1.438±0.042	2.435±0.080	2.946±0.023	3.379±0.321	2.329±0.064	2.901±0.072	2.845±0.462
	Intra Distance	1.626±0.044	1.622±0.007	2.234±0.019	1.976±0.230	2.424±0.024	1.921±0.052	1.953±0.244
	Inter Distance	5.073±0.124	5.832±0.015	5.334±0.112	7.533±0.244	7.325±0.045	7.243±0.235	6.329±0.576
	Avg. ARI	0.154	0.187	0.919	0.94	0.942	0.932	0.855
Compound (6)	Avg. Fitness	2.772±0.127	3.823±0.443	3.592±0.033	4.189±0.232	4.551±0.018	4.223±0.322	3.825±0.630
	Intra Distance	3.235±0.127	3.796±0.162	6.722±0.053	4.422±0.740	5.603±0.223	4.672±0.529	4.282±0.813
	Inter Distance	15.233±0.344	18.257±0.653	19.191±0.688	21.324±0.931	24.465±0.331	23.627±0.240	25.514±3.253
	Avg. ARI	0.784	0.828	0.928	0.912	0.936	0.940	0.840
R15 (15)	Avg. Fitness	1.382±0.353	1.923±0.029	16(271,12(3))	15(271,16(3))	15	15(261,16(3),17(1))	15(281,17(2))
	Intra Distance	1.694±0.134	1.773±0.038	0.446±0.032	1.325±0.932	0.539±0.074	0.573±0.064	0.533±0.012
	Inter Distance	8.176±0.420	9.328±0.016	5.854±0.035	8.226±2.453	5.923±0.161	5.374±0.123	6.268±0.132
	Avg. ARI	0.149	0.125	0.688	1.000	0.234	0.34	0.982
Dim2 (9)	Avg. Fitness	3.824±0.363	3.148±0.426	62.133±11.560	127.536±56.743	7.226±0.266	115.433±8.327	128.747±7.337
	Intra Distance	64322.933±171.383	74171.311±0.000	8548.503±230.599	13297.671±0.000	13297.671±0.000	13297.671±0.000	13297.671±0.000
	Inter Distance	632561.326±1215.367	627468.325±1513.632	612923.753±256.326	632648.468±5235.278	572125.333±524.685	612566.126±2215.332	617739.822±1292.223
	Avg. ARI	0.05	0.203	0.469	1.000	1.000	1.000	1.000
G2 (2)	Avg. Fitness	8.367±1.624	14.377±0.000	31.137±3.520	28.359±5.052	31.166±1.231	29.574±2.481	22.855±3.971
	Intra Distance	11.919±0.323	17.631±0.000	17.947±0.424	17.509±0.562	17.947±0.339	17.654±0.339	17.020±1.258
	Inter Distance	84.314±5.212	110.382±0.000	141.868±0.737	145.787±0.046	141.980±4.162	148.656±12.676	146.253±15.739
	Avg. ARI	0.106	0.143	1.000	1.000	1.000	1.000	1.000
Dim32 (16)	Avg. Fitness	0.149±0.071	0.366±0.057	2(181,5(8),6(4))	3(221,4(5),16(3))	3	6(271,7(3))	16
	Intra Distance	174.279±21.999	225.128±0.000	154.321±10.090	1.464±0.173	3.171±0.582	2.681±0.761	5.366±1.936
	Inter Distance	357.748±34.157	505.580±0.469	369.233±81.479	689.055±181.976	425.219±9.228	626.308±154.728	18.997±0.000
	Avg. ARI	0.177	0.313	0.444	0.131	1.000	0.141	438.504±7.786
Customer (6)	Avg. Fitness	2.857±1.174	8.683±2.584	2.356±0.049	3(251,6(4),7(1))	5(261,6(4))	6(271,7(3))	6(271,5(3))
	Intra Distance	20740.059±2465.931	22543.934±3215.112	23649.710±3060.238	31937.862±7486.158	61684.642±2429.622	33259.212±6538.291	31125.805±8602.349
	Inter Distance	42848.364±8671.446	68462.293±3843.452	44800.441±2459.566	99959.108±24283.764	129627.096±28596.391	126137.721±2334.650	97819.753±15557.605
	Avg. ARI	0.935	0.915	0.529	0.238	0.238	0.334	0.845
Unbalance (8)	Avg. Fitness	54.810±7.132	57.242±0.674	107.451±29.880	96.515±34.284	68.115±2.772	91.328±2.336	108.676±5.449
	Intra Distance	17480.289±1136.566	17502.468±113.471	24610.960±669.354	18479.130±2864.790	24681.605±1007.485	20578.820±3146.750	17912.095±3443.909
	Inter Distance	131277.934±10129.727	137691.053±1356.423	179398.534±966.687	204641.839±4844.924	176973.351±3714.590	213659.393±5748.629	230571.609±11320.226
	Avg. ARI	0.244	0.672	0.234	0.481	0.234	0.234	0.898
Spiral (3)	Avg. Fitness	0.616±0.108	1.886±0.329	0.733±0.078	2.118±0.348	1.977±0.043	1.863±0.082	5(15,4(8),3(7))
	Intra Distance	5.649±0.341	5.329±0.009	9.932±0.091	6.578±0.998	7.794±0.272	7.525±0.385	1.617±0.253
	Inter Distance	13.845±1.610	18.819±1.227	12.227±0.133	22.592±3.051	23.750±0.786	23.032±1.534	7.251±0.897
	Avg. ARI	0.114	0.093	0.379	0.267	0.683	0.442	20.918±2.560
Mockdata1 (4)	Avg. Fitness	0.874±0.023	2.439±0.037	1.081±0.034	2.318±0.066	2.079±0.044	2.158±0.048	2.010±0.018
	Intra Distance	3.999±0.001	3.324±0.061	3.910±0.074	7.110±0.082	8.045±0.093	7.481±0.059	5.515±0.042
	Inter Distance	11.164±0.187	15.449±0.302	11.366±0.175	22.578±0.139	22.585±0.370	22.280±0.340	17.695±0.203
	Avg. ARI	0.524	0.375	0.929	0.447	0.529	0.473	0.938
Mockdata2 (2)	Avg. Fitness	2.497±0.015	2.682±0.072	2.505±0.049	3.112±0.017	2.491±0.065	2.495±0.029	5(28,4(2))
	Intra Distance	3.452±0.028	3.647±0.036	4.036±0.067	4.796±0.088	8.793±0.056	4.286±0.055	2.569±0.020
	Inter Distance	14.127±0.139	16.864±0.415	14.174±0.233	17.280±0.257	27.135±0.228	17.083±0.728	4.058±0.036
	Avg. ARI	0.239	0.313	0.326	0.494	0.369	0.487	17.769±0.467
Mockdata3 (4)	Avg. Fitness	2.316±0.025	2.949±0.048	2.672±0.062	3.831±0.062	3.826±0.073	2.582±0.045	3.033±0.079
	Intra Distance	7.198±0.026	7.886±0.036	16.894±0.360	17.038±0.074	17.354±0.031	7.953±0.029	13.056±0.084
	Inter Distance	37.796±1.314	44.854±1.784	39.899±1.822	55.08±1.936	60.184±1.273	48.309±1.458	46.752±1.789
	Avg. ARI	0.472	0.502	0.389	0.624	0.704	0.968	0.971
Total CSR	N/A	N/A	32.632%	44.386%	48.070%	53.158%	72.982%	
Total ARI	0.370	0.431	0.550	0.586	0.650	0.632	0.752	

to our datasets in Table I. Unlike the traditional K-means algorithm, the BF-K-means is conducted with different cluster numbers and chooses the best result as the solution. The other four algorithms, namely, ACDE, MANM, TGCA, and the proposed E-DE, automatically adapt the cluster number during optimization. The best result of each dataset is selected and marked in boldface according to the correctness of #counts and the average DB index value. Meanwhile, the best result according to the ARI is also marked in boldface.

It can be observed that our E-DE algorithm exhibits better performance when comparing other data clustering algorithms on most datasets. E-DE can always find the correct number of clusters where it reaches approximately 81% clustering

accuracy on the 19 datasets. For the compared automatic data clustering algorithms, we can observe that the ACDE algorithm fails to find the correct cluster number in many cases. The ACDE is sometimes unstable that its clustering results in different runs vary widely, as can be seen from the #counts for the Vowel, Vote, Ecoli, Dim2, G2, and Customer datasets. For the MANM column in Table II, the memetic algorithm has similar clustering accuracy with the ACDE. Its local search strategies consume many fitness evaluations, which leads the algorithm to exploit too much at the beginning. For the two-stage TGCA algorithm, its performance is slightly better than the MANM. But the crossover way that was used still has the cross-dimension learning error that the clustering accuracy

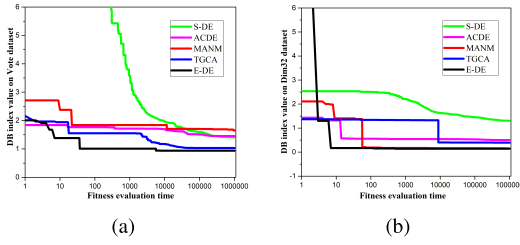


Fig. 5. DB index convergence curves of S-DE, ACDE, MANM, TGCA, and E-DE algorithms. (a) Vote dataset. (b) Dim32 dataset.

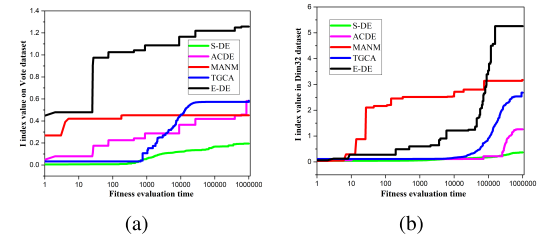


Fig. 6. I index convergence curves of S-DE, ACDE, MANM, TGCA, and E-DE algorithms. (a) Vote dataset. (b) Dim32 dataset.

cannot improve so much. Besides, although enumerating all possible cluster numbers, the general performance of the BF- K -means is lower than all evolutionary clustering algorithms. The reasons are: 1) the K -means algorithm only minimizes the intracluster distances, whereas the EC algorithms optimize the cluster validation index that considers both intracluster and intercluster distances and 2) the K -means is a local search algorithm that can find a local optimum only, whereas the EC algorithms try to find the global optimum of the problem.

For an evolutionary clustering algorithm, the dimension of the problem space is at least the number of features times the number of clusters. Commonly, in the field of EC, the problem possessing higher than 500 dimensions is considered as a high-dimensional problem. The Dim32 dataset contains 32 dimensions and 16 clusters, whose search space has at least $32 \times 16 = 512$ dimensions. It can be considered a high-dimensional problem in the EC field. From Table II, we can see that the E-DE algorithm performs well for the Dim32 dataset and surpasses the other six methods a lot. This is because the traditional evolutionary clustering algorithms endure the problems of encoding redundancy and/or cross-dimension learning error, whereas our elastic encoding and subspace crossover prevents the E-DE algorithm from these deficiencies. Nevertheless, for some extremely high-dimensional clustering situations like deep feature-based clustering, the dimension reduction or the cooperative evolution framework [62], [63] is required for the evolutionary clustering methods.

Further, we compare the convergence curves of the E-DE algorithm with those of the S-DE, ACDE, MANM, and TGCA algorithms. (We do not consider the K -means and BF- K -means algorithms because such methods terminate in a short time.) Taking the Vote and Dim32 datasets as examples, Fig. 5 compares the convergence curves of the DB index fitness values. It can be observed that the E-DE algorithm converges faster than S-DE, ACDE, MANM, and TGCA.

2) *Validation by the I Index:* Next, we use the I index to examine the seven clustering algorithms, for which the results are presented in Table III. It can be observed that our E-DE algorithm has an outstanding performance compared with other clustering algorithms for most of the datasets. E-DE tends to find the correct number of clusters so that it reaches an approximately 73% CSR on the 19 datasets. Meanwhile, the total ARI obtained by E-DE is much higher than the others. Taking the Vote and Dim32 datasets as examples to examine the convergence of the compared algorithms evaluated by the I index value, the curves are presented in Fig. 6. From

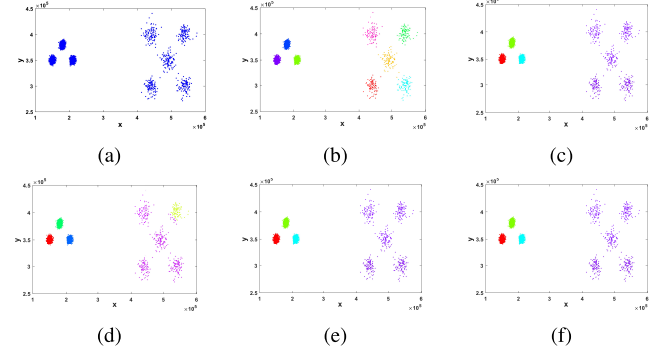


Fig. 7. Clustering results of Unbalance dataset using the I index: (a) is the original dataset while (b)–(f) are the clustering results of the E-DE, BF- K -means, ACDE, MANM, and TGCA algorithms.

Fig. 6(a), we can see that the E-DE converges much faster than the S-DE, ACDE, MANM, and TGCA. In Fig. 6(b), the E-DE algorithm eventually obtains better fitness than the other four clustering algorithms.

Further, the clustering results of the 2-D Unbalance are displayed in Fig. 7. Fig. 7(a) is the original dataset without adopting the clustering algorithm, while the remaining sub-figures show the results of the E-DE, BF- K -means, ACDE, MANM, and TGCA algorithms in turn. The visual results also validate the powerfulness of our E-DE algorithm. Particularly, it can be seen from Fig. 7(a) that, for the Unbalance dataset, the groups of points possess different densities. Then, as shown in Fig. 7(b), E-DE correctly partitions the original dataset into eight clusters, which proves the effectiveness of our elastic encoding strategy and the corresponding evolution operators. In comparison, in Fig. 7(d), the ACDE cannot divide the five clusters in the right part thoroughly, and, in Fig. 7(c)–(f), the BF- K -means, MANM, and TGCA wrongly regard the five different clusters as one cluster.

3) *Significance Tests:* We conduct statistical tests based on the fitness evaluation indices to ensure that our improvement of performance is statistically significant. Particularly, the Wilcoxon's rank-sum test is performed for one-on-one comparisons, while the Kruskal–Wallis (KW) test is conducted for multiple comparisons. The significance test results using DB and I indices are reported in Tables IV and V, respectively. Meanwhile, in the tables, the symbol (+) indicates that the performance of E-DE is significantly better than the compared algorithm(s), whereas the symbol (–) indicates the opposite. Considering Wilcoxon's rank-sum test results, we can observe

TABLE IV

P-VALUES OF WILCOXON RANK-SUM TESTS AND KW TESTS BETWEEN OUR METHOD AND THE OTHER SIX CLUSTERING ALGORITHMS OF THE DB INDEX ON 19 DATASETS. (+) INDICATES THAT THE PERFORMANCE OF E-DE IS SIGNIFICANTLY BETTER THAN THE COMPARED ALGORITHM WITH A CONFIDENCE LEVEL OF 95%, WHILE (−) REPRESENTS THE OPPOSITE

Method		Iris	Cancer	Glass	Vowel	Zoo	Vote	Ecoli	Seed	Compound
Wilcoxon	E-DE to K-means	0.008(+)	0.010(+)	0.002(+)	0.002(-)	0.002(+)	0.037(+)	0.160	0.045(+)	0.010(+)
	E-DE to S-DE	0.021(+)	0.010(+)	0.002(+)	0.002(-)	0.002(+)	0.002(+)	0.027(+)	0.037(+)	0.019(+)
	E-DE to BF-K-means	0.038(+)	0.002(+)	0.042(+)	0.013(+)	0.002(+)	0.002(+)	0.026(+)	0.010(+)	0.037(+)
	E-DE to ACDE	0.002(+)	0.002(+)	0.002(+)	0.084	0.019(-)	0.027(+)	0.003(+)	0.007(+)	0.023(+)
	E-DE to MANM	0.049(+)	0.028(+)	0.008(+)	0.017(-)	0.002(+)	0.002(+)	0.206	0.007(-)	0.010(+)
KW	E-DE to TGCA	0.016(+)	0.002(+)	0.034(+)	0.063	0.002(+)	0.002(+)	0.008(+)	0.002(-)	0.016(+)
	E-DE to Others	0.010(+)	0.007(+)	0.008(+)	0.178	0.002(+)	0.027(+)	0.010(+)	0.010(+)	0.003(+)

Method		R15	Dim2	G2	Dim32	Customer	Unbalance	Spiral	Mockdata1	Mockdata2	Mockdata3
Wilcoxon	E-DE to K-means	0.010(+)	0.010(+)	0.002(+)	0.002(+)	0.084(-)	0.002(+)	0.084	0.042(+)	0.002(+)	0.026(+)
	E-DE to S-DE	0.007(+)	0.010(+)	0.014(+)	0.008(+)	0.028(-)	0.009(+)	0.032(+)	0.016(+)	0.002(+)	0.002(+)
	E-DE to BF-K-means	0.002(+)	0.005(+)	0.004(+)	0.063	0.002(+)	0.002(+)	0.002(+)	0.028(-)	0.010(+)	0.042(+)
	E-DE to ACDE	0.002(+)	0.002(+)	0.002(+)	0.014(+)	0.047(+)	0.004(+)	0.005(+)	0.024(+)	0.017(-)	0.002(+)
	E-DE to MANM	0.002(+)	0.002(+)	0.013(-)	0.002(+)	0.008(+)	0.044(+)	0.004(-)	0.010(+)	0.002(-)	0.010(+)
KW	E-DE to TGCA	0.010(+)	0.014(+)	0.002(-)	0.013(+)	0.019(+)	0.002(+)	0.002(+)	0.032(+)	0.037(+)	0.002(+)
	E-DE to Others	0.012(+)	0.002(+)	0.006(+)	0.013(+)	0.002(+)	0.368	0.008(+)	0.002(+)	0.004(+)	0.025(+)

TABLE V

P-VALUES OF THE WILCOXON RANK-SUM TESTS AND KW TESTS BETWEEN OUR METHOD AND THE OTHER SIX CLUSTERING ALGORITHMS OF THE I INDEX ON 19 DATASETS. (+) INDICATES THAT THE PERFORMANCE OF E-DE IS SIGNIFICANTLY BETTER THAN THE COMPARED ALGORITHM WITH A CONFIDENCE LEVEL OF 95%, WHILE (−) REPRESENTS THE OPPOSITE

Method		Iris	Cancer	Glass	Vowel	Zoo	Vote	Ecoli	Seed	Compound
Wilcoxon	E-DE to K-means	0.002(+)	0.002(+)	0.002(-)	0.005(-)	0.002(+)	0.002(+)	0.010(+)	0.010(+)	0.002(+)
	E-DE to S-DE	0.010(+)	0.010(+)	0.010(-)	0.010(-)	0.010(+)	0.009(+)	0.030(-)	0.010(+)	0.008(+)
	E-DE to BF-K-means	0.005(-)	0.019(+)	0.044(+)	0.026(+)	0.002(+)	0.002(+)	0.002(+)	0.008(-)	0.020(+)
	E-DE to ACDE	0.002(-)	0.018(+)	0.010(+)	0.012(-)	0.010(+)	0.002(+)	0.005(-)	0.067	0.008(-)
	E-DE to MANM	0.028(+)	0.023(+)	0.037(+)	0.019(+)	0.010(+)	0.010(+)	0.006(-)	0.006(+)	0.007(+)
KW	E-DE to TGCA	0.002(+)	0.002(+)	0.005(+)	0.002(+)	0.004(-)	0.004(+)	0.001(-)	0.002(+)	0.002(-)
	E-DE to Others	0.013	0.037(+)	0.036(+)	0.028	0.033(+)	0.030(+)	0.003(-)	0.004(+)	0.018(+)

Method		R15	Dim2	G2	Dim32	Customer	Unbalance	Spiral	Mockdata1	Mockdata2	Mockdata3
Wilcoxon	E-DE to K-means	0.010(+)	0.002(+)	0.002(+)	0.010(+)	0.002(+)	0.002(+)	0.010(+)	0.033(+)	0.002(+)	0.002(+)
	E-DE to S-DE	0.010(+)	0.002(+)	0.010(+)	0.010(+)	0.067	0.010(+)	0.010(+)	0.010(+)	0.065	0.003(+)
	E-DE to BF-K-means	0.007(+)	0.033(+)	0.037(-)	0.009(+)	0.010(+)	0.002(+)	0.003(+)	0.001(+)	0.003(+)	0.010(+)
	E-DE to ACDE	0.003(+)	0.004(+)	0.005(-)	0.009(+)	0.001(+)	0.003(+)	0.003(+)	0.019(+)	0.010(-)	0.010(+)
	E-DE to MANM	0.013(-)	0.006(+)	0.002(-)	0.002(+)	0.013(+)	0.002(+)	0.002(-)	0.007(+)	0.005(+)	0.010(+)
KW	E-DE to TGCA	0.002(+)	0.010(+)	0.002(-)	0.002(+)	0.010(+)	0.005(+)	0.008(+)	0.009(+)	0.002(+)	0.084(+)
	E-DE to Others	0.010(+)	0.044(+)	0.004	0.002(+)	0.030(+)	0.008(+)	0.001(+)	0.003(+)	0.034(+)	0.042(+)

from Table IV that, in terms of the DB index, the differences between the results of E-DE and the other algorithms are significant (the p -values are smaller than 0.05) on almost all datasets. Meanwhile, in most cases, E-DE outperforms the others. Similar results are also observed from Table V, showing significantly better performance of the proposed method in the comparisons using the I index. In addition, the above KW test results prove that, on most datasets, our E-DE algorithm is able to achieve statistically better results than the other six algorithms.

C. Analysis of Time Complexity

For the evolutionary clustering algorithms, the computational costs involve two parts: 1) the fitness evaluation of the clustering index and 2) the evolution operation. The complexity of fitness evaluation is identical for different evolutionary clustering algorithms. Particularly, the data assignment consumes $O(N \times K)$ time in calculating the distance between data points and cluster centroids, and the fitness evaluation of DB and I indices both takes an additional $O(K^2)$ time to calculate the distance between cluster centroids. Therefore, the fitness evaluation step consumes $O((N + K) \times K \times np)$ time in each iteration. The E-DE algorithm does not bring additional distance calculations when adopting the new mutation and crossover operators. The complexity of the evolution operation is the same as the classical DE algorithm, which is $O(K \times np)$ in each iteration. The evolution operations of ACDE and TGCA have the equal time consumption with our E-DE. Besides, MANM has an additional $O((K + np) \times np)$ time cost, which is caused by the auxiliary merging and niching procedures. From the time complexity analysis, we can see that the evolution costs of the evolutionary clustering algorithms are far

less than the evaluation costs, which can be ignored, especially when the scale of datasets is large. Besides, as indicated in the experimental comparisons, the proposed E-DE converges faster than the other algorithms, which consumes fewer fitness evaluation times to obtain promising results. From this point of view, E-DE is more computationally efficient than the others.

V. CONCLUSION

In this article, we proposed an E-DE to tackle the automatic data clustering problem. First, an elastic individual encoding scheme was introduced to adapt the number of clusters without external guidance. Unlike the zero-one switch encoding scheme used in traditional algorithms, this elastic scheme enables the population of DE to contain individuals that have different dimensions and evolves this parameter (i.e., the number of dimensions) inherently during the evolution. Second, a two-phase mutation operator was designed, which subtly utilizes the information of three individuals with random selection, deletion, and permutation operations to generate mutant vectors. Through this mutation, the individual learns sufficient information about the data clustering structure from the current population in order to adapt the cluster numbers and centroids. Third, we developed an exponential subspace crossover strategy that exchanges partial cluster layouts between parameter vectors. This exchange method not only avoids the cross-dimension learning error encountered by the previous algorithms, but it also adjusts the number of clusters in a flexible way. The experimental results proved that our E-DE algorithm outperformed the other clustering algorithms in terms of both internal and external clustering indices.

The E-DE algorithm showed significant improvements over the existing methods, but some limitations still remain. The proposed algorithm belongs to centroid-based (partitioning) clustering. Although such methods have been widely applied in practice, they endure the performance degradation in handling the nonconvex datasets. So one of our future works is to extend the E-DE framework by incorporating some density or connectivity-based clustering methods for nonconvex datasets. In addition, the evolutionary clustering methods may encounter deficiency when handling the extremely high-dimensional datasets like the deep feature-based image database. Another possible extension is to make some adjustments, such as incorporating a cooperative coevolution framework, to further enhance the scalability of the E-DE algorithm. Besides, due to the good performance of E-DE, it would be appealing to apply the algorithm to some practical applications, such as image analysis, document clustering, and gene detection. To summarize, we believe that our E-DE clustering method is flexible and robust to solve the automatic data clustering problem and there remains lots of possibilities for it to extend to other subjects or reality fields.

REFERENCES

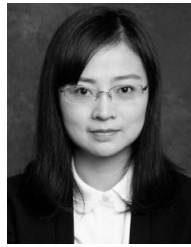
- [1] Q. Zhang, C. Zhu, L. T. Yang, Z. Chen, L. Zhao, and P. Li, "An incremental CFS algorithm for clustering large data in Industrial Internet of Things," *IEEE Trans. Ind. Inf.*, vol. 13, no. 3, pp. 1193–1201, Jun. 2017.
- [2] W. Bi, M. Cai, M. Liu, and G. Li, "A big data clustering algorithm for mitigating the risk of customer churn," *IEEE Trans. Ind. Inf.*, vol. 12, no. 3, pp. 1270–1281, Jun. 2016.
- [3] P. A. Jaskowiak, R. J. G. B. Campello, and I. G. Costa, "Proximity measures for clustering gene expression microarray data: A validation methodology and a comparative analysis," *IEEE/ACM Trans. Comput. Biol. Bioinf.*, vol. 10, no. 4, pp. 845–857, Jul./Aug. 2013.
- [4] P. Maji and S. Paul, "Rough-fuzzy clustering for grouping functionally similar genes from microarray data," *IEEE/ACM Trans. Comput. Biol. Bioinf.*, vol. 10, no. 2, pp. 286–299, Mar./Apr. 2013.
- [5] P. C. H. Ma, K. C. C. Chan, X. Yao, and D. K. Y. Chiu, "An evolutionary clustering algorithm for gene expression microarray data analysis," *IEEE Trans. Evol. Comput.*, vol. 10, no. 3, pp. 296–314, Jun. 2006.
- [6] R. C. Romero-Zalaz, C. Rubio-Escudero, J. P. Cobb, F. Herrera, Ó. Córdón, and I. Zwir, "A multiobjective evolutionary conceptual clustering methodology for gene annotation within structural databases: A case of study on the gene ontology database," *IEEE Trans. Evol. Comput.*, vol. 12, no. 6, pp. 679–701, Dec. 2008.
- [7] I. D. Gebreu, X. Alameda-Pineda, F. Forbes, and R. Horaud, "EM algorithms for weighted-data clustering with application to audio-visual scene analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 12, pp. 2402–2415, Dec. 2016.
- [8] N. Lu and H. Miao, "Clustering tree-structured data on manifold," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 1956–1968, Oct. 2016.
- [9] G. Folino, C. Pizzuti, and G. Spezzano, "Training distributed GP ensemble with a selective algorithm based on clustering and pruning for pattern classification," *IEEE Trans. Evol. Comput.*, vol. 12, no. 4, pp. 458–468, Aug. 2008.
- [10] J. A. Hartigan and M. A. Wong, "Algorithm as 136: A k -means clustering algorithm," *J. Roy. Stat. Soc. C (Appl. Stat.)*, vol. 28, no. 1, pp. 100–108, 1979.
- [11] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "SLIC superpixels compared to state-of-the-art superpixel methods," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2274–2282, Nov. 2012.
- [12] J. Qin, W. Fu, H. Gao, and W. X. Zheng, "Distributed k -means algorithm and fuzzy c -means algorithm for sensor networks based on multiagent consensus theory," *IEEE Trans. Cybern.*, vol. 47, no. 3, pp. 772–783, Mar. 2017.
- [13] D. S. Hochbaum, "Approximation algorithms for NP-hard problems," *ACM SIGACT News*, vol. 28, no. 2, pp. 40–52, 1997.
- [14] K. Krishna and M. N. Murty, "Genetic k -means algorithm," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 29, no. 3, pp. 433–439, Jun. 1999.
- [15] Y. Lu, S. Lu, F. Fotouhi, Y. Deng, and S. J. Brown, "FGKA: A fast genetic K -means clustering algorithm," in *Proc. ACM Symp. Appl. Comput.*, 2004, pp. 622–623.
- [16] P. Merz and A. Zell, "Clustering gene expression profiles with memetic algorithms," in *Proc. Int. Conf. Parallel Problem Solving Nat.*, 2002, pp. 811–820.
- [17] S. Das, A. Abraham, and A. Konar, "Automatic clustering using an improved differential evolution algorithm," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 38, no. 1, pp. 218–237, Jan. 2008.
- [18] C. Liu, A. Zhou, and G. Zhang, "Automatic clustering method based on evolutionary optimisation," *IET Comput. Vis.*, vol. 7, no. 4, pp. 258–271, Aug. 2013.
- [19] E. R. Hruschka and N. F. F. Ebecken, "A genetic algorithm for cluster analysis," *Intell. Data Anal.*, vol. 7, no. 1, pp. 15–25, 2003.
- [20] E. R. Hruschka, R. J. G. B. Campello, and L. N. De Castro, "Evolving clusters in gene-expression data," *Inf. Sci.*, vol. 176, no. 13, pp. 1898–1927, 2006.
- [21] S. J. Nanda and G. Panda, "A survey on nature inspired metaheuristic algorithms for partitioning clustering," *Swarm Evol. Comput.*, vol. 16, pp. 1–18, Jun. 2014.
- [22] A. Jos-García and W. Gmez-Flores, "Automatic clustering using nature-inspired metaheuristics: A survey," *Appl. Soft Comput.*, vol. 41, pp. 192–213, Apr. 2016.
- [23] E. Hancer and D. Karaboga, "A comprehensive survey of traditional, merge-split and evolutionary approaches proposed for determination of cluster number," *Swarm Evol. Comput.*, vol. 32, pp. 49–67, Feb. 2017.
- [24] W.-L. Xiang, N. Zhu, S.-F. Ma, X.-L. Meng, and M.-Q. An, "A dynamic shuffled differential evolution algorithm for data clustering," *Neurocomputing*, vol. 158, pp. 144–154, Jun. 2015.
- [25] R. Tinós, L. Zhao, F. Chicano, and D. Whitley, "NK hybrid genetic algorithm for clustering," *IEEE Trans. Evol. Comput.*, vol. 22, no. 5, pp. 748–761, Oct. 2018.
- [26] C. Guan, K. K. F. Yuen, and F. Coenen, "Particle swarm optimized density-based clustering and classification: Supervised and unsupervised learning approaches," *Swarm Evol. Comput.*, vol. 44, pp. 876–896, Feb. 2019.
- [27] A. Mukhopadhyay, U. Maulik, and S. Bandyopadhyay, "Multiobjective genetic algorithm-based fuzzy clustering of categorical attributes," *IEEE Trans. Evol. Comput.*, vol. 13, no. 5, pp. 991–1005, Oct. 2009.
- [28] I. Saha, D. Maity, and U. Maulik, "Categorical data analysis using multiobjective differential evolution based fuzzy clustering," in *Proc. IEEE Int. Conf. Adv. Comput. Commun. Informat. (ICACCI)*, 2013, pp. 2013–2017.
- [29] A. Gupta, S. Datta, and S. Das, "Fast automatic estimation of the number of clusters from the minimum inter-center distance for k -means clustering," *Pattern Recognit. Lett.*, vol. 116, pp. 72–79, Dec. 2018.
- [30] H.-H. Tam, S.-C. Ng, A. K. Lui, and M.-F. Leung, "Improved activation schema on automatic clustering using differential evolution algorithm," in *Proc. IEEE Congr. Evol. Comput. (CEC)*, 2017, pp. 1749–1756.
- [31] J. Arellano-Verdejo, E. Alba, and S. Godoy-Calderon, "Efficiently finding the optimum number of clusters in a dataset with a new hybrid differential evolution algorithm: DELA," *Soft Comput.*, vol. 20, no. 3, pp. 895–905, 2016.
- [32] J. Tvrđik and I. Křivý, "Hybrid differential evolution algorithm for optimal clustering," *Appl. Soft Comput.*, vol. 35, pp. 502–512, Oct. 2015.
- [33] U. Maulik and I. Saha, "Automatic fuzzy clustering using modified differential evolution for image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 9, pp. 3503–3510, Sep. 2010.
- [34] D. Chang and X. Zhang, "Dynamic niching genetic algorithm with data attraction for automatic clustering," *Tsinghua Sci. Technol.*, vol. 14, no. 6, pp. 718–724, Dec. 2009.
- [35] W. Sheng, S. Chen, M. Fairhurst, G. Xiao, and J. Mao, "Multilocal search and adaptive niching based memetic algorithm with a consensus criterion for data clustering," *IEEE Trans. Evol. Comput.*, vol. 18, no. 5, pp. 721–741, Oct. 2014.
- [36] W. Sheng, S. Chen, M. Sheng, G. Xiao, J. Mao, and Y. Zheng, "Adaptive multisubpopulation competition and multiniche crowding-based memetic algorithm for automatic data clustering," *IEEE Trans. Evol. Comput.*, vol. 20, no. 6, pp. 838–858, Dec. 2016.
- [37] H. He and Y. Tan, "A two-stage genetic algorithm for automatic clustering," *Neurocomputing*, vol. 81, pp. 49–59, Apr. 2012.
- [38] S. Chakraborty and S. Das, "Simultaneous variable weighting and determining the number of clusters—A weighted Gaussian means algorithm," *Stat. Probab. Lett.*, vol. 137, pp. 148–156, Jun. 2018.

- [39] S. Saha, A. Ekbal, and A. K. Alok, "Semi-supervised clustering using multiobjective optimization," in *Proc. IEEE 12th Int. Conf. Hybrid Intell. Syst. (HIS)*, 2012, pp. 360–365.
- [40] A. Mukhopadhyay, U. Maulik, and S. Bandyopadhyay, "Hybrid evolutionary multiobjective fuzzy c-medoids clustering of categorical data," in *Proc. IEEE Workshop Hybrid Intell. Models Appl. (HIMA)*, 2013, pp. 7–12.
- [41] Y. Wan, Y. Zhong, and A. Ma, "Fully automatic spectral-spatial fuzzy clustering using an adaptive multiobjective memetic algorithm for multispectral imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 4, pp. 2324–2340, Apr. 2019.
- [42] M. Garza-Fabre, J. Handl, and J. Knowles, "An improved and more scalable evolutionary approach to multiobjective clustering," *IEEE Trans. Evol. Comput.*, vol. 22, no. 4, pp. 515–535, Aug. 2018.
- [43] R. Xu, J. Xu, and D. C. Wunsch, "A comparison study of validity indices on swarm-intelligence-based clustering," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 42, no. 4, pp. 1243–1256, Aug. 2012.
- [44] O. Arbelaitz, I. Gurrutxaga, J. Muguerza, J. M. Pérez, and I. Perona, "An extensive comparative study of cluster validity indices," *Pattern Recognit.*, vol. 46, no. 1, pp. 243–256, 2013.
- [45] D. L. Davies and D. W. Bouldin, "A cluster separation measure," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. PAMI-1, no. 2, pp. 224–227, Apr. 1979.
- [46] S. Bandyopadhyay and U. Maulik, "Nonparametric genetic clustering: Comparison of validity indices," *IEEE Trans. Syst., Man, Cybern. C, Appl. Rev.*, vol. 31, no. 1, pp. 120–125, Feb. 2001.
- [47] L. Hubert and P. Arabie, "Comparing partitions," *J. Classification*, vol. 2, no. 1, pp. 193–218, Dec. 1985.
- [48] G. Wu, X. Shen, H. Li, H. Chen, A. Lin, and P. N. Suganthan, "Ensemble of differential evolution variants," *Inf. Sci.*, vol. 423, pp. 172–186, Jan. 2018.
- [49] X. Qiu, J.-X. Xu, Y. Xu, and K. C. Tan, "A new differential evolution algorithm for minimax optimization in robust design," *IEEE Trans. Cybern.*, vol. 48, no. 5, pp. 1355–1368, May 2018.
- [50] P. Ochoa, O. Castillo, J. Soria, and P. Cortes-Antonio, "Differential evolution algorithm using a dynamic crossover parameter with high-speed interval type 2 fuzzy system," in *Advances in Soft Computing*. Cham, Switzerland: Springer Int., 2018, pp. 369–378.
- [51] O. Castillo, F. Valdez, J. Soria, L. Amador-Angulo, P. Ochoa, and C. Peraza, "Comparative study in fuzzy controller optimization using bee colony, differential evolution, and harmony search algorithms," *Algorithms*, vol. 12, no. 1, p. 9, 2018.
- [52] (Aug. 31, 2019). *E-DE Source Code*. [Online]. Available: https://github.com/rich-lavender/E-DE_v1
- [53] S. Das and P. N. Suganthan, "Differential evolution: A survey of the state-of-the-art," *IEEE Trans. Evol. Comput.*, vol. 15, no. 1, pp. 4–31, Feb. 2011.
- [54] D. Dua and C. Graff. (2017). *UCI Machine Learning Repository*. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [55] C. T. Zahn, "Graph-theoretical methods for detecting and describing gestalt clusters," *IEEE Trans. Comput.*, vol. C-100, no. 1, pp. 68–86, Jan. 1971.
- [56] C. J. Veenman, M. J. T. Reinders, and E. Backer, "A maximum variance cluster algorithm," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 9, pp. 1273–1280, Sep. 2002.
- [57] I. Kärkkäinen and P. Fränti, "Gradual model generator for single-pass clustering," *Pattern Recognit.*, vol. 40, no. 3, pp. 784–795, 2007.
- [58] P. Fränti, R. Marescu-Istodor, and C. Zhong, "XNN graph," in *Proc. Joint IAPR Int. Workshops Stat. Techn. Pattern Recognit. Struct. Syn. Pattern Recognit.*, 2016, pp. 207–217.
- [59] P. Franti, O. Virtajoki, and V. Hautamaki, "Fast agglomerative clustering using a k -nearest neighbor graph," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 11, pp. 1875–1881, Nov. 2006.
- [60] M. Rezaei and P. Fränti, "Set matching measures for external cluster validity," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 8, pp. 2173–2186, Aug. 2016.
- [61] H. Chang and D.-Y. Yeung, "Robust path-based spectral clustering," *Pattern Recognit.*, vol. 41, no. 1, pp. 191–203, 2008.
- [62] M. N. Omidvar, X. Li, Y. Mei, and X. Yao, "Cooperative co-evolution with differential grouping for large scale optimization," *IEEE Trans. Evol. Comput.*, vol. 18, no. 3, pp. 378–393, Jun. 2014.
- [63] X. Y. Zhang, Y. J. Gong, Y. Lin, J. Zhang, S. Kwong, and J. Zhang, "Dynamic cooperative coevolution for large scale optimization," *IEEE Trans. Evol. Comput.*, to be published.



Jun-Xian Chen received the bachelor's degree from Jinan University, Guangzhou, China, in 2017. He is currently pursuing the master's degree with the School of Computer Science and Engineering, South China University of Technology, Guangzhou.

His current research interests include evolutionary computation, swarm intelligence, and data mining.



Yue-Jiao Gong (S'10–M'15–SM'19) received the B.S. and Ph.D. degrees in computer science from Sun Yat-sen University, Guangzhou, China, in 2010 and 2014, respectively.

She is currently a Full Professor with the School of Computer Science and Engineering, South China University of Technology, Guangzhou. Her current research interests include evolutionary computation, swarm intelligence, and their applications to intelligent transportation and smart city scheduling. She has published over 80 papers, including over 30

IEEE TRANSACTIONS papers in the above areas.



Wei-Neng Chen (S'07–M'12–SM'17) received the bachelor's and Ph.D. degrees in computer science from Sun Yat-sen University, Guangzhou, China, in 2006 and 2012, respectively.

He is currently a Professor with the School of Computer Science and Engineering, South China University of Technology, Guangzhou. His current research interests include swarm intelligence algorithms and their applications on cloud computing, operations research, and software engineering.

Prof. Chen was a recipient of the IEEE Computational Intelligence Society Outstanding Dissertation Award in 2016 for his Ph.D. thesis and the National Science Fund for Excellent Young Scholars in 2016.



Mengting Li received the M.S. and Ph.D. degrees in computer science from Sun Yat-sen University, Guangzhou, China.

She is currently a Postdoctoral Research Fellow with the School of Data and Computer Science, Sun Yat-sen University. Her current research interests include data mining and machine learning.



Jun Zhang (M'02–SM'08–F'17) received the Ph.D. in electrical engineering degree from the City University of Hong Kong, Hong Kong, in 2002.

He is currently a Visiting Scholar with Victoria University, Melbourne, VIC, Australia. His current research interests include computational intelligence, cloud computing, high-performance computing, and operations research.

Dr. Zhang is currently an Associate Editor of the IEEE TRANSACTIONS ON EVOLUTIONARY COMPUTATION, the IEEE TRANSACTIONS ON CYBERNETICS, and the IEEE TRANSACTIONS ON INDUSTRIAL ELECTRONICS.