

# Toward Multidiversified Ensemble Clustering of High-Dimensional Data: From Subspaces to Metrics and Beyond

Dong Huang<sup>ID</sup>, *Member, IEEE*, Chang-Dong Wang<sup>ID</sup>, *Member, IEEE*, Jian-Huang Lai<sup>ID</sup>, *Senior Member, IEEE*, and Chee-Keong Kwoh, *Senior Member, IEEE*

**Abstract**—The rapid emergence of high-dimensional data in various areas has brought new challenges to current ensemble clustering research. To deal with the curse of dimensionality, recently considerable efforts in ensemble clustering have been made by means of different subspace-based techniques. However, besides the emphasis on subspaces, rather limited attention has been paid to the potential diversity in similarity/dissimilarity metrics. It remains a surprisingly open problem in ensemble clustering how to create and aggregate a large population of diversified metrics, and furthermore, how to jointly investigate the multilevel diversity in the large populations of metrics, subspaces, and clusters in a unified framework. To tackle this problem, this article proposes a novel multidiversified ensemble clustering approach. In particular, we create a large number of diversified metrics by randomizing a scaled exponential similarity kernel, which are then coupled with random subspaces to form a large set of metric-subspace pairs. Based on the similarity matrices derived from these metric-subspace pairs, an ensemble of diversified base clusterings can be thereby constructed. Furthermore, an entropy-based criterion is utilized to explore the cluster wise diversity in ensembles, based on which three specific ensemble clustering algorithms are presented by incorporating three types of consensus functions. Extensive experiments are conducted on 30 high-dimensional datasets, including 18 cancer gene expression datasets and 12 image/speech datasets, which demonstrate the superiority of our algorithms over the state of the art. The source code is available at <https://github.com/huangdonghere/MDEC>.

**Index Terms**—Consensus clustering, diversified metrics, ensemble clustering, high-dimensional data, random subspaces.

## I. INTRODUCTION

THE LAST decade has witnessed significant progress in the development of the ensemble clustering technique [1]–[22], which is typically featured by its ability of combining multiple base clusterings into a probably better and more robust consensus clustering and has recently shown promising advantages in discovering clusters of arbitrary shapes, dealing with noisy data, coping with data from multiple sources, and producing robust clustering results [7].

In recent years, with high-dimensional data widely appearing in various areas, new challenges have been brought to the conventional ensemble clustering algorithms, which, however, often lack the ability to well address the high-dimensional issues. As is called the curse of dimensionality, it is highly desired but very difficult to find the inherent cluster structure hidden in the huge dimensions, especially when it is frequently coupled with quite low sample size. Recently some efforts have been devoted to ensemble clustering of high-dimensional data, which typically exploit some subspace-based (or feature-based) techniques (such as random subspace sampling [13], [23]–[26], stratified subspace sampling [6], and subspace projection [27]) to explore the diversity in high dimensionality. Inherently, these subspace-based techniques select or linearly combine data features into different subsets (i.e., subspaces) by a variety of strategies to seek more perspectives for finding cluster structures.

Besides the issue of subspaces (or features), the choice of similarity/dissimilarity metrics is another crucial factor in dealing with high-dimensional data [28], [29]. The existing ensemble clustering methods generally adopt one or a few preselected metrics, which are often selected explicit or implicitly based on the expert's knowledge or some prior assumptions. However, few, if not none, of them have considered the potentially huge opportunities hidden in randomized metric spaces. On the one hand, it is very difficult to select or learn an optimal metric for a given dataset without human supervision or implicit assumptions. On the other hand, with different metrics capable of reflecting different perspectives on data, the joint use of a large number of randomized/diversified metrics may reveal more valuable information hidden in high dimensionality. However, it is surprisingly still an open problem in ensemble clustering how to produce and aggregate a large number of diversified metrics to enhance the consensus performance. Furthermore,

Manuscript received 22 September 2020; revised 30 November 2020; accepted 3 January 2021. Date of publication 7 May 2021; date of current version 17 October 2022. This work was supported in part by NSFC under Grant 61976097, Grant 61876193, and Grant 61876104; and in part by A\*STAR-NTU-SUTD AI Partnership under Grant RGANS1905. This article was recommended by Associate Editor X. Wang. (*Corresponding author: Chang-Dong Wang.*)

Dong Huang is with the College of Mathematics and Informatics, South China Agricultural University, Guangzhou 510642, China, and also with Pazhou Lab, Guangzhou 510335, China, and also with the Guangzhou Key Laboratory of Smart Agriculture, Guangzhou 510006, China (e-mail: [huangdonghere@gmail.com](mailto:huangdonghere@gmail.com)).

Chang-Dong Wang and Jian-Huang Lai are with the School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou 510006, China, also with the Guangdong Key Laboratory of Information Security Technology, Guangzhou 510006, China, and also with the Key Laboratory of Machine Intelligence and Advanced Computing, Ministry of Education, Beijing 510006, China (e-mail: [changdongwang@hotmail.com](mailto:changdongwang@hotmail.com); [stsljh@mail.sysu.edu.cn](mailto:stsljh@mail.sysu.edu.cn)).

Chee-Keong Kwoh is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore (e-mail: [aseckkwoh@ntu.edu.sg](mailto:aseckkwoh@ntu.edu.sg)).

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.2021.3049633>.

Digital Object Identifier 10.1109/TCYB.2021.3049633

starting from the metric diversification problem, another crucial challenge arises as to how to jointly exploit multiple levels of diversity in the large number of metrics, subspaces, and clusters, in a unified ensemble clustering framework.

To tackle the above-mentioned problem, this article presents a novel ensemble clustering approach called multidiversified ensemble clustering (MDEC) by jointly exploiting large populations of diversified metrics, random subspaces, and weighted clusters. Particularly, a scaled exponential similarity kernel is utilized as the seed kernel, which has advantage in parameter flexibility and neighborhood adaptivity and is randomized to breed a large set of diversified metrics. Then, the set of diversified metrics is coupled with random subspaces to form a large number of metric-subspace pairs, which contribute to the jointly randomized ensemble generation process where the set of diversified base clusterings can thereby be produced with the help of spectral clustering (SC). Furthermore, to exploit the clusterwise diversity in the multiple base clusterings, an entropy-based cluster validity strategy is adopted to evaluate and weight each base cluster by considering the distribution of clusters in the entire ensemble, based on which three specific ensemble clustering algorithms are, therefore, proposed by incorporating three types of consensus functions. Extensive experiments are conducted on 30 real-world high-dimensional datasets, including 18 cancer gene expression datasets and 12 image/speech datasets. The experimental results have shown the superiority of our approach against the state-of-the-art ensemble clustering approaches for clustering high-dimensional data.

For clarity, the main contributions of this work are summarized as follows.

- 1) This article for the first time, to our knowledge, shows that the joint use of a large population of randomly diversified metrics can significantly benefit the ensemble clustering of high-dimensional data in a purely unsupervised manner.
- 2) A new metric diversification strategy is proposed by randomizing the scaled exponential similarity kernel with both parameter flexibility and neighborhood adaptivity considered, which is further coupled with random subspace sampling for the jointly randomized generation of base clusterings.
- 3) A new ensemble clustering framework called MDEC is presented, which is capable of simultaneously exploiting large populations of diversified metrics, random subspaces, and weighted clusters in a unified model for high-dimensional data.
- 4) Three specific algorithms are designed under the proposed MDEC framework by incorporating three types of consensus functions. Experiments conducted on a variety of real-world high-dimensional datasets have confirmed the advantages of the proposed algorithms over the state of the art.

The remainder of this article is organized as follows. The related work is reviewed in Section II. The proposed ensemble clustering framework is described in Section III. The experimental results are reported in Section IV. Finally, this article is concluded in Section V.

## II. RELATED WORK

Due to its ability of combining multiple base clusterings into a probably better and more robust consensus clustering, the ensemble clustering technique has been receiving increasing attention in recent years. Many ensemble clustering algorithms have been developed from different technical perspectives [1]–[5], [7]–[12], [15], [30]–[42], which can be classified into three main categories, namely, the pairwise co-occurrence-based methods, the graph partitioning-based methods, and the median partition-based methods.

The pairwise co-occurrence-based methods [?], [30] typically construct a co-association matrix by considering the frequency that two data samples occur in the same cluster among the multiple base clusterings. The co-association matrix is then used as the similarity matrix for the data samples, upon which some clustering algorithms can thereby be performed to obtain the final clustering result. Fred and Jain [30] first introduced the concept of the co-association matrix and proposed the evidence accumulation clustering (EAC) method, which applied a hierarchical agglomerative clustering algorithm [43] on the co-association matrix to build the consensus clustering. To extend the EAC method, Wang *et al.* [37] took the cluster sizes into consideration and proposed the probability accumulation method. Yi *et al.* [38] dealt with the uncertain entries in the co-association matrix by first labeling them as unobserved, and then recovering the unobserved entries by the matrix completion technique. Liu *et al.* [11] proved that the SC of the co-association matrix is equivalent to a weighted version of  $K$ -means, and proposed the spectral ensemble clustering (SEC) method to effectively and efficiently obtain the consensus result.

The graph partitioning-based methods [8], [39], [40] generally construct a graph model for the ensemble of multiple base clusterings, and then partition the graph into several disjoint subsets to obtain the final clustering result. Strehl and Ghosh [39] solved the ensemble clustering problem by using three graph partitioning-based algorithms, namely, the cluster-based similarity partitioning algorithm (CSPA), hypergraph partitioning algorithm (HGPA), and metaclustering algorithm (MCLA). Fern and Brodley [40] formulated a bipartite graph (BG) model by treating both clusters and data samples as nodes, and partitioned the graph by the METIS algorithm [44] to obtain the consensus result. Huang *et al.* [8] dealt with the ensemble clustering problem by sparse graph representation and random walk trajectory analysis, and presented the probability trajectory-based graph partitioning (PTGP) method.

The median partition-based methods [9], [41], [42] typically formulate the ensemble clustering problem into an optimization problem, which aims to find the median partition such that the similarity between the base partitions (i.e., base clusterings) and the median partition is maximized. The median partition problem is NP-hard [41]. To find an approximate solution, Topchy *et al.* [41] casted the median partition problem into a maximum-likelihood problem and solved it by the EM algorithm. Franek and Jiang [42] reduced the ensemble clustering problem to an Euclidean median problem and solved it by the Weiszfeld algorithm [45]. Huang *et al.* [9] formulated the ensemble clustering problem

into a binary linear programming problem and obtained an approximate solution based on the factor graph model and the max-product belief propagation [46].

Although in recent years, significant advances have been made in the research of ensemble clustering [1]–[5], [7]–[12], [15], [30]–[42], yet the existing methods are mostly devised for general-purpose scenarios and lack the desirable ability to appropriately address the clustering problem of high-dimensional data. More recently, some efforts have been made to deal with the curse of dimensionality, where subspace-based (or feature-based) techniques are often exploited. Jing *et al.* [6] adopted stratified feature sampling to generate a set of subspaces, which are further incorporated into several ensemble clustering algorithms to build the consensus clustering for high-dimensional data. Yu *et al.* [24] proposed a novel subspace-based ensemble clustering framework called AP<sup>2</sup>CE, which integrates random subspaces, affinity propagation, normalized cut, and five candidate distance metrics. Furthermore, Yu *et al.* [25] proposed a semisupervised subspace-based ensemble clustering framework by incorporating random subspaces, constraint propagation, incremental ensemble selection, and normalized cut into the framework. Fern and Brodley [27] exploited random subspace projection to build a set of subspaces, which, in fact, are obtained by (randomly) linear combination of features (or feature sets). Chu *et al.* [47] used three clustering results (generated by three different clustering algorithms) to guide the feature learning process. These methods [6], [24], [25], [27], [47] typically exploit the diversity in high dimensionality by various subspace-based or feature learning-based techniques, but few of them have fully considered the potentially huge diversity in metric spaces. The existing methods [6], [24], [25], [27], [47] generally use one or a few preselected similarity/dissimilarity metrics, which are selected implicitly based on the expert's knowledge or some prior assumptions. Although the method in [24] proposed to randomly select a metric out of the five candidate metrics at each time, yet it still failed to go beyond a few metrics to explore the huge potential hidden in a large number of diversified metrics, which may play a crucial role in clustering high-dimensional data. The key challenge here lies in how to create such a large number of highly diversified metrics, and further how to jointly exploit the diversity in the large number of metrics, together with subspacewise diversity and clusterwise diversity, to achieve a unified ensemble clustering framework for high-dimensional data.

### III. PROPOSED FRAMEWORK

This section describes the overall algorithm of the proposed ensemble clustering framework. A brief overview is provided in Section III-A. The metric diversification process is presented in Section III-B. The jointly randomized ensemble generation is introduced in Section III-C. Finally, the consensus functions are given in Section III-D.

#### A. Brief Overview

In this article, we propose a novel MDEC framework (as shown in Fig. 1). First, we create a large number of diversified

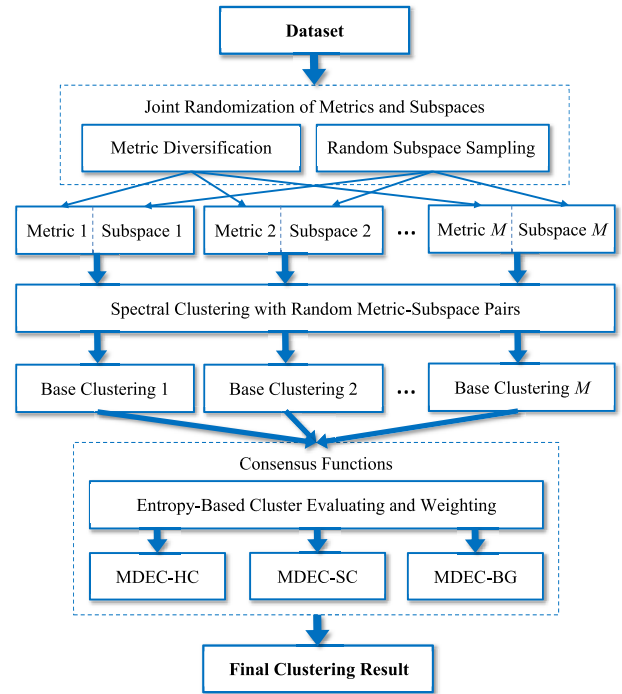


Fig. 1. Flow diagram of the proposed MDEC approach.

metrics by randomizing a scaled exponential similarity kernel, and combine the diversified metrics with the random subspaces to form a large set of random metric-subspace pairs. Second, with each random metric-subspace pair, we construct a similarity matrix for the data samples. Then, SC is performed on these similarity matrices derived from metric-subspace pairs to obtain an ensemble of base clusterings. Third, to exploit the clusterwise diversity in the ensemble of multiple base clusterings, we adopt an entropy-based criterion to evaluate and weight the clusters by considering the distribution of cluster labels in the entire ensemble. With the weighted clusters, a locally weighted co-association matrix is further constructed to serve as a summary of the ensemble. Finally, to obtain the consensus clustering, three types of consensus functions are presented based on hierarchical clustering (HC), SC, and BG model, respectively, leading to three specific ensemble clustering algorithms called MDEC-HC, MDEC-SC, and MDEC-BG, respectively. It is noteworthy that we simultaneously incorporate three levels of diversity, that is, metricwise diversity, subspacewise diversity, and clusterwise diversity, in a unified framework, which has shown its advantages in dealing with high-dimensional data when compared to the state-of-the-art ensemble clustering algorithms. In the following sections, we will further introduce each component of the proposed framework.

#### B. Diversification of Metrics

The choice of similarity/dissimilarity metrics plays a crucial role in the field of data mining and machine learning [48]–[53], especially for high-dimensional data analysis where the high dimensionality further complicates the use of



metrics [49], [53]. Some supervised or semisupervised learning techniques have been developed to learn a suitable metric for some specific applications [48]–[53], but in unsupervised scenarios, it is still very difficult to discover one or a few proper metrics given a task without prior knowledge.

Instead of relying on one or a few selected or learned metrics, this article proposes to jointly use a large number of randomly diversified metrics, which will be further coupled with a large number of random subspaces, in a unified ensemble clustering framework for high-dimensional data. Here, two subproblems are involved, namely, how to create a large number of diversified metrics and how to collectively exploit them in ensemble clustering.

To create the diversified metrics, we take advantage of the kernel trick with randomization incorporated. The kernel similarity metrics have been proved to be a powerful tool for clustering complex data [48], [51], which, however, suffer from the difficulties in selecting proper kernel parameters. The kernel parameters can be learned by some metric learning techniques [48], [51] with supervision or semisupervision. But without human supervision, it is often extremely difficult to decide proper kernel parameters. This is a critical disadvantage of kernel methods for conventional (unsupervised) applications, which, nevertheless, just becomes an important advantage in our situation where what is highly desired is not the selection of a good kernel similarity metric but the flexibility to create a large number of diversified ones.

Specifically, in this article, we adopt a scaled exponential similarity kernel [54] as the seed kernel, which has advantage in parameter flexibility and neighborhood adaptivity and is then randomized to breed a large population of diversified metrics. Given a set of  $N$  data samples  $\mathcal{X} = \{x_1, \dots, x_N\}$ , where  $x_i \in \mathbb{R}^D$  is the  $i$ th sample and  $D$  is the number of features. The kernel similarity between samples  $x_i$  and  $x_j$  is defined as:

$$\phi^{\mathcal{X}}(x_i, x_j) = \begin{cases} \exp\left(-\frac{d(x_i, x_j)}{\mu \varepsilon_{ij}}\right), & \text{if } x_i \in N_k(x_j) \\ & \text{or } x_j \in N_k(x_i) \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where  $\mu$  is a hyperparameter,  $\varepsilon_{ij}$  is a scaling term,  $d(x_i, x_j)$  is the Euclidean distance between  $x_i$  and  $x_j$ , and  $N_k(x_i)$  denotes the set of the  $k$  nearest neighbors of  $x_i$ . The average distance between  $x_i$  and its  $k$  nearest neighbors can be computed as

$$\rho_k(x_i) = \frac{1}{k} \cdot \sum_{x_p \in N_k(x_i)} d(x_i, x_p). \quad (2)$$

Then, as suggested in [54], to simultaneously take into consideration the neighborhood of  $x_i$ , the neighborhood of  $x_j$ , and their distance, the scaling term  $\varepsilon_{ij}$  is defined as the average of  $\rho_k(x_i)$ ,  $\rho_k(x_j)$ , and  $d(x_i, x_j)$ . That is

$$\varepsilon_{ij} = \frac{\rho_k(x_i) + \rho_k(x_j) + d(x_i, x_j)}{3}. \quad (3)$$

The above kernel is a variant of the Gaussian kernel. It has two free parameters, that is, the hyperparameter  $\mu$  and the number of nearest neighbors  $k$ . The motivation to adopt the SES kernel as the seed kernel in our approach is twofold. First, with the influence of the scaling term where the

$k$ -nearest-neighbors' information is incorporated, the kernel has the adaptivity to the neighborhood structure among data samples. Moreover, with each  $k$  value corresponding to a specific neighborhood size, by randomizing the parameter  $k$ , multiscale neighborhood information can be explored to enhance the diversity. Second, the two free parameters  $k$  and  $\mu$  in the kernel provide high flexibility for adjusting the influence of the kernel and can contribute to the desired diversity of the generated metrics by randomly perturbing the two parameters.

Specifically, we randomly select the two parameters  $\mu \in [\mu_{\min}, \mu_{\max}]$  and  $k \in [k_{\min}, k_{\max}]$ , respectively, as follows:

$$\mu = \mu_{\min} + \sigma_1(\mu_{\max} - \mu_{\min}) \quad (4)$$

$$k = k_{\min} + \lfloor \sigma_2(k_{\max} - k_{\min}) \rfloor \quad (5)$$

where  $\sigma_1 \in [0, 1]$  and  $\sigma_2 \in [0, 1]$  are two uniform random variables, and  $\lfloor \cdot \rfloor$  outputs the floor of a real number.

Note that our objective is not to find a good pair of parameters  $\mu$  and  $k$  but to randomize them to yield a large population of diversified metrics. The parameters  $\mu$  and  $k$  are suggested to be randomly selected in a wide range to enhance the diversity.

By performing the random selection  $M$  times, a set of  $M$  pairs of  $\mu$  and  $k$  are obtained, which correspond to  $M$  randomized kernel similarity metrics for the dataset  $\mathcal{X}$ , denoted as

$$\phi_{\mu_1, k_1}^{\mathcal{X}}(\cdot, \cdot), \phi_{\mu_2, k_2}^{\mathcal{X}}(\cdot, \cdot), \dots, \phi_{\mu_M, k_M}^{\mathcal{X}}(\cdot, \cdot) \quad (6)$$

where  $\mu_i$  and  $k_i$  are the  $i$ th pair of randomized parameters.

### C. Ensemble Generation by Joint Randomization

In high-dimensional data, the cluster structures are often hidden in different low-dimensional subspaces [49], [53]. Besides the subspaces, different metrics may also provide complementary information for investigating the high-dimensional space. In this section, with the set of diversified metrics generated, we proceed to couple the large number of diversified metrics with random subspaces to explore the rich information in various subspaces in high-dimensional data.

Specifically, let  $\mathcal{F} = \{f_1, \dots, f_D\}$  be the set of features in the dataset  $\mathcal{X}$ , where  $f_i$  denotes the  $i$ th feature. A random subspace is a set of a certain number of features that are randomly sampled from the original feature set. The cluster structure of high-dimensional data may be hidden in different feature subspaces as well as in different metric spaces. In this article, we propose to jointly exploit large populations of diversified metrics and random subspaces. Specifically, we perform random subspace sampling  $M$  times to obtain  $M$  random subspaces, denoted as  $\mathcal{F}_1, \dots, \mathcal{F}_M$ , which lead to  $M$  component datasets, denoted as  $\mathcal{X}_1, \dots, \mathcal{X}_M$ . Note that each component dataset  $\mathcal{X}_i$  has the same number of data samples as the original dataset  $\mathcal{X}$ , but its feature set  $\mathcal{F}_i$  only consists of  $d \leq D$  attributes that are randomly sampled from  $\mathcal{F}$  with a sampling ratio  $\tau \in (0, 1]$ . Obviously, if  $\tau = 1$ , then it means every subspace is in fact the original feature space, that is, no subsampling actually happens. Here, with the random subspaces generated, we can couple each of them with a randomly diversified metric (as described in Section III-B), and thus obtain

$M$  random metric-subspace pairs, denoted as

$$\varphi_{\mu_1, k_1}^{\mathcal{X}_1}(\cdot, \cdot), \varphi_{\mu_2, k_2}^{\mathcal{X}_2}(\cdot, \cdot), \dots, \varphi_{\mu_M, k_M}^{\mathcal{X}_M}(\cdot, \cdot). \quad (7)$$

In terms of the  $m$ th metric-subspace pair  $\varphi_{\mu_m, k_m}^{\mathcal{X}_m}(\cdot, \cdot)$ , the similarity between samples  $x_i$  and  $x_j$  is computed by first mapping  $x_i$  and  $x_j$  onto the subspace associated with the component dataset  $\mathcal{X}_m$  and then computing their kernel similarity with the randomly selected parameters  $\mu_m$  and  $k_m$ . Thus, we can obtain  $M$  similarity matrices in terms of the  $M$  metric-subspace pairs as follows:

$$\mathcal{S} = \{S^{(1)}, S^{(2)}, \dots, S^{(M)}\} \quad (8)$$

where the  $m$ th similarity matrix [i.e.,  $S^{(m)}$ ] is constructed with respect to the  $m$ th metric-subspace pair  $\varphi_{\mu_m, k_m}^{\mathcal{X}_m}(\cdot, \cdot)$ , denoted as

$$S^{(m)} = \{S_{ij}^{(m)}\}_{N \times N} \quad (9)$$

where

$$S_{ij}^{(m)} = \varphi_{\mu_m, k_m}^{\mathcal{X}_m}(x_i, x_j) \quad (10)$$

is the  $(i, j)$ th entry in  $S^{(m)}$ . Obviously, according to the definition of the kernel, it holds that  $S_{ij}^{(m)} \in (0, 1]$  for any  $x_i, x_j \in \mathcal{X}$ . If samples  $x_i$  and  $x_j$  have the same feature values in the subspace associated with  $\mathcal{X}_m$ , then their similarity  $S_{ij}^{(m)}$  reaches its maximum 1.

Having constructed  $M$  similarity matrices with diversified metric-subspace pairs, we then exploit the SC [55] to construct the ensemble of base clusterings. Specifically, for the  $m$ th similarity matrix  $S^{(m)}$ , we treat each data sample as a graph node and build a similarity graph as follows:

$$G^{(m)} = (V, E^{(m)}) \quad (11)$$

where  $V = \mathcal{X}$  is the node set, and  $E^{(m)}$  is the edge set. The edge weights are decided by the similarity matrix  $S^{(m)}$ , that is, for any  $x_i, x_j \in \mathcal{X}$ , we have  $E_{ij}^{(m)} = S_{ij}^{(m)}$ . Let  $K_m$  denote the number of clusters in the  $m$ -base clustering. The objective of SC is to partition the graph  $G^{(m)}$  into  $K_m$  disjoint subsets. To this end, we construct the normalized graph Laplacian  $L_m$  as follows:

$$L_m = I - D_m^{-1/2} S^{(m)} D_m^{-1/2} \quad (12)$$

where the degree matrix  $D_m \in \mathbb{R}^{N \times N}$  is a diagonal matrix with its  $(i, i)$ th entry defined as the sum of the  $i$ th row of  $S^{(m)}$ . The eigenvectors corresponding to the first  $K_m$  eigenvalues of  $L_m$  are computed and then stacked to form a new matrix  $U_m \in \mathbb{R}^{N \times K_m}$ , where the  $i$ th column of  $U_m$  is the eigenvector corresponding to the  $i$ th eigenvalue of  $L_m$ . Thereafter, the matrix  $T_m \in \mathbb{R}^{N \times K_m}$  can be obtained from  $U_m$  by normalizing the rows to norm 1. By treating each row of  $T_m$  as a data point in  $\mathbb{R}^{K_m}$ , we can cluster the rows into  $K_m$  clusters by  $K$ -means discretization, and thereby obtain the  $m$ th base clustering based on the similarity matrix  $S^{(m)}$ . Formally, the  $m$ th base clustering is denoted as

$$\pi^m = \{C_1^m, C_2^m, \dots, C_{K_m}^m\} \quad (13)$$

where  $C_i^m$  is the  $i$ th cluster in  $\pi^m$ . It is obvious that the  $K_m$  clusters in a base clustering cover the entire

dataset, that is,  $\bigcup_{i=1}^{K_m} C_i^m = \mathcal{X}$ , and two clusters in the same base clustering will not overlap with each other, that is  $\forall i \neq j, C_i^m \cap C_j^m = \emptyset$ .

Finally, based on the  $M$  diversified similarity matrices, we can construct an ensemble of  $M$  base clusterings, denoted as

$$\Pi = \{\pi^1, \pi^2, \dots, \pi^M\} \quad (14)$$

where  $\pi^m$  is the  $m$ th base clustering in the ensemble  $\Pi$ .

#### D. Consensus Functions

This article aims to integrate multiple levels of diversity in a unified ensemble clustering framework for high-dimensional data. With the metricwise and subspacewise diversity exploited in the jointly randomized ensemble generation, in this section, we proceed to explore the clusterwise diversity in ensembles and incorporate three types of consensus functions to combine the multiple base clusterings into the final consensus clustering result.

With each base clustering consisting of a certain number of clusters, the entire ensemble can also be viewed as a large set of clusters from different base clusterings. To exploit the different reliability of different clusters and incorporate the clusterwise diversity in the consensus function, here, we adopt a local weighting strategy [15] to evaluate and weight the base clusters by jointly considering the distribution of cluster labels in the entire ensemble using an entropic criterion. Formally, we denote the ensemble of clusters as

$$\mathcal{C} = \{C_1, C_2, \dots, C_{N_c}\} \quad (15)$$

where  $C_i$  is the  $i$ th cluster and  $N_c$  is the total number of clusters in the ensemble  $\Pi$ . Note that  $N_c = \sum_{m=1}^M K_m$ .

Each cluster is a set of data samples. To estimate the uncertainty of different clusters, the concept of entropy is utilized here [15]. Given a cluster  $C_i \in \mathcal{C}$  and a base clustering  $\pi^m \in \Pi$ , the uncertainty (or entropy) of  $C_i$  with respect to  $\pi^m$  can be computed as

$$H^{\pi^m}(C_i) = - \sum_{j=1}^{K_m} p(C_i, C_j^m) \log_2 p(C_i, C_j^m) \quad (16)$$

where

$$p(C_i, C_j^m) = \frac{|C_i \cap C_j^m|}{|C_i|} \quad (17)$$

is the proportion of data samples in  $C_i$  that also appear in  $C_j^m$ . It is obvious that  $p(C_i, C_j^m) \in [0, 1]$ , which leads to  $H^{\pi^m}(C_i) \in [0, +\infty)$ . If and only if all data samples in  $C_i$  also occur in the same cluster in  $\pi^m$ , the uncertainty of  $C_i$  with respect to  $\pi^m$  reaches its minimum 0.

With the general assumption that the set of base clusterings is independent of each other, we can obtain the uncertainty (or entropy) of  $C_i$  with respect to the entire ensemble  $\Pi$  as follows:

$$H^\Pi(C_i) = \sum_{m=1}^M H^{\pi^m}(C_i). \quad (18)$$

Intuitively, higher uncertainty indicates lower reliability for a cluster, which implies that the ensemble of base clusterings tend to disagree with the cluster and accordingly a smaller weight can be associated with it [15]. In particular, we proceed to compute a reliability index from the above-mentioned uncertainty measure, and exploit it as a cluster weighting term in our consensus function. The experimental analysis about the efficacy of the cluster weighting term will also be provided in Section IV-E. Specifically, the ensemble-driven cluster index (ECI) is computed as an indication for the reliability of each cluster in the ensemble, which is defined as follows:

$$\text{ECI}(C_i) = \exp\left(-\frac{H^\Pi(C_i)}{M}\right). \quad (19)$$

Obviously, for any  $\pi^m \in \Pi$ , it holds that  $H^{\pi^m}(C_i) \in [0, +\infty)$ , then we have  $H^\Pi(C_i) \in [0, +\infty)$  and thereby  $\text{ECI}(C_i) \in (0, 1]$ . Note that a larger value of ECI is associated with a cluster of lower uncertainty (i.e., higher reliability). If and only if the data samples in  $C_i$  appear in the same cluster in all of the  $M$  base clusterings (i.e., all base clusterings agree that the data samples in  $C_i$  should belong to the same cluster), the uncertainty of  $C_i$  with respect to  $\Pi$  reaches its minimum 0 and the ECI of  $C_i$  reaches its maximum 1.

The ECI measure serves as a reliability index for different clusters in the ensemble  $\Pi$ . By using ECI as a cluster-weighting term, the locally weighted co-association matrix can be obtained as follows:

$$A = \{A_{ij}\}_{N \times N} \quad (20)$$

$$A_{ij} = \frac{1}{M} \cdot \sum_{m=1}^M w_i^m \cdot \delta_{ij}^m \quad (21)$$

$$w_i^m = \text{ECI}(\text{Cls}^m(x_i)) \quad (22)$$

$$\delta_{ij}^m = \begin{cases} 1, & \text{if } \text{Cls}^m(x_i) = \text{Cls}^m(x_j) \\ 0, & \text{otherwise} \end{cases} \quad (23)$$

where  $\text{Cls}^m(x_i)$  denotes the cluster in  $\pi^m$  that  $x_i$  belongs to. Note that  $w_i^m$  is the cluster weighting term that weights each cluster according to its ECI value, while  $\delta_{ij}^m$  is the pairwise co-occurrence term that indicates whether two samples occur in the same cluster in a base clustering  $\pi^m$ .

With the clusterwise diversity explored, we further exploit three types of consensus functions to obtain the consensus clustering result, and thereby present three specific ensemble clustering algorithms, called MDEC-HC, MDEC-SC, and MDEC-BG, respectively, under our MDEC framework.

1) *MDEC-HC*: In MDEC-HC, an HC-based consensus function is incorporated, which constructs a dendrogram (i.e., a hierarchical cluster tree) by iterative region merging. Specifically, given a dataset  $\mathcal{X}$ , the  $N$  data samples in  $\mathcal{X}$  are treated as the set of initial regions, denoted as

$$\mathcal{R}^{(0)} = \{R_1^{(0)}, \dots, R_{|\mathcal{R}^{(0)}|}^{(0)}\} \quad (24)$$

where  $R_i^{(0)} = \{x_i\}$  is the  $i$ th initial region and  $|\mathcal{R}^{(0)}| = N$  is the number of initial regions. Here, the locally weighted co-association matrix  $A$  in (20) is used as the similarity matrix for the initial regions, denoted as  $Z^{(0)} = A$ , upon which the iterative region merging can then be performed. In each

iteration, the two regions with the highest similarity will be merged into a new and larger region.

Let  $\mathcal{R}^{(t)} = \{R_1^{(t)}, \dots, R_{|\mathcal{R}^{(t)}|}^{(t)}\}$  denote the set of regions after the  $t$ -th iteration, whose similarity matrix can be updated according to the average-linkage criterion after the region merging process, that is

$$Z^{(t)} = \{Z_{ij}^{(t)}\}_{|\mathcal{R}^{(t)}| \times |\mathcal{R}^{(t)}|} \quad (25)$$

$$Z_{ij}^{(t)} = \frac{1}{|R_i^{(t)}| \cdot |R_j^{(t)}|} \sum_{x_u \in R_i^{(t)}, x_v \in R_j^{(t)}} a_{uv} \quad (26)$$

where  $|\mathcal{R}^{(t)}|$  denotes the number of regions in the region set  $\mathcal{R}^{(t)}$ , and  $|R_i^{(t)}|$  denotes the number of data samples in the region  $R_i^{(t)}$ . Note that in each iteration, the number of regions decrements by one. Obviously, it holds that  $|\mathcal{R}^{(t)}| = N - t$ , and the dendrogram will be completed after exactly  $N - 1$  iterations. Each level of the dendrogram corresponds to a clustering result with a certain number of clusters, and, therefore, the final consensus clustering can be obtained by specifying a level in the dendrogram.

2) *MDEC-SC*: In MDEC-SC, an SC-based consensus function is exploited. We first construct a graph with the data samples treated as graph nodes and the locally weighted co-association matrix  $A$  used as the adjacent matrix, denoted as

$$G = (V, E) \quad (27)$$

where  $V = \mathcal{X}$  is the node set, and  $E$  is the edge set. The edge weight between nodes  $x_i$  and  $x_j$  is defined as  $E_{ij} = A_{ij}$ . The normalized graph Laplacian is computed as

$$L = I - D^{-1/2} A D^{-1/2} \quad (28)$$

where  $D \in \mathbb{R}^{N \times N}$  is the degree matrix. Let  $K$  be the number of clusters that we aim to obtain. Then, we compute the eigenvectors corresponding to the first  $K$  eigenvalues of  $L$ , which are further stacked to form a new matrix  $U \in \mathbb{R}^{N \times K}$  with each eigenvector being a column of it. Thereafter, the matrix  $U$  will be row normalized, upon which  $K$ -means discretization will be performed to obtain the final clustering result.

3) *MDEC-BG*: In MDEC-BG, a BG-based consensus function is presented. Different from MDEC-SC, we build a BG with both data samples and base clusters treated as graph nodes, and then perform BG partitioning to obtain the consensus clustering. Specifically, the BG is defined as

$$\tilde{G} = (U, V, \tilde{E}) \quad (29)$$

where  $U \cup V$  is the node set (with  $U = \mathcal{X}$  and  $V = \mathcal{C}$ ), and  $\tilde{E}$  is the edge set. An edge between two nodes exists if and only if one of them is a data sample and the other one is a base cluster that contains this data sample. Typically, given two nodes  $u_i \in \mathcal{X}$  and  $v_j \in \mathcal{C}$ , their edge weight will be decided by two factors, that is, their belonging-to relationship and the reliability of the connected cluster, which can be estimated by the ECI index. Formally, the edge weight between  $u_i \in \mathcal{X}$  and  $v_j \in \mathcal{C}$  is defined as

$$\tilde{E}_{ij} = \begin{cases} \text{ECI}(v_j), & \text{if } u_i \in v_j \\ 0, & \text{otherwise.} \end{cases} \quad (30)$$

TABLE I  
DESCRIPTION OF THE EIGHTEEN CANCER GENE EXPRESSION DATASETS

Dataset	Abbr.	#Sample	Dimension	#Class
Bhattacharjee-2001	<i>DS-1</i>	203	1,543	5
Bredel-2005	<i>DS-2</i>	50	1,739	3
Chowdary-2006	<i>DS-3</i>	104	182	2
Garber-2001	<i>DS-4</i>	66	4,553	4
Golub-1999-v1	<i>DS-5</i>	72	1,877	2
Golub-1999-v2	<i>DS-6</i>	72	1,877	3
Nutt-2003-v1	<i>DS-7</i>	50	1,377	4
Pomeroy-2002-v2	<i>DS-8</i>	42	1,379	5
Su-2001	<i>DS-9</i>	174	1,571	10
West-2001	<i>DS-10</i>	49	1,198	2
Yeoh-2002-v1	<i>DS-11</i>	248	2,526	2
Yeoh-2002-v2	<i>DS-12</i>	248	2,526	6
Alizadeh-2000-v2	<i>DS-13</i>	62	2,093	3
Alizadeh-2000-v3	<i>DS-14</i>	62	2,093	4
Armstrong-2002-v2	<i>DS-15</i>	72	2,194	3
Ramaswamy-2001	<i>DS-16</i>	190	1,363	14
Risinger-2003	<i>DS-17</i>	42	1,771	4
Tomlins-2006-v1	<i>DS-18</i>	104	2,315	5

Note that the ECI value reflects the reliability of a cluster with respect to the entire ensemble of base clusterings [as shown in (16)–(19)]. By taking advantage of ECI to define the edge weights in the BG  $\tilde{G}$ , the clusters with higher reliability can exhibit greater influence on the BG and thus play a bigger role in the consensus process. Then, with the BG constructed, we further adopt the transfer cut [56] to efficiently partition the graph nodes into several subsets. By treating the data samples in each subset as a final cluster, the consensus clustering result can be obtained.

#### IV. EXPERIMENTS

In this section, we conduct experiments on a variety of real-world high-dimensional datasets to compare the proposed three ensemble clustering algorithms against nine state-of-the-art algorithms.

##### A. Datasets and Evaluation Measures

In our experiments, 30 real-world high-dimensional datasets are used, including 18 cancer gene expression datasets [57] and 12 image/speech datasets. Among the 12 speech/image datasets, ISOLET [58] is a speech recognition dataset, while the other 11 datasets, including COIL20 [59], MSRA25 [60], Binary Alphadigits [61], YaleB32 [60], Semeion [58], UMist [62], Multiple Features [58], Flowers17 [63], MNIST [61], USPS [61], and Gisette [58], are image datasets. To simplify the description, in the following, 30 benchmark datasets will be abbreviated as DS-1 to DS-30, respectively (as shown in Tables I and II).

To evaluate the quality of the clustering result, two widely used evaluation measures are adopted, namely, normalized mutual information (NMI) [39] and adjusted rand index (ARI) [64]. Note that greater values of NMI and ARI indicate better clustering results.

##### B. Baseline Methods and Experimental Settings

In our experiments, we compare the proposed ensemble clustering algorithms with nine baseline algorithms, that is, stratified sampling-based CSPA (SSCSPA) [6], stratified

TABLE II  
DESCRIPTION OF THE TWELVE IMAGE/SPEECH DATASETS

Dataset	Abbr.	#Sample	Dimension	#Class
ISOLET	<i>DS-19</i>	7,797	617	26
COIL20	<i>DS-20</i>	1,440	1,024	20
MSRA25	<i>DS-21</i>	1,799	256	12
Binary Alphadigits	<i>DS-22</i>	1,404	320	36
YaleB32	<i>DS-23</i>	2,414	1,024	38
Semeion	<i>DS-24</i>	1,593	256	10
UMist	<i>DS-25</i>	575	10,304	20
Multiple Features	<i>DS-26</i>	2,000	649	10
Flowers17	<i>DS-27</i>	1,360	30,000	17
MNIST	<i>DS-28</i>	5,000	784	10
Gisette	<i>DS-29</i>	7,000	5,000	2
USPS	<i>DS-30</i>	11,000	256	10

sampling-based HGPA (SSHGPA) [6], stratified sampling-based MCLA (SSMCLA) [6],  $K$ -means-based consensus clustering (KCC) [7], PTGP [8], SEC [11], entropy-based consensus clustering (ECC) [12], locally weighted graph partitioning (LWGP) [15], and self-paced clustering ensemble (SPCE) [65].

For the baseline algorithms, the parameters will be set as suggested by their corresponding papers. For the proposed algorithms, to produce a set of diversified metrics, the two kernel parameters  $\mu$  and  $k$  are randomized in the ranges of [0.2, 0.8] and [5, 20], respectively. To generate the ensemble of base clusterings, the ensemble size  $M = 100$  and the sampling ratio  $\tau = 0.5$  are used. The number of clusters in each base clustering is randomly selected in the range of [2,  $\sqrt{N}$ ]. Furthermore, the performances of our algorithms with respect to different ensemble sizes  $M$  and different sampling ratios  $\tau$  will also be evaluated in Sections IV-D and IV-E, respectively.

##### C. Comparison With Other Ensemble Clustering Methods

In this section, we evaluate the proposed ensemble clustering algorithms, namely, MDEC-HC, MDEC-SC, and MDEC-BG, against nine baseline ensemble clustering algorithms. For each test algorithm, the number of clusters is set to the true number of classes on the dataset, which is a commonly adopted experimental protocol in [11] and [12]. If an algorithm is not computationally feasible on a dataset due to the out-of-memory error, its score on this dataset will be labeled as “OM.” For each dataset, we run every test algorithm 100 times, and report their average NMI and ARI scores in Tables III and IV, respectively.

In terms of NMI, as shown in Table III, the proposed algorithms outperform the baseline algorithms on most of the benchmark datasets. The SPCE achieves the best NMI scores on 5 out of the 30 datasets, but the three proposed algorithms outperform SPCE on most of the other datasets. It is noteworthy that the proposed MDEC-HC algorithm achieves the best NMI scores on 16 out of the 30 datasets. Also, the proposed MDEC-HC, MDEC-SC, and MDEC-BG algorithms rank in the top three positions 26, 27, and 25 times, respectively, out of the total of 30 datasets, while the fourth best algorithm (i.e., SPCE) ranks in the top three positions just seven times. In terms of ARI, as shown in Table III, our three algorithms also exhibit clear advantages over the baselines, ranking in the



TABLE III

AVERAGE PERFORMANCES (WITH RESPECT TO NMI) OVER 100 RUNS BY DIFFERENT ENSEMBLE CLUSTERING ALGORITHMS. ON EACH DATASET, THE HIGHEST THREE SCORES ARE HIGHLIGHTED IN **BOLD**, WHILE THE HIGHEST ONE IN [**BOLD AND BRACKETS**]

Dataset	SSCSPA	SSHGPA	SSMCLA	KCC	PTGP	SEC	ECC	LWGP	SPCE	MDEC-HC	MDEC-SC	MDEC-BG
DS-1	20.63 $\pm$ 2.11	19.06 $\pm$ 2.15	26.71 $\pm$ 2.95	26.88 $\pm$ 6.93	39.86 $\pm$ 1.55	27.56 $\pm$ 3.04	35.90 $\pm$ 1.23	42.10 $\pm$ 4.82	[53.53 $\pm$ 2.27]	52.41 $\pm$ 1.68	51.72 $\pm$ 1.25	51.64 $\pm$ 1.11
DS-2	<b>40.53</b> $\pm$ 2.79	32.54 $\pm$ 7.33	33.88 $\pm$ 2.94	36.12 $\pm$ 2.79	31.50 $\pm$ 3.62	36.01 $\pm$ 3.67	36.52 $\pm$ 5.94	31.56 $\pm$ 3.28	<b>40.64</b> $\pm$ 5.27	[45.23 $\pm$ 3.35]	38.44 $\pm$ 3.19	38.58 $\pm$ 3.02
DS-3	60.26 $\pm$ 4.47	73.58 $\pm$ 4.94	47.99 $\pm$ 17.36	60.13 $\pm$ 3.26	8.15 $\pm$ 0.00	50.27 $\pm$ 11.41	61.18 $\pm$ 1.25	8.15 $\pm$ 0.00	23.83 $\pm$ 18.80	[85.97 $\pm$ 0.00]	[85.97 $\pm$ 0.00]	[85.97 $\pm$ 0.00]
DS-4	15.02 $\pm$ 0.84	17.03 $\pm$ 1.87	14.49 $\pm$ 2.40	15.34 $\pm$ 4.08	14.20 $\pm$ 2.49	13.57 $\pm$ 2.67	14.89 $\pm$ 3.15	12.44 $\pm$ 2.59	[27.76 $\pm$ 1.53]	21.77 $\pm$ 0.82	<b>22.38</b> $\pm$ 4.05	<b>22.86</b> $\pm$ 0.79
DS-5	42.42 $\pm$ 5.77	46.91 $\pm$ 6.64	62.53 $\pm$ 9.31	55.65 $\pm$ 3.73	57.11 $\pm$ 6.67	49.37 $\pm$ 6.05	53.46 $\pm$ 2.11	56.45 $\pm$ 2.69	51.02 $\pm$ 0.60	[74.92 $\pm$ 13.84]	73.75 $\pm$ 7.77	71.79 $\pm$ 11.81
DS-6	48.35 $\pm$ 3.59	51.73 $\pm$ 7.02	42.31 $\pm$ 15.19	53.68 $\pm$ 8.42	67.32 $\pm$ 5.32	59.31 $\pm$ 6.68	67.51 $\pm$ 4.18	65.98 $\pm$ 3.90	45.22 $\pm$ 2.62	[76.17 $\pm$ 5.42]	75.14 $\pm$ 3.33	76.04 $\pm$ 2.86
DS-7	40.61 $\pm$ 3.50	37.22 $\pm$ 6.00	36.71 $\pm$ 3.90	46.31 $\pm$ 3.48	<b>48.40</b> $\pm$ 0.98	46.53 $\pm$ 2.67	43.21 $\pm$ 1.58	44.79 $\pm$ 0.60	[49.78 $\pm$ 2.27]	48.26 $\pm$ 0.82	45.50 $\pm$ 3.89	47.36 $\pm$ 1.04
DS-8	50.30 $\pm$ 4.17	56.29 $\pm$ 5.27	40.85 $\pm$ 12.12	57.55 $\pm$ 5.32	58.67 $\pm$ 2.73	52.92 $\pm$ 8.03	56.58 $\pm$ 3.09	55.53 $\pm$ 2.78	63.20 $\pm$ 2.72	[66.16 $\pm$ 2.44]	65.53 $\pm$ 2.77	64.80 $\pm$ 1.14
DS-9	52.05 $\pm$ 1.40	52.72 $\pm$ 3.49	44.66 $\pm$ 10.09	57.22 $\pm$ 2.78	57.02 $\pm$ 3.24	56.91 $\pm$ 4.53	56.95 $\pm$ 2.76	53.15 $\pm$ 2.26	58.26 $\pm$ 7.31	66.57 $\pm$ 2.27	[69.55 $\pm$ 0.93]	69.41 $\pm$ 1.32
DS-10	29.05 $\pm$ 4.68	[35.75 $\pm$ 1.25]	29.21 $\pm$ 4.49	31.17 $\pm$ 3.09	31.18 $\pm$ 1.33	28.29 $\pm$ 4.29	31.22 $\pm$ 0.00	4.08 $\pm$ 0.00	21.37 $\pm$ 5.62	32.03 $\pm$ 3.92	32.45 $\pm$ 2.21	34.29 $\pm$ 2.37
DS-11	20.20 $\pm$ 0.30	18.60 $\pm$ 1.37	4.56 $\pm$ 4.57	79.18 $\pm$ 11.23	84.56 $\pm$ 2.01	56.61 $\pm$ 35.73	85.59 $\pm$ 6.47	1.92 $\pm$ 10.32	41.78 $\pm$ 22.53	[94.35 $\pm$ 2.49]	92.91 $\pm$ 2.50	93.64 $\pm$ 2.13
DS-12	31.48 $\pm$ 2.51	32.65 $\pm$ 4.85	9.42 $\pm$ 5.32	33.40 $\pm$ 1.31	34.93 $\pm$ 4.53	32.85 $\pm$ 2.87	36.09 $\pm$ 3.73	26.44 $\pm$ 1.05	41.48 $\pm$ 1.17	45.98 $\pm$ 5.39	50.81 $\pm$ 5.34	[57.01 $\pm$ 3.79]
DS-13	57.11 $\pm$ 1.02	48.88 $\pm$ 3.43	56.71 $\pm$ 2.28	65.40 $\pm$ 18.11	<b>97.56</b> $\pm$ 7.89	59.52 $\pm$ 13.51	70.32 $\pm$ 20.54	71.92 $\pm$ 3.87	91.10 $\pm$ 9.48	[98.34 $\pm$ 3.44]	94.86 $\pm$ 4.31	93.38 $\pm$ 3.74
DS-14	42.77 $\pm$ 3.38	49.40 $\pm$ 2.65	49.02 $\pm$ 3.82	52.26 $\pm$ 4.44	<b>63.97</b> $\pm$ 2.07	50.11 $\pm$ 9.60	63.47 $\pm$ 2.07	53.01 $\pm$ 4.60	[65.99 $\pm$ 2.49]	65.28 $\pm$ 1.46	63.61 $\pm$ 1.63	63.02 $\pm$ 1.42
DS-15	66.82 $\pm$ 2.16	49.13 $\pm$ 6.83	66.67 $\pm$ 3.09	50.34 $\pm$ 5.45	75.00 $\pm$ 4.99	41.77 $\pm$ 7.09	72.08 $\pm$ 5.32	58.02 $\pm$ 9.73	61.86 $\pm$ 8.91	[80.90 $\pm$ 2.71]	80.42 $\pm$ 2.11	78.56 $\pm$ 2.16
DS-16	50.44 $\pm$ 1.65	54.22 $\pm$ 1.67	18.71 $\pm$ 4.94	45.13 $\pm$ 2.79	38.42 $\pm$ 2.16	38.05 $\pm$ 5.40	43.57 $\pm$ 8.85	29.19 $\pm$ 0.86	52.51 $\pm$ 1.15	65.56 $\pm$ 1.53	66.60 $\pm$ 1.18	[66.64 $\pm$ 1.16]
DS-17	30.98 $\pm$ 5.61	30.90 $\pm$ 4.37	16.20 $\pm$ 9.87	29.09 $\pm$ 3.41	31.32 $\pm$ 2.60	28.12 $\pm$ 3.68	31.93 $\pm$ 4.47	15.57 $\pm$ 0.79	32.05 $\pm$ 1.04	[36.79 $\pm$ 2.93]	32.77 $\pm$ 1.44	33.13 $\pm$ 1.17
DS-18	38.12 $\pm$ 1.90	35.95 $\pm$ 3.24	42.68 $\pm$ 3.19	44.86 $\pm$ 5.88	44.66 $\pm$ 1.78	37.92 $\pm$ 5.01	40.50 $\pm$ 1.83	42.61 $\pm$ 2.06	46.57 $\pm$ 2.01	[61.11 $\pm$ 2.59]	53.30 $\pm$ 2.24	54.30 $\pm$ 2.04
DS-19	69.12 $\pm$ 1.52	60.32 $\pm$ 1.56	73.81 $\pm$ 1.00	68.43 $\pm$ 1.04	72.96 $\pm$ 0.70	68.71 $\pm$ 1.25	70.06 $\pm$ 0.88	74.91 $\pm$ 0.53	72.35 $\pm$ 1.86	[77.09 $\pm$ 0.77]	76.41 $\pm$ 0.29	76.48 $\pm$ 0.50
DS-20	78.42 $\pm$ 1.17	76.77 $\pm$ 2.20	75.38 $\pm$ 4.97	72.33 $\pm$ 1.98	67.00 $\pm$ 2.65	72.06 $\pm$ 3.01	74.26 $\pm$ 1.37	80.16 $\pm$ 1.10	86.01 $\pm$ 1.28	[91.66 $\pm$ 0.75]	90.15 $\pm$ 1.06	91.25 $\pm$ 0.78
DS-21	59.41 $\pm$ 1.50	61.66 $\pm$ 2.84	65.07 $\pm$ 2.07	61.05 $\pm$ 1.61	63.67 $\pm$ 1.97	59.79 $\pm$ 3.59	62.70 $\pm$ 2.39	66.29 $\pm$ 1.59	69.31 $\pm$ 1.56	70.37 $\pm$ 1.05	[71.78 $\pm$ 1.64]	71.03 $\pm$ 1.28
DS-22	55.31 $\pm$ 0.97	46.99 $\pm$ 1.35	54.45 $\pm$ 1.35	54.21 $\pm$ 0.73	55.06 $\pm$ 0.57	56.41 $\pm$ 0.60	54.79 $\pm$ 0.59	51.67 $\pm$ 0.87	<b>60.57</b> $\pm$ 0.65	58.06 $\pm$ 1.31	[62.54 $\pm$ 0.62]	62.33 $\pm$ 0.62
DS-23	10.63 $\pm$ 0.31	11.86 $\pm$ 0.52	11.10 $\pm$ 0.50	10.11 $\pm$ 0.40	10.15 $\pm$ 0.40	10.39 $\pm$ 0.45	10.47 $\pm$ 0.34	9.40 $\pm$ 0.33	27.29 $\pm$ 1.23	34.12 $\pm$ 1.36	35.40 $\pm$ 0.98	[37.54 $\pm$ 1.25]
DS-24	55.69 $\pm$ 3.29	35.82 $\pm$ 2.07	57.51 $\pm$ 2.33	53.22 $\pm$ 1.77	62.64 $\pm$ 0.69	54.39 $\pm$ 3.28	53.61 $\pm$ 1.74	61.93 $\pm$ 3.57	62.27 $\pm$ 3.03	[71.33 $\pm$ 0.64]	67.85 $\pm$ 1.41	68.51 $\pm$ 1.25
DS-25	57.93 $\pm$ 1.22	65.52 $\pm$ 2.53	58.69 $\pm$ 3.61	60.85 $\pm$ 1.74	62.67 $\pm$ 1.11	60.49 $\pm$ 1.39	61.37 $\pm$ 1.20	62.50 $\pm$ 1.53	69.84 $\pm$ 1.34	[83.97 $\pm$ 1.77]	82.24 $\pm$ 2.73	83.15 $\pm$ 2.45
DS-26	81.34 $\pm$ 3.10	69.94 $\pm$ 3.07	86.33 $\pm$ 2.11	73.08 $\pm$ 3.97	83.45 $\pm$ 2.14	64.14 $\pm$ 6.06	75.24 $\pm$ 1.57	87.02 $\pm$ 1.64	82.00 $\pm$ 4.02	90.68 $\pm$ 2.09	[94.19 $\pm$ 0.30]	93.99 $\pm$ 0.26
DS-27	21.75 $\pm$ 0.75	18.75 $\pm$ 1.24	23.74 $\pm$ 0.59	24.15 $\pm$ 0.47	24.96 $\pm$ 0.39	24.18 $\pm$ 0.49	24.16 $\pm$ 0.39	21.60 $\pm$ 0.75	[37.09 $\pm$ 0.38]	25.44 $\pm$ 0.64	27.56 $\pm$ 0.39	27.14 $\pm$ 0.38
DS-28	53.13 $\pm$ 3.25	38.54 $\pm$ 2.66	58.34 $\pm$ 2.20	53.46 $\pm$ 3.18	64.53 $\pm$ 1.64	53.10 $\pm$ 2.26	53.42 $\pm$ 2.23	60.13 $\pm$ 1.22	51.77 $\pm$ 8.95	[78.47 $\pm$ 2.19]	76.04 $\pm$ 1.20	77.92 $\pm$ 0.96
DS-29	20.44 $\pm$ 6.59	19.28 $\pm$ 10.95	47.98 $\pm$ 4.69	41.41 $\pm$ 9.87	49.82 $\pm$ 0.87	10.73 $\pm$ 6.22	46.07 $\pm$ 2.01	5.62 $\pm$ 14.21	9.23 $\pm$ 0.54	67.23 $\pm$ 5.28	67.43 $\pm$ 0.68	[67.94 $\pm$ 0.34]
DS-30	47.98 $\pm$ 2.88	34.33 $\pm$ 2.08	53.71 $\pm$ 2.89	51.17 $\pm$ 1.52	65.00 $\pm$ 1.84	47.13 $\pm$ 2.49	52.74 $\pm$ 1.77	59.30 $\pm$ 1.41	OM	[77.58 $\pm$ 1.37]	74.91 $\pm$ 1.26	75.33 $\pm$ 2.27
Avg. score	44.94	42.74	43.65	48.77	52.19	44.91	51.33	43.78	-	64.79	64.07	64.50
Avg. rank	8.80	9.13	8.73	8.00	6.13	9.13	7.00	8.40	5.53	2.00	2.63	2.37

\* Note that OM indicates the out-of-memory error.

top three positions 28, 25, and 25 times, respectively, while the fourth best algorithm only ranks in the top three positions four times.

Furthermore, to provide a summary view across the 30 benchmark datasets, we report the *average score* and *average rank* of different algorithms in the last two rows in Tables III and IV. Note that the average score (across 30 datasets) is computed by taking the average on NMI (or ARI) scores, while the average rank is obtained by taking the average on the ranking positions, for each algorithm across all the benchmark datasets. Note that if an algorithm is not computationally feasible on all datasets, it will not have an average score, but will still have an average rank, where the *lost* score on a dataset (due to the out-of-memory error) will be considered as the last position on this dataset. As can be seen in Table III, the proposed three algorithms achieve average NMI(%) scores of 64.79, 64.07, and 64.50, respectively, which are significantly higher than the fourth best average score of 52.19. In terms of the ranking positions in Table III, the proposed three algorithms obtain average ranks of 2.00, 2.63, and 2.37, respectively, while the fourth best algorithm only obtains an average rank of 5.53. Similar advantages can also be observed in terms of ARI. The average ARI scores and the average ranks (across the 30 datasets) of the proposed algorithms significantly outperform the nine baseline algorithms (as shown in Table IV).

#### D. Robustness to Ensemble Sizes

In this section, we compare the performances of different ensemble clustering algorithms with varying ensemble sizes. Specifically, we perform the test algorithms on the datasets

with the ensemble size varying from 20 to 300, and report their average NMI and ARI scores in Figs. 2 and 3, respectively.

In terms of NMI, as can be seen in Fig. 2, the proposed algorithms yield stably high performances across the 30 benchmark datasets, with varying ensemble sizes. Although the SPCE algorithm outperforms our algorithms on DS-7 and DS-27 datasets, yet on most of the other datasets, our algorithms achieve better or comparable performance when compared to the baseline ensemble clustering algorithms. Especially, on DS-3, DS-4, DS-5, DS-6, DS-9, DS-16, DS-18, DS-20, DS-23, DS-24, DS-25, DS-26, DS-28, DS-29, and DS-30 datasets, our three proposed algorithms have shown clear advantages over the baseline algorithms as the ensemble size goes from 20 to 300. Similarly, in terms of ARI, our algorithms also exhibit very competitive performances on most of the benchmark datasets with varying ensemble sizes (as can be seen in Fig. 3).

Furthermore, in Fig. 4, we illustrate the average NMI and ARI curves (across the 30 benchmark datasets) by different ensemble clustering algorithms with varying ensemble sizes. Specifically, Fig. 4(a) is obtained by taking the average of 30 subfigures in Fig. 2, while Fig. 4(b) by taking the average of 30 subfigures in Fig. 3. As can be observed in Fig. 4, the proposed algorithms achieve significantly better performances (with respect to both NMI and ARI) than the baseline ensemble clustering algorithms across the 30 benchmark datasets. Even when compared to the best two baseline algorithms, that is, PTGP and ECC, the proposed three algorithms can still consistently achieve approximately 20% higher average NMI/ARI scores.



TABLE IV

AVERAGE PERFORMANCES (WITH RESPECT TO ARI) OVER 100 RUNS BY DIFFERENT ENSEMBLE CLUSTERING ALGORITHMS. ON EACH DATASET, THE HIGHEST THREE SCORES ARE HIGHLIGHTED IN **BOLD**, WHILE THE HIGHEST ONE IN **[BOLD AND BRACKETS]**

Dataset	SSCSPA	SSHGPA	SSMCLA	KCC	PTGP	SEC	ECC	LWGP	SPCE	MDEC-HC	MDEC-SC	MDEC-BG
DS-1	8.25±1.68	7.01±1.67	12.26±3.33	13.74±7.37	23.96±1.57	12.35±3.34	22.08±1.33	30.83±10.37	<b>[51.88±4.14]</b>	<b>50.42±2.51</b>	<b>46.94±2.15</b>	46.71±1.81
DS-2	36.46±2.66	29.64±6.81	30.95±3.64	35.63±2.62	33.18±6.35	33.71±5.18	35.52±4.90	22.98±7.25	<b>47.03±9.54</b>	<b>[53.89±4.35]</b>	46.82±2.67	<b>46.99±2.42</b>
DS-3	65.20±3.14	80.99±3.44	52.52±20.09	69.40±4.03	6.57±0.00	58.71±13.38	71.06±1.35	6.57±0.00	29.80±28.79	<b>[92.38±0.00]</b>	<b>[92.38±0.00]</b>	<b>[92.38±0.00]</b>
DS-4	8.97±0.50	9.38±1.89	<b>17.26±6.14</b>	14.44±8.72	8.85±2.55	11.18±7.49	14.64±6.83	15.99±3.22	<b>[21.94±4.88]</b>	<b>18.42±0.95</b>	16.89±3.03	16.98±0.78
DS-5	44.73±4.34	45.15±10.61	73.83±8.22	64.01±4.32	65.53±6.71	57.42±5.94	61.72±2.50	65.00±2.85	59.14±0.65	<b>80.99±16.82</b>	<b>[82.41±9.41]</b>	<b>78.53±15.58</b>
DS-6	46.12±3.25	49.59±10.14	39.78±18.56	51.89±14.45	68.29±6.54	60.02±9.10	72.63±5.64	65.13±3.73	56.21±5.55	<b>79.45±6.05</b>	<b>80.14±3.37</b>	<b>[80.92±3.17]</b>
DS-7	24.54±2.60	23.19±5.12	20.48±5.47	34.18±4.19	<b>36.43±1.79</b>	35.02±3.46	26.20±2.62	35.56±0.59	32.87±0.69	<b>[39.86±1.24]</b>	34.12±5.93	<b>38.57±1.57</b>
DS-8	36.14±6.23	44.96±7.56	24.64±13.25	47.44±6.69	47.29±4.63	42.17±10.56	46.25±4.74	45.62±3.86	50.18±3.11	<b>[58.47±2.85]</b>	<b>57.68±2.86</b>	<b>57.04±1.31</b>
DS-9	36.17±1.94	35.40±4.20	28.27±11.94	42.58±3.68	41.44±4.38	42.03±5.91	40.05±4.08	38.63±2.83	38.14±8.43	<b>52.65±2.57</b>	<b>55.98±1.60</b>	<b>[56.48±2.35]</b>
DS-10	35.96±5.65	<b>[43.82±1.40]</b>	30.87±5.69	38.30±3.28	38.67±1.66	33.22±6.87	38.75±0.00	0.01±0.00	25.41±10.12	39.53±4.62	<b>40.12±2.49</b>	<b>42.19±2.65</b>
DS-11	12.36±0.14	8.65±3.13	13.12±8.07	88.94±14.60	93.09±1.10	56.16±49.17	93.54±3.52	0.22±12.44	53.28±35.18	<b>[97.90±1.04]</b>	<b>97.27±1.13</b>	<b>97.60±0.94</b>
DS-12	23.12±2.21	24.41±4.01	6.17±5.31	18.46±2.35	18.15±5.31	15.03±3.78	20.46±3.96	17.53±0.64	17.68±1.57	<b>30.83±6.96</b>	<b>32.85±5.88</b>	<b>[41.49±3.12]</b>
DS-13	43.00±0.74	39.27±2.45	43.39±2.01	55.44±23.52	<b>97.68±9.43</b>	48.03±17.25	61.88±26.89	76.94±9.01	94.16±6.26	<b>[99.00±2.08]</b>	<b>96.88±2.61</b>	95.98±2.27
DS-14	26.77±4.98	33.34±3.38	35.32±3.29	40.71±7.03	41.90±2.85	31.78±7.48	41.49±1.70	42.84±3.49	<b>[50.15±1.88]</b>	<b>45.52±1.30</b>	<b>43.09±1.23</b>	42.87±1.22
DS-15	68.50±2.50	48.65±7.02	65.12±3.81	46.04±6.69	76.81±6.68	34.81±8.69	72.54±7.32	56.48±10.58	65.09±12.86	<b>[84.57±2.99]</b>	<b>84.10±2.24</b>	<b>82.13±2.33</b>
DS-16	28.99±1.67	36.03±1.93	7.45±3.67	17.68±4.70	6.17±1.50	10.47±4.51	16.87±7.24	3.07±0.38	6.86±0.68	<b>[57.54±3.33]</b>	<b>48.82±2.27</b>	<b>48.74±3.14</b>
DS-17	17.65±6.12	<b>21.11±3.90</b>	5.48±9.14	10.69±5.22	13.90±3.04	12.29±5.34	14.27±3.02	-6.19±0.48	8.74±2.76	<b>[22.16±3.96]</b>	18.87±2.57	<b>19.48±2.14</b>
DS-18	24.15±1.51	24.24±3.64	29.01±4.07	30.23±5.14	29.89±3.72	22.30±4.04	25.08±2.68	27.22±2.30	26.13±2.58	<b>[51.33±4.08]</b>	<b>40.39±3.17</b>	<b>41.91±2.49</b>
DS-19	46.24±2.84	33.48±2.31	<b>53.59±1.74</b>	44.11±2.07	46.03±2.12	43.64±2.18	44.63±1.85	<b>54.18±1.68</b>	44.69±5.44	<b>[56.57±1.90]</b>	49.05±0.60	50.49±1.68
DS-20	64.16±1.84	59.10±3.94	56.54±9.01	52.50±3.59	37.76±4.80	51.94±5.44	54.40±2.30	64.33±1.99	64.99±3.39	<b>[81.58±1.30]</b>	<b>77.21±2.72</b>	<b>80.21±1.68</b>
DS-21	39.90±1.88	41.42±2.85	43.39±2.54	38.72±1.92	37.71±2.89	35.75±5.65	39.99±3.26	42.29±2.83	42.64±4.96	<b>44.16±1.61</b>	<b>[46.13±2.49]</b>	<b>45.08±1.85</b>
DS-22	25.06±1.23	15.59±1.15	24.82±1.88	25.01±1.21	27.02±0.52	27.16±0.77	25.46±1.04	24.94±0.94	24.21±2.45	<b>30.56±1.61</b>	<b>34.27±0.90</b>	<b>[34.62±0.93]</b>
DS-23	0.63±0.09	1.05±0.24	0.91±0.14	0.61±0.07	0.84±0.09	0.80±0.12	0.64±0.15	0.72±0.06	0.41±0.08	7.77±0.58	<b>11.23±0.54</b>	<b>[13.11±0.89]</b>
DS-24	45.44±4.83	21.30±1.95	46.18±3.41	39.94±3.42	48.25±1.13	39.33±4.61	38.02±3.31	51.21±3.99	49.28±3.89	<b>[60.23±1.13]</b>	<b>53.71±2.35</b>	<b>55.74±2.64</b>
DS-25	27.98±1.33	35.84±3.76	27.62±4.08	30.77±2.06	31.17±1.11	29.95±1.88	31.24±1.45	33.48±1.69	37.60±3.33	<b>[67.41±2.82]</b>	<b>61.32±6.10</b>	<b>63.91±5.02</b>
DS-26	78.94±4.63	61.55±4.35	84.82±3.40	62.73±6.98	79.62±4.36	48.52±8.63	65.22±2.48	85.68±2.94	77.17±7.15	<b>88.10±4.07</b>	<b>[94.46±0.44]</b>	<b>94.19±0.40</b>
DS-27	7.98±0.40	6.88±0.85	8.97±0.51	9.51±0.32	9.34±0.24	9.22±0.37	<b>9.70±0.44</b>	9.69±0.47	4.03±0.59	9.48±0.38	<b>10.08±0.30</b>	<b>[10.17±0.35]</b>
DS-28	42.13±4.44	25.53±3.13	46.71±3.13	40.52±3.60	53.84±2.85	40.97±3.11	40.84±3.37	48.54±1.94	34.18±12.43	<b>[70.13±4.54]</b>	<b>65.83±2.05</b>	<b>69.39±1.88</b>
DS-29	26.82±8.18	24.92±13.52	58.48±5.61	50.80±11.04	60.54±0.92	11.89±6.18	56.48±2.15	6.77±17.60	5.27±2.76	<b>76.30±4.91</b>	<b>77.15±0.64</b>	<b>[77.59±0.29]</b>
DS-30	33.95±4.16	19.23±1.81	38.28±3.47	33.41±1.81	52.57±3.11	30.04±3.59	37.33±2.13	43.99±1.67	OM	<b>[67.35±2.13]</b>	<b>61.86±2.41</b>	<b>62.52±3.90</b>
Avg. score	34.21	31.69	34.21	38.28	41.08	32.86	40.63	33.68	-	<b>57.15</b>	<b>55.27</b>	<b>56.13</b>
Avg. rank	8.60	8.87	8.23	7.77	6.57	9.03	7.00	7.33	7.53	<b>1.93</b>	<b>2.73</b>	<b>2.27</b>

### E. Influence of Metrics, Subspaces, and Clusters

In this article, we jointly exploit large populations of diversified metrics, random subspaces, and weighted clusters in a unified ensemble clustering framework, and propose three specific ensemble clustering algorithms, called MDEC-HC, MDEC-SC, and MDEC-BG, respectively. In this section, we will evaluate the influence of the three key factors (i.e., diversified metrics, random subspaces, and weighted clusters) upon the proposed three algorithms.

First, we compare the diversified metrics with several widely used similarity metrics, namely, cosine similarity, correlation coefficient, and Spearman correlation coefficient. For each of the three proposed algorithms, we test its performance of using the diversified metrics against the above three metrics. Taking MDEC-HC as an example, as shown in Fig. 5, using diversified metrics leads to an average NMI(%) score 67.49, while using the three conventional similarity metrics leads to average NMI(%) scores of 50.99, 48.23, and 41.34, respectively. Similar advantages of using diversified metrics can also be observed in the cases of MDEC-SC and MDEC-BG, which confirm that the use of diversified metrics in the proposed algorithms can substantially benefit the consensus clustering performance (as shown in Fig. 5).

Second, we evaluate the performances of the proposed ensemble clustering algorithms with different subspace sampling ratio  $\tau$ , which varies from 0.1 to 1. As can be observed in Figs. 6 and 7, moderate values of  $\tau$  generally lead to better consensus clustering performance. When the sampling ratio  $\tau$  goes from 0.8 to 1, the performance declines, which suggests that the use of random subspaces exhibits a positive

influence when compared to using the full feature sets (by setting  $\tau = 1$ ). At the other extreme, when setting  $\tau$  to very small values, for example, in the range of [0.1, 0.3], the performance also declines, due to the fact that the subspaces generated by a very small sampling ratio may not well reflect the underlying structure of the dataset. Empirically, it is suggested that the sampling ratio  $\tau$  be set in the range of [0.4, 0.8], which strikes a balance between diversity and quality. In this article, the sampling ratio  $\tau = 0.5$  is used in the experiments on all the benchmark datasets.

Third, we evaluate the performances of the proposed algorithms *with* and *without* the weighted clusters. Note that the performances of the proposed algorithms without weighted clusters are obtained by setting all cluster weights equal to one. As shown in Figs. 6 and 7, in terms of both NMI and ARI, the proposed three algorithms with weighted clusters exhibit consistently better average performances (across the 30 benchmark datasets) than that without weighted clusters.

From the comparison results in Figs. 5, 6, and 7, we can have two main observations: 1) the performance of our approach benefits from the use of diversified metrics, random subspaces, and weighted clusters and 2) out of the three beneficial factors, the diversified metrics play the possibly most important role in the consensus clustering performance, with consideration to the clear improvement (with respect to both NMI and ARI) that they lead to.

### F. Execution Time

In this section, we evaluate the efficiency of different ensemble clustering algorithms and report their execution times on the benchmark datasets in Table V. The experiments are

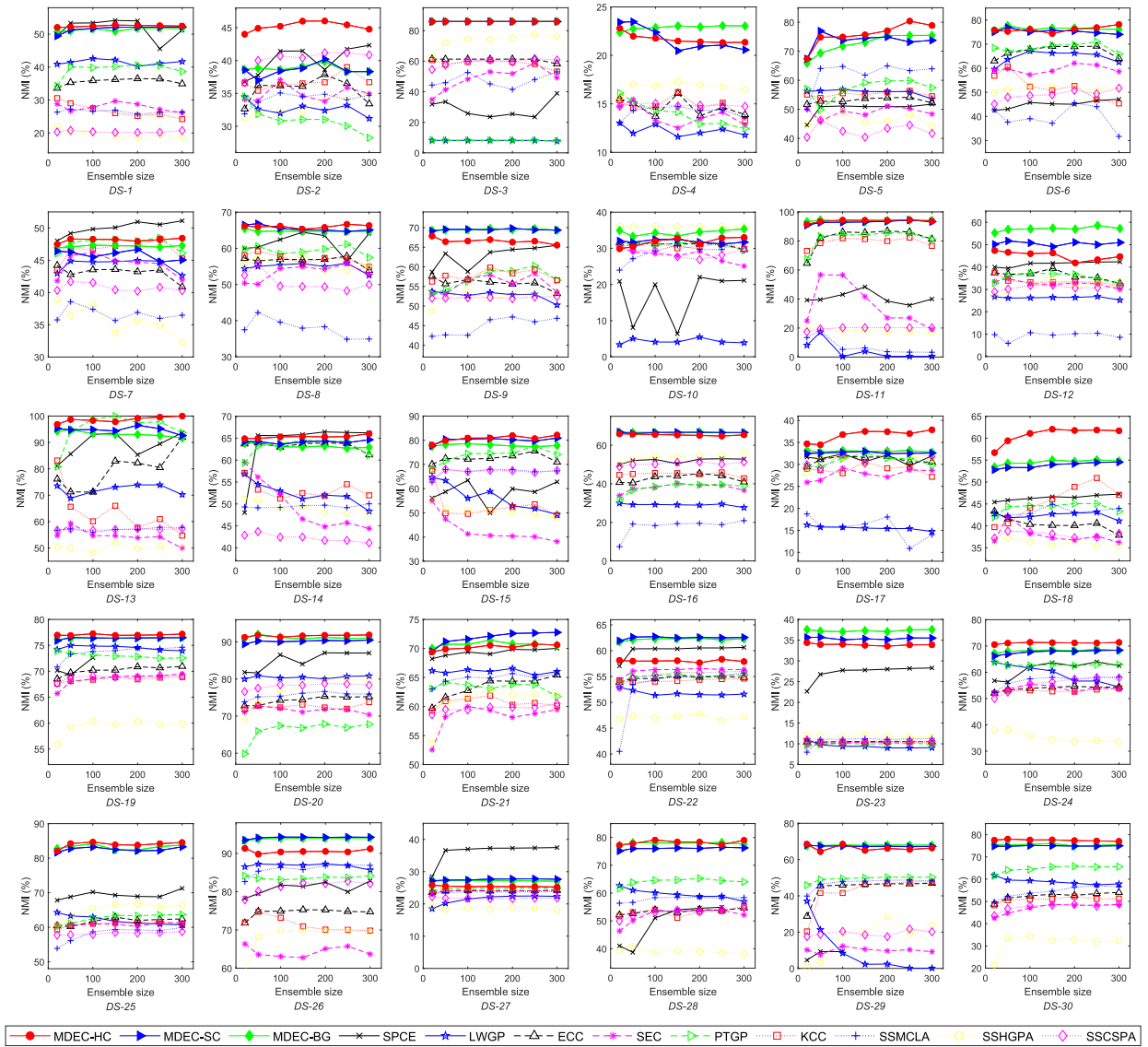


Fig. 2. Average performances (with respect to NMI) over 20 runs by different ensemble clustering algorithms with varying ensemble sizes  $M$ . Note that if an algorithm is not feasible (due to the out-of-memory error) for some large ensemble sizes on a dataset, then its curve will stop earlier than other algorithms.

conducted in MATLAB R2016a 64-bit on a workstation with Intel i9-7940X CPU and 64-GB memory.

In general, larger dimensions and larger sample sizes lead to greater computational costs for the ensemble clustering algorithms. As shown in Table V, the proposed MDEC-HC, MDEC-SC, and MDEC-BG algorithms each consume less than 1 s of time on 14 out of the totally 18 cancer gene expression datasets. On the 12 image/speech datasets, the proposed algorithms also show comparable time efficiency with the other ensemble clustering algorithms.

To summarize, as can be seen in Tables III–V and Figs. 2–4, the proposed three ensemble clustering algorithms have shown considerable advantages in clustering effectiveness while exhibiting competitive time efficiency when compared with the state-of-the-art ensemble clustering algorithms.

## V. CONCLUSION

In this article, we presented a new ensemble clustering approach called MDEC, which is capable of jointly exploiting large populations of diversified metrics, random subspaces, and weighted clusters in a unified ensemble clustering framework. Specifically, a large number of diversified metrics are generated by randomizing a scaled exponential similarity kernel. The diversified metrics are then coupled with the random subspaces to form a large set of metric-subspace pairs. Upon the similarity matrices derived from the metric-subspace pairs, the SC algorithm is performed to construct an ensemble of diversified base clusterings. With the base clusterings generated, an entropy-based cluster validity strategy is utilized to evaluate and weight the clusters with consideration to the distribution of the cluster labels in the entire ensemble. Finally, based on diversified metrics, random subspaces, and weighted clusters, three specific ensemble clustering algorithms are

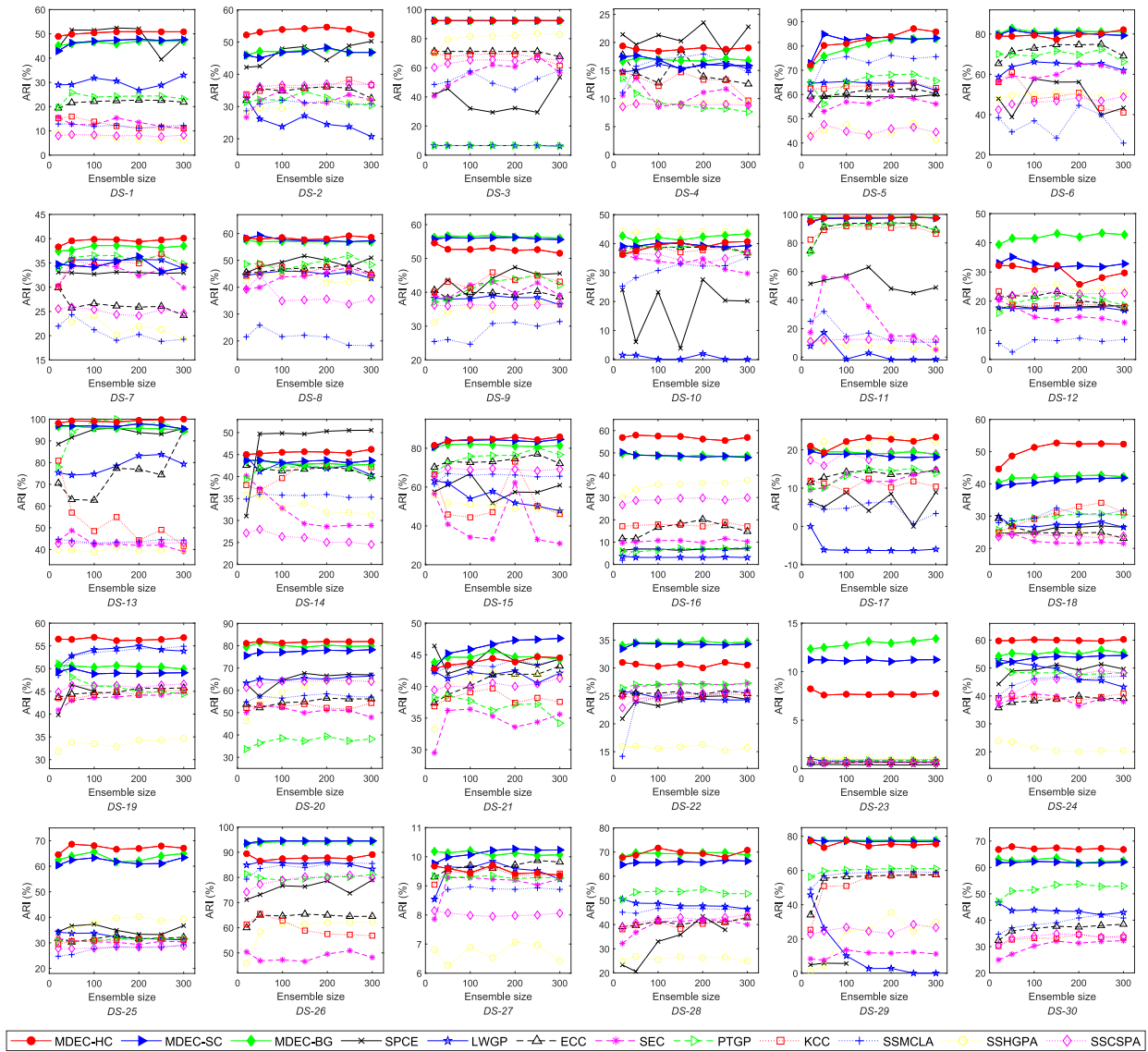


Fig. 3. Average performances (with respect to ARI) over 20 runs by different ensemble clustering algorithms with varying ensemble sizes  $M$ . Note that if an algorithm is not feasible (due to the out-of-memory error) for some large ensemble sizes on a dataset, then its curve will stop earlier than other algorithms.

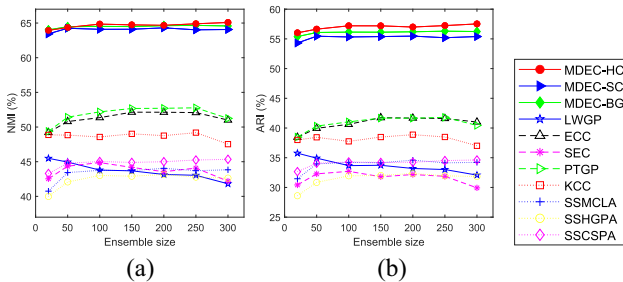


Fig. 4. Average curves (across 30 datasets) by different algorithms with varying ensemble size  $M$ . Note that (a) is obtained by averaging the 30 subfigures in Fig. 2, and (b) 30 subfigures in Fig. 3.

devised by incorporating three types of consensus functions. Extensive experiments are conducted on 30 real-world high-dimensional datasets (including 18 cancer gene expression datasets and 12 image/speech datasets), which have

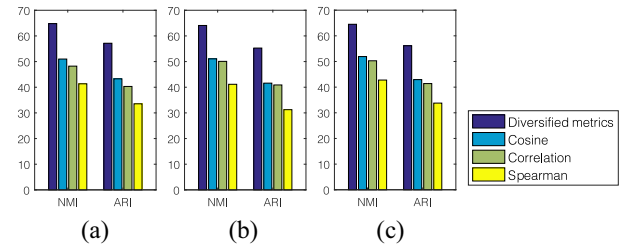


Fig. 5. Average performance (across 30 datasets) of the proposed algorithms using diversified metrics against using other similarity metrics. (a) MDEC-HC. (b) MDEC-SC. (c) MDEC-BG.

demonstrated the advantages of the proposed algorithms over the state-of-the-art ensemble clustering algorithms.

Note that in the practical application process of the proposed algorithms (as well as other ensemble clustering algorithms), two types of balance should be taken into consideration, that is, the balance between diversity and quality (which involves



TABLE V  
EXECUTION TIMES (IN SECONDS) OF DIFFERENT ENSEMBLE CLUSTERING ALGORITHMS (WITH BOTH ENSEMBLE GENERATION AND CONSENSUS FUNCTION INCLUDED)

Dataset	SSCSPA	SSHGPA	SSMCLA	KCC	PTGP	SEC	ECC	LWGP	SPCE	MDEC-HC	MDEC-SC	MDEC-BG
DS-1	1.12	1.51	1.38	0.87	1.25	0.80	0.95	1.65	1.80	1.00	1.01	1.07
DS-2	0.82	1.02	1.01	0.44	0.62	0.42	0.45	1.38	0.53	0.75	0.76	0.77
DS-3	0.91	1.04	1.08	0.45	0.69	0.42	0.46	1.10	0.79	0.77	0.78	0.80
DS-4	1.12	1.31	1.31	0.64	0.81	0.61	0.68	1.24	1.05	0.87	0.88	0.90
DS-5	0.86	1.02	1.06	0.52	0.73	0.49	0.52	1.14	0.71	0.78	0.79	0.80
DS-6	0.86	1.08	1.06	0.54	0.71	0.50	0.56	1.44	0.79	0.78	0.78	0.80
DS-7	1.13	1.00	0.99	0.46	0.64	0.43	0.48	1.04	0.58	0.73	0.74	0.75
DS-8	0.80	0.97	0.97	0.46	0.61	0.43	0.48	1.01	0.56	0.71	0.72	0.72
DS-9	1.06	1.52	1.28	0.79	1.04	0.68	0.90	1.47	1.72	0.90	0.91	0.95
DS-10	0.80	0.93	0.96	0.45	0.61	0.43	0.45	1.01	0.51	0.72	0.73	0.74
DS-11	1.68	1.78	1.89	1.31	1.74	1.23	1.37	2.15	4.01	1.14	1.15	1.21
DS-12	1.52	1.81	1.75	1.34	1.68	1.20	1.50	2.10	3.94	1.15	1.17	1.23
DS-13	0.85	1.09	1.03	0.49	0.67	0.47	0.50	1.07	0.63	0.74	0.74	0.76
DS-14	0.86	1.09	1.03	0.52	0.69	0.49	0.54	1.09	0.64	0.74	0.75	0.76
DS-15	0.89	1.12	1.09	0.52	0.71	0.49	0.54	1.14	0.76	1.16	1.17	1.19
DS-16	1.07	1.48	1.40	0.91	1.75	0.75	1.06	2.26	1.97	0.87	0.89	0.93
DS-17	0.82	0.98	0.99	0.44	0.58	0.41	0.45	0.99	0.55	0.73	0.74	0.75
DS-18	1.02	1.33	1.22	0.66	0.89	0.62	0.70	1.31	1.21	0.83	0.84	0.87
DS-19	41.32	90.94	26.43	61.71	156.60	43.51	87.18	59.94	4032.33	137.40	144.16	139.01
DS-20	3.43	4.22	3.09	5.32	4.89	4.15	6.58	7.41	114.84	5.44	5.71	5.84
DS-21	3.43	3.58	2.91	2.82	2.72	1.48	4.26	5.39	178.47	6.89	7.18	7.37
DS-22	2.98	43.10	3.37	3.52	3.53	2.04	5.78	5.72	106.74	6.09	6.56	6.48
DS-23	9.55	16.17	6.84	17.40	19.09	13.82	20.82	18.93	354.36	15.66	16.92	16.32
DS-24	3.42	5.85	3.54	3.78	4.54	2.37	5.80	6.33	90.75	7.23	7.47	7.64
DS-25	5.92	7.04	6.18	13.46	13.43	13.11	13.90	15.21	29.72	4.72	4.79	4.89
DS-26	5.58	5.36	4.33	6.80	7.66	4.73	8.87	9.37	140.78	9.49	9.82	10.06
DS-27	44.27	46.46	44.12	204.92	204.92	203.34	206.45	206.91	299.38	31.26	31.51	31.64
DS-28	33.05	31.17	18.75	43.00	69.10	33.34	53.65	44.77	1115.70	66.18	68.29	67.67
DS-29	117.37	95.22	97.91	450.27	602.59	441.79	464.34	459.80	3590.84	238.26	240.38	239.67
DS-30	74.94	79.57	37.49	65.89	590.31	43.67	109.65	70.75	OM	302.52	312.08	305.36

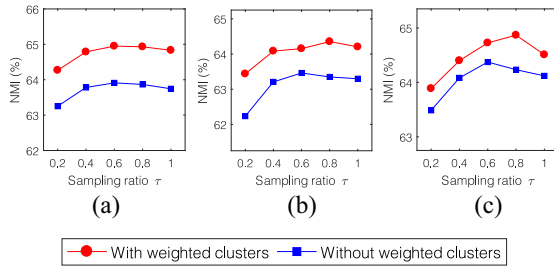


Fig. 6. Average NMI scores (across 30 datasets) of the proposed algorithms with and without weighted clusters using varying sampling ratios  $\tau$ . (a) MDEC-HC. (b) MDEC-SC. (c) MDEC-BG.

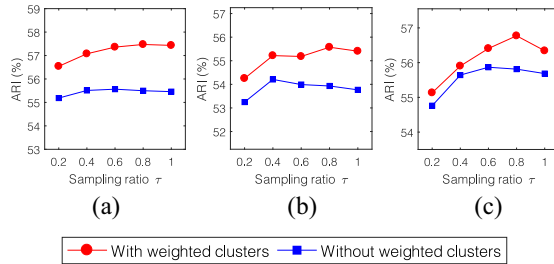


Fig. 7. Average ARI scores (across 30 datasets) of the proposed algorithms with and without weighted clusters using varying sampling ratios  $\tau$ . (a) MDEC-HC. (b) MDEC-SC. (c) MDEC-BG.

enhancing different levels of diversity while maintaining the quality of the individuals) and the balance between effectiveness and efficiency (which involves the choice of base clusterers and the size of the ensembles).

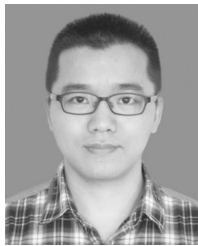
In the future work, there are several possible research directions. First, starting from this work, more metric diversification strategies can be studied to seek more opportunities in metric spaces for ensemble clustering of high-dimensional data. Second, different consensus methods for exploiting multilevel diversity in ensembles can be investigated to deal with different types of complex data. Last but not least, the diversification-and-fusion strategy may be generalized to other unsupervised learning tasks (such as unsupervised feature selection and unsupervised outlier detection) to alleviate the difficulties in parameter selection, metric selection, and model selection.

## REFERENCES

- [1] T. Li and C. Ding, "Weighted consensus clustering," in *Proc. SIAM Int. Conf. Data Min. (SDM)*, 2008, pp. 798–809.
- [2] N. Iam-On, T. Boongoen, S. Garrett, and C. Price, "A link-based approach to the cluster ensemble problem," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 12, pp. 2396–2409, Dec. 2011.
- [3] T. Wang, "CA-Tree: A hierarchical structure for efficient and scalable coassociation-based cluster ensembles," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 41, no. 3, pp. 686–698, Jun. 2011.
- [4] N. Li and L. J. Latecki, "Clustering aggregation as maximum-weight independent set," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2012, pp. 782–790.
- [5] L. Zheng, T. Li, and C. Ding, "A framework for hierarchical ensemble clustering," *ACM Trans. Knowl. Disc. Data*, vol. 9, no. 2, pp. 1–23, 2014.
- [6] L. Jing, K. Tian, and J. Z. Huang, "Stratified feature sampling method for ensemble clustering of high dimensional data," *Pattern Recognit.*, vol. 48, no. 11, pp. 3688–3702, 2015.
- [7] J. Wu, H. Liu, H. Xiong, J. Cao, and J. Chen, "K-means-based consensus clustering: A unified view," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 1, pp. 155–169, Jan. 2015.

- [8] D. Huang, J.-H. Lai, and C.-D. Wang, "Robust ensemble clustering using probability trajectories," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 5, pp. 1312–1326, May 2016.
- [9] D. Huang, J. Lai, and C.-D. Wang, "Ensemble clustering using factor graph," *Pattern Recognit.*, vol. 50, pp. 131–142, Feb. 2016.
- [10] H. Liu, M. Shao, S. Li, and Y. Fu, "Infinite ensemble clustering," *Data Min. Knowl. Disc.*, vol. 32, pp. 385–416, Mar. 2018.
- [11] H. Liu, J. Wu, T. Liu, D. Tao, and Y. Fu, "Spectral ensemble clustering via weighted  $k$ -means: Theoretical and practical evidence," *IEEE Trans. Knowl. Data Eng.*, vol. 29, no. 5, pp. 1129–1143, May 2017.
- [12] H. Liu, R. Zhao, H. Fang, F. Cheng, Y. Fu, and Y.-Y. Liu, "Entropy-based consensus clustering for patient stratification," *Bioinformatics*, vol. 33, no. 17, pp. 2691–2698, 2017.
- [13] Z. Yu *et al.*, "Semi-supervised ensemble clustering based on selected constraint projection," *IEEE Trans. Knowl. Data Eng.*, vol. 30, no. 12, pp. 2394–2407, Oct. 2018.
- [14] Y. Shi *et al.*, "Transfer clustering ensemble selection," *IEEE Trans. Cybern.*, vol. 50, no. 6, pp. 2872–2885, Jun. 2020.
- [15] D. Huang, C. D. Wang, and J. H. Lai, "Locally weighted ensemble clustering," *IEEE Trans. Cybern.*, vol. 48, no. 5, pp. 1460–1473, May 2018.
- [16] D. Huang, C.-D. Wang, H. Peng, J.-H. Lai, and C.-K. Kwok, "Enhanced ensemble clustering via fast propagation of cluster-wise similarities," *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. 51, no. 1, pp. 508–520, Jan. 2021.
- [17] D. Huang, C.-D. Wang, J.-S. Wu, J.-H. Lai, and C.-K. Kwok, "Ultra-scalable spectral clustering and ensemble clustering," *IEEE Trans. Knowl. Data Eng.*, vol. 32, no. 6, pp. 1212–1226, Jun. 2020.
- [18] L. Bai, J. Liang, and Y. Guo, "An ensemble clusterer of multiple fuzzy  $k$ -means clusterings to recognize arbitrarily shaped clusters," *IEEE Trans. Fuzzy Syst.*, vol. 26, no. 6, pp. 3524–3533, Dec. 2018.
- [19] L. Bai, J. Liang, H. Du, and Y. Guo, "An information-theoretical framework for cluster ensemble," *IEEE Trans. Knowl. Data Eng.*, vol. 31, no. 8, pp. 1464–1477, Aug. 2019.
- [20] Z. Tao, H. Liu, S. Li, Z. Ding, and Y. Fu, "Marginalized multiview ensemble clustering," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 2, pp. 600–611, Feb. 2020.
- [21] X. Yu, G. Yu, J. Wang, and C. Domeniconi, "Co-clustering ensembles based on multiple relevance measures," *IEEE Trans. Knowl. Data Eng.*, early access, Sep. 17, 2019, doi: [10.1109/TKDE.2019.2942029](https://doi.org/10.1109/TKDE.2019.2942029).
- [22] F. Li, Y. Qian, J. Wang, C. Dang, and L. Jing, "Clustering ensemble based on sample's stability," *Artif. Intell.*, vol. 273, pp. 37–55, Aug. 2019.
- [23] Z. Yu, H.-S. Wong, and H. Wang, "Graph-based consensus clustering for class discovery from gene expression data," *Bioinformatics*, vol. 23, no. 21, pp. 2888–2896, 2007.
- [24] Z. Yu, L. Li, J. Liu, J. Zhang, and G. Han, "Adaptive noise immune cluster ensemble using affinity propagation," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 12, pp. 3176–3189, Dec. 2015.
- [25] Z. Yu *et al.*, "Incremental semi-supervised clustering ensemble for high dimensional data clustering," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 3, pp. 701–714, Mar. 2016.
- [26] Z. Yu *et al.*, "Adaptive ensembling of semi-supervised clustering solutions," *IEEE Trans. Knowl. Data Eng.*, vol. 29, no. 8, pp. 1577–1590, Aug. 2017.
- [27] X. Z. Fern and C. E. Brodley, "Random projection for high dimensional data clustering: A cluster ensemble approach," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2003, pp. 186–193.
- [28] J. H. Lee, K. T. McDonnell, A. Zelenyuk, D. Imre, and K. Mueller, "A structure-based distance metric for high-dimensional space exploration with multidimensional scaling," *IEEE Trans. Vis. Comput. Graphics*, vol. 20, no. 3, pp. 351–364, Mar. 2014.
- [29] C. M. Hsu and M. S. Chen, "On the design and applicability of distance functions in high-dimensional data space," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 4, pp. 523–536, Apr. 2009.
- [30] A. L. N. Fred and A. K. Jain, "Combining multiple clusterings using evidence accumulation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 6, pp. 835–850, Jun. 2005.
- [31] J. Wu, H. Liu, H. Xiong, and J. Cao, "A theoretic framework of  $k$ -means-based consensus clustering," in *Proc. Int. Joint Conf. Artif. Intell.*, 2013, pp. 1799–1805.
- [32] C. Zhong, X. Yue, Z. Zhang, and J. Lei, "A clustering ensemble: Two-level-refined co-association matrix with path-based transformation," *Pattern Recognit.*, vol. 48, no. 8, pp. 2699–2709, 2015.
- [33] D. Huang, J.-H. Lai, and C.-D. Wang, "Combining multiple clusterings via crowd agreement estimation and multi-granularity link analysis," *Neurocomputing*, vol. 170, pp. 240–250, Dec. 2015.
- [34] Y. Fan, N. Li, C. Li, Z. Ma, L. J. Latecki, and K. Su, "Restart and random walk in local search for maximum vertex weight cliques with evaluations in clustering aggregation," in *Proc. Int. Joint Conf. Artif. Intell. (IJCAI)*, 2017, pp. 622–630.
- [35] M. Yousefnezhad, S. J. Huang, and D. Zhang, "WoCE: A framework for clustering ensemble by exploiting the wisdom of crowds theory," *IEEE Trans. Cybern.*, vol. 48, no. 2, pp. 486–499, Feb. 2018.
- [36] D. Huang, J.-H. Lai, C.-D. Wang, and P. C. Yuen, "Ensembling over-segmentations: From weak evidence to strong segmentation," *Neurocomputing*, vol. 207, pp. 416–427, Sep. 2016.
- [37] X. Wang, C. Yang, and J. Zhou, "Clustering aggregation by probability accumulation," *Pattern Recognit.*, vol. 42, no. 5, pp. 668–675, 2009.
- [38] J. Yi, T. Yang, R. Jin, and A. K. Jain, "Robust ensemble clustering by matrix completion," in *Proc. IEEE Int. Conf. Data Min. (ICDM)*, 2012, pp. 1176–1181.
- [39] A. Strehl and J. Ghosh, "Cluster ensembles: A knowledge reuse framework for combining multiple partitions," *J. Mach. Learn. Res.*, vol. 3, pp. 583–617, Feb. 2003.
- [40] X. Z. Fern and C. E. Brodley, "Solving cluster ensemble problems by bipartite graph partitioning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, p. 36, 2004.
- [41] A. Topchy, A. K. Jain, and W. Punch, "Clustering ensembles: Models of consensus and weak partitions," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 12, pp. 1866–1881, Dec. 2005.
- [42] L. Franek and X. Jiang, "Ensemble clustering by means of clustering embedding in vector spaces," *Pattern Recognit.*, vol. 47, no. 2, pp. 833–842, 2014.
- [43] A. K. Jain, "Data clustering: 50 years beyond  $k$ -means," *Pattern Recognit. Lett.*, vol. 31, no. 8, pp. 651–666, 2010.
- [44] G. Karypis and V. Kumar, "A fast and high quality multilevel scheme for partitioning irregular graphs," *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 359–392, 1998.
- [45] E. Weiszfeld and F. Plastria, "On the point for which the sum of the distances to  $n$  given points is minimum," *Ann. Oper. Res.*, vol. 167, no. 1, pp. 7–41, 2009.
- [46] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 498–519, Feb. 2001.
- [47] J. Chu, H. Wang, J. Liu, Z. Gong, and T. Li, "Unsupervised feature learning architecture with multi-clustering integration RBM," *IEEE Trans. Knowl. Data Eng.*, early access, Aug. 12, 2020, doi: [10.1109/TKDE.2020.3015959](https://doi.org/10.1109/TKDE.2020.3015959).
- [48] Q. Wang, P. C. Yuen, and G. Feng, "Semi-supervised metric learning via topology preserving multiple semi-supervised assumptions," *Pattern Recognit.*, vol. 46, no. 9, pp. 2576–2587, 2013.
- [49] F. Xiong, M. Kam, L. Hrebien, B. Wang, and Y. Qi, "Kernelized information-theoretic metric learning for cancer diagnosis using high-dimensional molecular profiling data," *ACM Trans. Knowl. Disc. Data*, vol. 10, no. 4, p. 38, 2016.
- [50] J. Li, A. J. Ma, and P. C. Yuen, "Semi-supervised region metric learning for person re-identification," *Int. J. Comput. Vis.*, vol. 126, pp. 855–874, Mar. 2018.
- [51] X. Liu *et al.*, "Absent multiple kernel learning algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 6, pp. 1303–1316, Jun. 2020.
- [52] X. Chen, W. Sun, B. Wang, Z. Li, X. Wang, and Y. Ye, "Spectral clustering of customer transaction data with a two-level subspace weighting method," *IEEE Trans. Cybern.*, vol. 49, no. 9, pp. 3230–3241, Sep. 2019.
- [53] H. Liu, Z. Han, Y.-S. Liu, and M. Gu, "Fast low-rank metric learning for large-scale and high-dimensional data," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 819–829.
- [54] B. Wang *et al.*, "Similarity network fusion for aggregating data types on a genomic scale," *Nat. Methods*, vol. 11, pp. 333–337, Jan. 2014.
- [55] U. von Luxburg, "A tutorial on spectral clustering," *Stat. Comput.*, vol. 17, no. 4, pp. 395–416, 2007.
- [56] Z. Li, X.-M. Wu, and S.-F. Chang, "Segmentation using superpixels: A bipartite graph partitioning approach," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2012, pp. 789–796.
- [57] M. C. de Souto, I. G. Costa, D. S. de Araujo, T. B. Ludermir, and A. Schliep, "Clustering cancer gene expression data: A comparative study," *BMC Bioinform.*, vol. 9, no. 1, p. 497, 2008.
- [58] K. Bache and M. Lichman. (2017). *UCI Machine Learning Repository*. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [59] S. A. Nene *et al.*, "Columbia object image library (COIL-20)," CUCS, New York, NY, USA, Rep. CUCS-005-96, 1996.
- [60] F. Nie, Z. Wang, R. Wang, and X. Li, "Submanifold-preserving discriminant analysis with an auto-optimized graph," *IEEE Trans. Cybern.*, vol. 50, no. 8, pp. 3682–3695, Aug. 2020.

- [61] S. Roweis. *Data for MATLAB Hackers*. [Online]. Available: <http://www.cs.nyu.edu/~roweis/data.html>
- [62] D. B. Graham and N. M. Allinson, "Characterising virtual eigensignatures for general purpose face recognition," in *Face Recognition*. Heidelberg, Germany: Springer, 1998, pp. 446–456.
- [63] M.-E. Nilsback and A. Zisserman, "A visual vocabulary for flower classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 2, 2006, pp. 1447–1454.
- [64] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance," *J. Mach. Learn. Res.*, vol. 11, no. 11, pp. 2837–2854, 2010.
- [65] P. Zhou, L. Du, X. Liu, Y.-D. Shen, M. Fan, and X. Li, "Self-paced clustering ensemble," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Apr. 20, 2020, doi: [10.1109/TNNLS.2020.2984814](https://doi.org/10.1109/TNNLS.2020.2984814).



**Dong Huang** (Member, IEEE) received the B.S. degree in computer science from the South China University of Technology, Guangzhou, China, in 2009, and the M.Sc. and Ph.D. degrees in computer science from Sun Yat-sen University, Guangzhou, in 2011 and 2015, respectively.

He joined South China Agricultural University, Guangzhou, in 2015, where he is currently working as an Associate Professor and the Deputy Head with the Department of Computer Science. From 2017 to 2018, he was a Visiting Fellow with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. He has published more than 50 papers in refereed journals and conferences, including IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, IEEE TRANSACTIONS ON CYBERNETICS, IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS, the *ACM Transactions on Knowledge Discovery From Data*, *Pattern Recognition*, *KBS*, *AAAI*, and *ICDM*. His current research interests include data mining and machine learning, and more specifically focus on ensemble clustering, multiview clustering, and large-scale clustering.



**Chang-Dong Wang** (Member, IEEE) received the B.S. degree in applied mathematics, the M.Sc. degree in computer science, and the Ph.D. degree in computer science from Sun Yat-sen University, Guangzhou, China, in 2008, 2010, and 2013, respectively.

He was a visiting student with the University of Illinois at Chicago, Chicago, IL, USA, from January 2012 to November 2012. He is currently an Associate Professor with the School of Computer Science and Engineering, Sun Yat-sen University. He has published more than 100 scientific papers in international journals and conferences, such as IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, IEEE TRANSACTIONS ON CYBERNETICS, the *ACM Transactions on Knowledge Discovery From Data*, *Pattern Recognition*, *Knowledge and Information Systems*, *IJCAI*, *AAAI*, *KDD*, *ICDM*, and *SDM*. His current research interests include machine learning and data mining.

Dr. Wang won the Honorable Mention for Best Research Paper Awards for his *ICDM* 2010 paper. He was Awarded Chinese Association for Artificial Intelligence Outstanding Dissertation in 2015.



**Jian-Huang Lai** (Senior Member, IEEE) received the M.Sc. degree in applied mathematics and the Ph.D. degree in mathematics from Sun Yat-sen University, Guangzhou, China, in 1989 and 1999, respectively.

In 1989, he joined Sun Yat-sen University as an Assistant Professor, where he is currently a Professor with the School of Computer Science and Engineering. He has published more than 200 scientific papers in the international journals and conferences on image processing and pattern recognition, such as the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, IEEE TRANSACTIONS ON NEURAL NETWORKS, IEEE TRANSACTIONS ON IMAGE PROCESSING, *Pattern Recognition*, *ICCV*, *CVPR*, *IJCAI*, *ICDM*, and *SDM*. His current research interests include the areas of digital image processing, pattern recognition, multimedia communication, and wavelet and its applications.

Prof. Lai serves as a Standing Member of the Image and Graphics Association of China, and also serves as a Standing Director of the Image and Graphics Association of Guangdong.



**Chee-Keong Kwoh** (Senior Member, IEEE) received the bachelor's (First Class) degree in electrical engineering and the master's degree in industrial system engineering from the National University of Singapore, Singapore, in 1987 and 1991, respectively, and the Ph.D. degree from the Imperial College of Science, Technology and Medicine, University of London, London, U.K., in 1995.

He has been with the School of Computer Science and Engineering, Nanyang Technological University, Singapore, since 1993, where he is the Programme Director of the M.Sc. in Bioinformatics Programme. His research interests include data mining, soft computing, and graph-based inference; applications areas include bioinformatics and biomedical engineering. He has done significant research work in his research areas and has published many quality international conferences and journal papers.

Dr. Kwoh is an editorial board member of the *International Journal of Data Mining and Bioinformatics*; *Scientific World Journal*; *Network Modeling and Analysis in Health Informatics and Bioinformatics*; *Theoretical Biology Insights*; and *Bioinformation*. He has been a Guest Editor for many journals, such as the *Journal of Mechanics in Medicine and Biology* and the *International Journal on Biomedical and Pharmaceutical Engineering*. He is a member of the Association for Medical and Bio-Informatics, Imperial College Alumni Association of Singapore. He has provided many services to professional bodies in Singapore and was conferred the Public Service Medal by the President of Singapore in 2008.