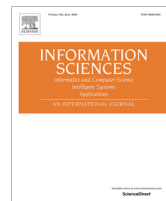




Contents lists available at ScienceDirect

Information Sciences

journal homepage: [www.elsevier.com/locate/ins](http://www.elsevier.com/locate/ins)



# Simultaneous learning coefficient matrix and affinity graph for multiple kernel clustering

Zhenwen Ren<sup>a,b</sup>, Haoyun Lei<sup>a</sup>, Quansen Sun<sup>b,\*</sup>, Chao Yang<sup>a,\*</sup>

<sup>a</sup> School of National Defense Science and Technology, Southwest University of Science and Technology, Mianyang 621010, China

<sup>b</sup> School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China

## ARTICLE INFO

### Article history:

Received 10 January 2020

Received in revised form 11 June 2020

Accepted 15 August 2020

Available online xxxx

### Keywords:

Subspace clustering

Kernel method

Multiple kernel learning

Affinity matrix

## ABSTRACT

Due to the effectiveness of handling non-linear data and avoiding kernel customization, multiple kernel clustering (MKC) has been widely investigated and achieved promising results for challenging non-linear clustering tasks. Generally, the existing MKC methods mainly consist of first learning a coefficient matrix by leveraging multiple kernel learning (MKL) and sample self-expressiveness property, and then constructing an affinity graph relying on this coefficient matrix to accomplish spectral clustering. However, the quality of the affinity matrix (graph) is largely determined by the learned coefficient matrix, thus the two independent steps are not conducive to learn an optimal affinity graph. To tackle this problem, in this paper, we propose a new MKC method that uses one-step learning scheme rather than two-step learning scheme to learn an affinity graph, termed SLMKC. Specifically, SLMKC bridges the relationship between the affinity matrix and coefficient matrix by an adaptive local structure learning strategy, so it can simultaneously learn both in a mutual promotion manner. Furthermore, a self-weighted MKL strategy is introduced to learn an optimal consensus kernel, which can avoid selecting a specific kernel function and tuning its associated parameters. Extensive experiments validate that our SLMKC outperforms the state-of-the-art MKC competitors significantly.

© 2020 Published by Elsevier Inc.

## 1. Introduction

Clustering is one of the most fundamental and important research topic in the realm of machine learning and data mining, which aims to partition the given unlabeled samples into disjoint groups [8]. To date, clustering problem has been widely researched in the last few decades and many typical methods have been surged in several studies [10,45,10,21]. In general, the existing clustering methods can be roughly classified into five categories, including iterative based methods [11], statistical based methods [12], algebraic based methods [20], matrix factorization based methods [26], and spectral clustering based methods [20].

Among these methods, spectral clustering based methods have become extremely popular for challenging clustering tasks, thanks to their capabilities to handle arbitrarily shaped clusters and their well-defined mathematical principles. Traditionally, such type of methods splits the clustering process into two independent steps. In the first one, an affinity graph/-matrix is constructed by exploiting the given unlabeled data. In the second one, spectral clustering algorithm is employed to obtain the final clustering assignments relying on this affinity graph. There is no arguing that the first step is the most impor-

\* Corresponding authors.

E-mail addresses: [sunquansen@njust.edu.cn](mailto:sunquansen@njust.edu.cn) (Q. Sun), [ycho1983@126.com](mailto:ycho1983@126.com) (C. Yang).

<https://doi.org/10.1016/j.ins.2020.08.056>

0020-0255/© 2020 Published by Elsevier Inc.

tant one since the performance of spectral clustering largely depends on constructing a desired affinity graph. Although these spectral clustering based methods can produce good performance, there are two great challenges that limit their applications in real-world settings: (1) how to handle non-linear data when they are not strictly collected from linearly independent subspaces; (2) how to learn an optimal affinity graph for spectral clustering purpose.

In order to tackle the challenging problem of how to handle non-linear data, single kernel learning (SKL) method has been in common use [39,14]. Theoretically, SKL maps the non-linear data from the original space into the Reproducing Kernel Hilbert Space (RKHS) where the linear property of the data is satisfied [2]. In practice, the kernel trick [3] is applied to avoid the explicit mapping, which is often computationally cheaper than the explicit computation of the projected coordinates. Although promising performance has been reported in a number of tasks, SKL methods have to require the user to select and tune a pre-defined kernel [16] since their performance is largely determined by the choice of kernel function. This is not user-friendly since the most suitable kernel function for a specific task or dataset is usually challenging to decide. Therefore, how to select an appropriate pre-defined single kernel function and how to tune its associated kernel parameters are worthy to study. In recent years, a powerful and flexible learning model, multiple kernel learning (MKL) [42,7,46,35], has attracted more and more attention to avoid kernel customization by constructing a consensus kernel from multiple candidate base kernels. More importantly, MKL has the great potential to fully exploit complementary information between these candidate kernels. Generally, MKL has been widely used in clustering [14], object recognition [7], and classification [31]. For illustrative purposes, the thoughts of the SKL and MKL are illustrated in Fig. 1(a) and (b), respectively.

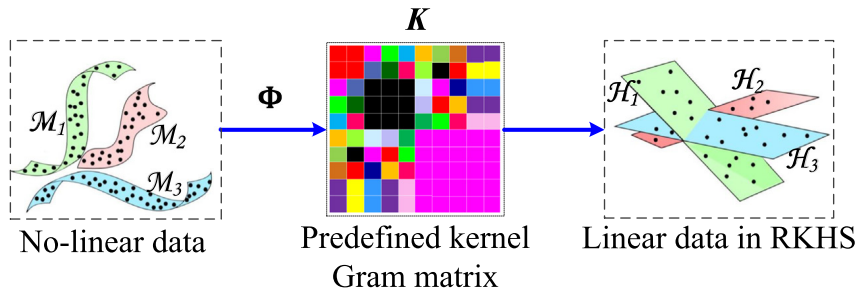
In addition, to learn a better affinity graph for spectral clustering [18], many scholars have proposed different methods to learn an expected affinity graph [39,29]. Overall, the mainstream technologies can be typically divided into three main categories. The first category is to construct a predefined similarity graph as affinity graph, such as complete graph and  $k$ -nearest graph. Usually, the complete graph is constructed based on binary (0 – 1) similarity, cosine similarity, or gaussian kernel similarity [41], while  $k$ -nearest graph selects its neighbor nodes by using  $k$  Nearest Neighbor algorithm. The second category is the adaptive neighbors graph learning [33,17], which builds a graph by assigning a probability for each sample as the neighborhood of another samples. Accordingly, the homogeneous samples have high affinity (similarity) values, while those heterogeneous samples have low affinity values. The third category is based on the sample self-expressiveness property [29], which reconstructs every data point by a linear combination of all other data points, and then produces a coefficient matrix that is used to construct an affinity graph. To the best of our knowledge, the self-expressiveness property based technology has been the most widely studied and obtained impressive performance in a wide variety of applications. This paper focuses on the sample self-expressiveness property [29] to construct an affinity graph used for spectral clustering.

Combining both MKL and self-expressiveness, many multiple kernel clustering (MKC) methods have been proposed [27,14,9,13,16,19]. Overall, these existing MKC methods usually first learn a self-expressiveness coefficient matrix and then construct an affinity matrix to accomplish spectral clustering independently. However, the quality of affinity matrix is largely determined by the coefficient matrix. Obviously, such two-step learning scheme is not conducive to learn an optimal affinity matrix, thus impairs the final clustering performance greatly.

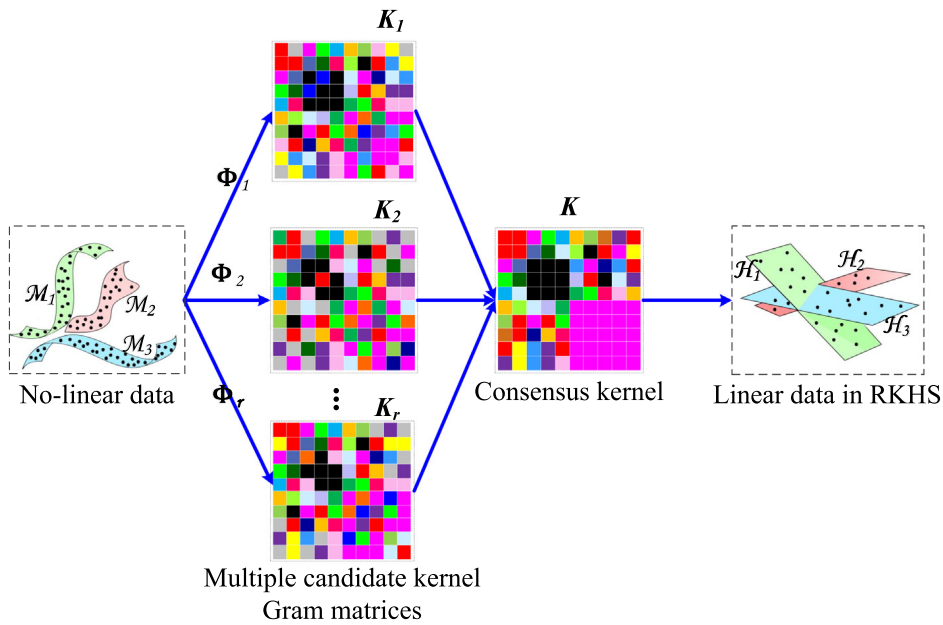
To effectively handle non-linear data and learn an optimal affinity matrix for challenging clustering, this paper proposes a novel method, namely Simultaneous Learning Self-expressiveness Coefficients and Affinity Matrix Multiple Kernel Clustering (SLMKC). Different from the existing MKC methods, SLMKC proposes a unified optimization objective function that uses one-step learning scheme rather than two-step learning scheme to learn an optimal affinity matrix in multiple kernels scene. Specifically, the relationship between the coefficient matrix and the affinity matrix is bridged via adaptive local structure learning. Meanwhile, SLMKC fuses multiple candidate base kernels to learn an optimal consensus kernel. Furthermore, to ensure that the learned affinity graph contains exactly  $c$  connected components, a constrained Laplacian rank term is introduced; to remove the negative effects of redundancy affinity graph edges, a top- $k$  probabilistic neighbors selection strategy is also employed. Consequently, a high-quality affinity graph is subsequently obtained. In summary, the contributions of this paper are summarized below.

- In order to more effectively handle the non-linear data and learn an optimal affinity graph for clustering purpose, SLMKC can address the difficult problem of how to select an appropriate single kernel function and tune the associated parameters of the selected kernel function.
- Different from existing MKC methods, which obey a two-step learning scheme to learn the coefficient matrix and the affinity matrix independently, SLMKC proposes a unified optimization model to simultaneously learn the consensus kernel, the self-expressiveness coefficient matrix, and the affinity matrix. Thus, SLMKC can address the challenging problem of performance degradation.
- Comprehensive experimental results accomplished on nine benchmark datasets (i.e., six image datasets and three text datasets) demonstrate that SLMKC substantially improves the clustering performance and computational cost, compared with the existing state-of-the-art MKC methods.

The rest of this paper is structured as follows. Section 2 revisits the related self-expressiveness learning model and MKC methods for clustering. Then, the proposed SLMKC method and its optimization algorithm and complexity are given in Section 3. Finally, Section 4 and Section 5 offer the sufficient experimental results and conclusion, respectively.



(a) Illustration of SKL based methods.



(b) Illustration of MKL based methods

**Fig. 1.** Illustrations of SKL based methods and MKL based methods. Different from SKL based methods, which use only one kernel function to map non-linear data from the original space to RKHS, MKL based methods integrate multiple candidate kernels to learn a consensus kernel. By doing so, MKL can avoid selecting and tune a single pre-defined kernel function; meanwhile, can strongly promote to effectively use the complementary information of these candidate kernels.

## 2. Related works

In this section, we first present the notations used throughout the paper, and then briefly revisit the self-expressiveness learning model and the related MKC methods.

### 2.1. Notations

Throughout this paper, we denote scalars by normal italic letters, matrices by boldface uppercase letters, and vectors by boldface lowercase letters, respectively. Some important notations are summarized in Table 1. Note here that the data matrix drawn from  $c$  independent subspaces is denoted by  $\mathbf{X} = \{\mathbf{X}_i | \mathbf{X}_i \in \mathbb{R}^{d \times n_i}\}_{i=1}^c$ , where  $\mathbf{X}_i = [x_{i1}, x_{i2}, \dots, x_{in_i}]$  and  $\sum_{i=1}^c n_i = n$ .

### 2.2. Self-expressiveness learning model

As we know, self-expressiveness property is derived based on the fact that every sample can be represented by a linear combination of all other samples [29], and subspace clustering methods rely on such property to construct a self-expressiveness coefficient matrix. Mathematically, the self-expressiveness learning model is defined as follows:

**Table 1**  
Notations used throughout the paper.

Notation	Definition
$\mathbf{X} \in \mathbb{R}^{d \times n}$	The column-wise data matrix
$x_j$ or $\mathbf{X}_j$	The $j$ th column of $\mathbf{X}$
$x_{ij}$ or $\mathbf{X}_{ij}$	The $(i,j)$ th entry of $\mathbf{X}$
$\mathbf{Z} \in \mathbb{R}^{n \times n}$	The self-expressiveness coefficient matrix
$\mathbf{A} \in \mathbb{R}^{n \times n}$	The affinity matrix
$\mathbf{P} \in \mathbb{R}^{n \times c}$	The clustering indicator matrix
$\mathbf{K}^i \in \mathbb{R}^{n \times n}$	The $i$ th candidate base kernel
$\mathbf{K} \in \mathbb{R}^{n \times n}$	The consensus kernel
$\mathbf{1}$	The all-one column vector
$\mathbf{I}$	The identity matrix
$\ \cdot\ _F$	The Frobenius norm
$\text{Tr}(\cdot)$	The trace operator
$\text{diag}(\cdot)$	Vectorize the diagonal elements of a matrix
$\text{Diag}(\cdot)$	Rearrange a diagonal matrix to a vector
$[\cdot]_+$	Give the nonnegative part of a matrix or vector

$$\min_{\mathbf{Z}} \frac{1}{2} \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \lambda \mathcal{R}(\mathbf{Z}) \quad (1)$$

where  $\lambda > 0$  is a trade-off parameter,  $\mathbf{Z}$  is the self-expressiveness coefficient matrix, and  $\mathcal{R}(\mathbf{Z})$  is a specific regularization term or constraint term [24], such as sparse constraint [10], low-rank constraint [25], adaptive weighted representation [43], continuous label learning [18], block diagonal representation [29], and structure preserving [46], etc. After obtaining matrix  $\mathbf{Z}$ , the undirected affinity matrix (also known as affinity graph) for spectral clustering is constructed by  $(|\mathbf{Z}|^T + |\mathbf{Z}|)/2$ .

The main challenge of problem (1) is that it can not effectively handle non-linear data. By taking this point into consideration, the non-linear data always can be mapped into RKHS through a kernel function. Then, based upon the aforementioned self-expressiveness property, an affinity matrix in RKHS can be learned by upgrading problem (1) to

$$\begin{aligned} & \min_{\mathbf{Z}} \frac{1}{2} \|\phi(\mathbf{X}) - \phi(\mathbf{X})\mathbf{Z}\|_F^2 + \lambda \mathcal{R}(\mathbf{Z}) \\ &= \min_{\mathbf{Z}} \text{Tr}(\phi(\mathbf{X})^T \phi(\mathbf{X}) - 2\phi(\mathbf{X})^T \phi(\mathbf{X})\mathbf{Z} + \mathbf{Z}^T \phi(\mathbf{X})^T \phi(\mathbf{X})\mathbf{Z}) + \lambda \mathcal{R}(\mathbf{Z}) \\ &= \min_{\mathbf{Z}} \text{Tr}(\mathbf{K} - 2\mathbf{KZ} + \mathbf{Z}^T \mathbf{KZ}) + \lambda \mathcal{R}(\mathbf{Z}) \end{aligned} \quad (2)$$

where  $\phi$  is a specific mapping function and  $\mathbf{K}$  is the kernel Gram matrix. Note that the kernel trick is used here. Such a model is also called as kernel self-expressiveness learning model. For SKL methods,  $\mathbf{K}$  is fixed, while for MKL methods,  $\mathbf{K}$  is an unknown consensus kernel matrix that is desired to be learned from multiple candidate kernels.

### 2.3. Related MKC methods

Based on self-expressiveness learning model [46],  $k$ -means algorithm [14], or graph learning technology [17], many previous MKC methods have made attempts to address the non-linear clustering problem [36,38]. For example, multiple kernel  $k$ -means (MKKM) [14] extends the traditional  $k$ -means algorithm into a multiple kernel setting for challenging clustering. To enhance robustness of MKKM, an extended version of MKKM, robust multiple kernel  $k$ -means (RMKKM) [9], has been proposed to simultaneously find the best clustering indicator label and the optimal combination of the given multiple kernels. Regarding the similarity between each candidate kernel matrix and the resulting affinity matrix, affinity aggregation for spectral clustering (AASC) [13] can be deemed as a MKL version of spectral clustering, which jointly replaces the single affinity matrix with multiple kernel matrices for spectral clustering. For learning a better consensus kernel, spectral clustering with multiple kernels (SCMK) [18] constructs a suitable consensus kernel from a linear combination of multiple candidate kernels. Eventually, SCMK proposes to learn kernel weights, cluster indicator, and similarity simultaneously. By considering the structure character of the consensus kernel, low-rank kernel learning graph based clustering (LKGr) [19] and sparse kernel learning graph based clustering (LKGs) [19] seek an optimal consensus kernel from a linear combination of candidate kernels by imposing a low-rank constraint and a sparse constraint on this consensus kernel matrix, respectively. Later, local structural graph and low-rank consensus multiple kernel learning (LLMKL) [35] suggests that the low-rank kernel used in LKGr is actually a Frobenius norm, and thus LLMKL proposes to use a substitute matrix to replace the consensus kernel matrix, and then imposes a low-rank constraint on this substitute matrix rather than consensus kernel matrix itself. Different from SCMK, LKGr, and LKGs, which use linear weighted strategy to learn a consensus kernel, self-weighted multiple kernel learning (SMKL) [16] introduces a novel self-weighted MKL model for clustering and semi-supervised learning, which learns a suitable consensus kernel based on the assumption that the consensus kernel is a neighbor of all candidate kernels

and the important kernel should receive a relatively large weight, and vice versa. By aligning each base kernel with the corresponding local optimal similarity matrix, local kernel alignment maximization (LKAMKC) [22] improves the clustering performance for multiple kernel clustering. Based on LKAMKC, global and local structure alignment model for multiple kernel clustering (GLSAMKC) [40] well considers both the alignment between the global and local structure of data with the same optimal similarity matrix. Neighbor-Kernel MKC [48] proposes an effective neighbor-kernel-based multiple kernels clustering method by considering the intrinsic neighborhood structure among base kernels. Recently, to improve the robustness and the quality of the consensus kernel, joint robust multiple kernel subspace clustering (JMKSC) [46] proposes an effective method for data clustering, which combines the subspace self-expressiveness, a correntropy induced kernel weighting strategy, and the block diagonal regularizer into a unified objective function, as a result, it can efficiently handle non-linear data and defy the impulsive noise and non-Gaussian noise.

However, these existing spectral clustering based MKC methods usually adopt a two-step learning scheme, which have to face the problem of not being able to learn the optimal affinity graph.

### 3. Simultaneous learning scheme for multiple kernel clustering

We first introduce the formulation of the proposed method and then present an effective algorithm to solve the final objective function.

#### 3.1. Problem formulation

Based on the kernel self-expressiveness learning model and considering  $\mathcal{R}(\mathbf{Z}) = \|\mathbf{Z}\|_F^2$ , an affinity graph can be learned in kernel space, i.e.,

$$\min_{\mathbf{Z}} \text{Tr}(\mathbf{K} - 2\mathbf{KZ} + \mathbf{Z}^T \mathbf{KZ}) + \lambda \|\mathbf{Z}\|_F^2 \quad (3)$$

where the regularization term,  $\|\mathbf{Z}\|_F^2$ , can avoid the trivial solution (i.e.,  $\mathbf{Z} = \mathbf{I}$ ) and avoid scale change [23], and meanwhile enhance the convexness [29] and grouping effect [30]. Existing MKC methods always use problem (3) to learn a self-expressiveness coefficient matrix  $\mathbf{Z}$  and require an additional step to construct an affinity matrix. Whereas, these is a drawback (presented in the Section 1) that needs to be defeated.

Now, we aim to bridge the relationship between coefficient matrix and affinity matrix to learn both simultaneously. As we know, the self-expressiveness coefficient matrix  $\mathbf{Z}$  can be deemed as the substitute of data matrix  $\mathbf{X}$  because each entry of  $\mathbf{Z}$  can quantify the similarity between two samples in  $\mathbf{X}$ ; that is to say, if two samples are close to each other in original space, the new representations of the two samples in new space must be similar to each other, too. Additionally, the adaptive local structure learning (ALSL) [47,15],  $\min_{\mathbf{G}} \|\mathbf{X}_i - \mathbf{X}_j\|_2^2 \mathbf{G}_{ij}$ , is a novel graph learning technology to learn a similarity graph  $\mathbf{G}$  that can preserve the local manifold structure of data, which is based on the fact that if two samples (e.g.,  $\mathbf{X}_i$  and  $\mathbf{X}_j$ ) are close to each other, they will have higher similarity (e.g.,  $\mathbf{G}_{ij}$ ), and vice versa. Motivated by ALSL, if the similarity between  $\mathbf{Z}_i$  and  $\mathbf{Z}_j$  is smaller, the corresponding affinity value (also known as probability) between them is higher, and vice versa. Consequently, we bridge the relationship between affinity matrix  $\mathbf{A}$  and coefficient matrix  $\mathbf{Z}$  by solving the problem defined as follows:

$$\min_{\mathbf{A}} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{Z}_i - \mathbf{Z}_j\|_F^2 \mathbf{A}_{ij} \quad (4)$$

Moreover, many recent studies [37] have shown that sparse representation is robust to noise and outliers. To improve robustness, we integrate sparse representation into problem (4) and then have

$$\min_{\mathbf{A}} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{Z}_i - \mathbf{Z}_j\|_F^2 \mathbf{A}_{ij} + \lambda \sum_{i=1}^n \|\mathbf{A}_i\|_1 \quad (5)$$

where  $\|\mathbf{A}_i\|_1$  encourages to obtain a sparse  $\mathbf{A}_i$ . Indeed, when  $\mathbf{A}_i$  is normalized by  $\mathbf{A}_i^T \mathbf{1} = 1$ , it exactly forces term  $\sum_{i=1}^n \|\mathbf{A}_i\|_1$  to be constant. That is, the constraint  $\mathbf{A}_i^T \mathbf{1} = 1$  can be deemed as the sparse constraint on  $\mathbf{A}_i$ ; meanwhile, the probabilistic property also should be held (i.e.,  $\mathbf{A}_i \succeq 0$  and  $\mathbf{A}_i^T \mathbf{1} = 1$ ). Then, problem (5) is transformed into

$$\min_{\mathbf{A}} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{Z}_i - \mathbf{Z}_j\|_F^2 \mathbf{A}_{ij} \quad \text{s.t. } \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A}_i \succeq 0 \quad (6)$$

The proof that problem (6) can lead to a sparse matrix  $\mathbf{A}$  is given in Proposition 1. Proposition 1 Problem (6) can lead to a sparse matrix  $\mathbf{A}$ .



Define that the  $(i, j)$ -th entry of  $\mathbf{Q}$  is denoted as  $\mathbf{Q}_{ij} = \|\mathbf{Z}_i - \mathbf{Z}_j\|_F^2$ , then problem (6) can be divided into  $n$  independent subproblems algebraically, i.e.,  $\min_{\mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A}_i \succeq 0} \|\mathbf{Q}_i \odot \mathbf{A}_i\|_1$ . Then, it is easy to see that each subproblem can be considered as a special case of sparse representation [28]. Obviously, this is to say that problem (6) can lead to a sparse matrix  $\mathbf{A}$ . The proof is completed. ■.

However, problem (6) inevitably has trivial solution. i.e., only one edge of graph  $\mathbf{A}$  has the affinity value 1, while all the other edges have the affinity value 0. Accordingly, a regularization term,  $\|\mathbf{A}_i\|_2^2$ , is appended to problem (6) to prevent from trivial solution, and then problem (6) can be reformulated as

$$\min_{\mathbf{A}} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{Z}_i - \mathbf{Z}_j\|_F^2 \mathbf{A}_{ij} + \lambda \sum_{i=1}^n \|\mathbf{A}_i\|_2^2 \quad \text{s.t.} \quad \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A}_i \succeq 0 \quad (7)$$

Mathematically, problem (7) can be transformed into

$$\min_{\mathbf{A}} \text{Tr}(\mathbf{Z} \mathbf{L}_A \mathbf{Z}^T) + \lambda \|\mathbf{A}\|_F^2, \quad \text{s.t.} \quad \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A}_i \succeq 0 \quad (8)$$

where  $\mathbf{L}_A$  is the Laplacian matrix corresponding to  $\mathbf{A}$ , i.e.,  $\mathbf{L}_A = \text{Diag}(\mathbf{A} \mathbf{1}) - \mathbf{A}$ . By using problem (8), the affinity graph  $\mathbf{A}$  is learned, however, this affinity graph may contain too many redundant edges. Subsequently, a top- $k$  probabilistic neighbors selection strategy is introduced to preserve the important edges and remove the unimportant edges so as to obtain a sparse affinity graph and to reduce computing overhead. The detail of such a strategy is given in Section 3.2.

Ideally, it is often more desirable to achieve an affinity graph that consists of exactly  $c$  graph connected components (sub-graphs) [33,47], where  $c$  is the numbers of clusters. That is, the matrix  $\mathbf{A}$  is block diagonal matrix with proper permutations associated to  $c$ . Obviously, it is difficult to directly obtain the expected affinity graph since the graph produced by problem (8) is always a strongly connected graph. Fortunately, Theorem 1 can tackle this issue. Theorem 1 [29] For any affinity graph  $\mathbf{A}$ , the multiplicity  $c$  of the eigenvalue 0 of the corresponding Laplacian matrix  $\mathbf{L}_A$  equals the number of connected components in  $\mathbf{A}$ .

Motivated by such theorem, it can be proved that if the Laplacian rank constraint,  $\text{rank}(\mathbf{L}_A) = n - c$ , is employed, the affinity graph  $\mathbf{A}$  will contain exact  $c$  connected components, i.e., the Laplacian rank constraint is very necessary to help to partition the learned affinity matrix into  $c$  connected sub-graphs. Here, in order to give more freedom to the affinity matrix  $\mathbf{A}$ , instead of using affinity matrix  $\mathbf{A}$ , we impose the Laplacian rank constraint on the Laplacian matrix  $\mathbf{L}_Z$  produced by coefficient matrix  $\mathbf{Z}$ , since  $\mathbf{Z} = (\mathbf{Z}^T + \mathbf{Z})/2$  also can be deemed as an affinity matrix. Mathematically, according to rank theory and Fan's theorem [34], we have

$$\text{rank}(\mathbf{L}_Z) = n - c \Rightarrow \min_{\mathbf{P} \in \mathbb{R}^{n \times c}, \mathbf{P}^T \mathbf{P} = \mathbf{I}} \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) \quad (9)$$

where  $\mathbf{P}^T = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n\}$  is the cluster indicator matrix.

In addition, SKL is difficult to select a pre-defined kernel function and tuning its kernel parameters [36]. Inspired by the auto-weighted multi-view learning model [32], we introduce a nearest neighbor based auto-weighted MKL model by integrating  $r$  candidate base kernels  $\{\mathbf{K}^i\}_{i=1}^r$  so as to learn an optimal consensus kernel  $\mathbf{K}$ , i.e.,

$$\min_{\mathbf{K}} \sum_{i=1}^r w_i \|\mathbf{K}^i - \mathbf{K}\|_F^2 \quad (10)$$

where  $w_i$  indicates the weight of the  $i$ -th candidate kernel  $\mathbf{K}^i$ . For ease of exploration, the detailed derivation of  $\mathbf{w} = [w_1, w_2, \dots, w_r]$  is presented in Section 3.3.

Regarding the thoughts mentioned above, we formulate a unified objective function to learn an optimal affinity graph for MKC, which is formulated as

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{K}, \mathbf{A}, \mathbf{P}} & \text{Tr}(\mathbf{K} - 2\mathbf{K}\mathbf{Z} + \mathbf{Z}^T \mathbf{K}\mathbf{Z}) + \lambda_1 (\|\mathbf{Z}\|_F^2 + \|\mathbf{A}\|_F^2) + \lambda_2 \text{Tr}(\mathbf{Z} \mathbf{L}_A \mathbf{Z}^T) + \lambda_3 \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) + \lambda_4 \sum_{i=1}^r w_i \|\mathbf{K}^i - \mathbf{K}\|_F^2 \\ \text{s.t.} & \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A} \succeq 0, \mathbf{P} \in \mathbb{R}^{n \times c}, \mathbf{P}^T \mathbf{P} = \mathbf{I} \end{aligned} \quad (11)$$

where  $\lambda_1, \lambda_2, \lambda_3$ , and  $\lambda_4$  are positive trade-off parameters. In fact, only two parameters  $\lambda_2$  and  $\lambda_3$  are need to be tuned.

### 3.2. Optimization

In this section, we solve the variables  $\mathbf{Z}, \mathbf{A}, \mathbf{P}$ , and  $\mathbf{K}$  in problem (11) based on the alternating direction method of multipliers (ADMM) algorithm [37]. Due to the requirement of ADMM, one auxiliary variable  $\mathbf{B}$  is first introduced to make the problem (11) separable, then we have

$$\min_{\mathbf{Z}, \mathbf{A}, \mathbf{B}, \mathbf{P}, \mathbf{K}} \text{Tr}(\mathbf{K} - 2\mathbf{KZ} + \mathbf{Z}^T \mathbf{KZ}) + \lambda_1 (\|\mathbf{Z}\|_F^2 + \|\mathbf{A}\|_F^2) + \lambda_2 \text{Tr}(\mathbf{B} \mathbf{L}_A \mathbf{B}^T) + \lambda_3 \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) + \lambda_4 \sum_i w_i \|\mathbf{K}^i - \mathbf{K}\|_F^2 \quad (12)$$

$$\text{s.t. } \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A} \succeq 0, \mathbf{B} = \mathbf{Z}, \mathbf{P}^T \mathbf{P} = \mathbf{I}$$

The augmented Lagrangian function of problem (12) is given by

$$\begin{aligned} \mathcal{L} \{ \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A} \succeq 0, \mathbf{P} \in \mathbb{R}^{n \times c}, \mathbf{P}^T \mathbf{P} = \mathbf{I} \} \\ = \text{Tr}(\mathbf{K} - 2\mathbf{KZ} + \mathbf{Z}^T \mathbf{KZ}) + \lambda_1 (\|\mathbf{Z}\|_F^2 + \|\mathbf{A}\|_F^2) + \lambda_2 \text{Tr}(\mathbf{B} \mathbf{L}_A \mathbf{B}^T) + \lambda_3 \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) \\ + \lambda_4 \sum_i w_i \|\mathbf{K}^i - \mathbf{K}\|_F^2 + \frac{\mu}{2} \|\mathbf{B} - \mathbf{Z} + \frac{\mathbf{Y}}{\mu}\|_F^2 \end{aligned} \quad (13)$$

where  $\mathbf{Y}$  and  $\mu$  are the Lagrangian multiplier matrix and the penalty parameter, respectively.

**(1) Update  $\mathbf{Z}$  as given  $\mathbf{A}, \mathbf{B}, \mathbf{P}$  and  $\mathbf{K}$ :**

The subproblem with respect to self-expressiveness coefficient matrix  $\mathbf{Z}$  can be written as

$$\min_{\mathbf{Z}} \text{Tr}(\mathbf{K} - 2\mathbf{KZ} + \mathbf{Z}^T \mathbf{KZ}) + \lambda_1 \|\mathbf{Z}\|_F^2 + \lambda_3 \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) + \frac{\mu}{2} \|\mathbf{B} - \mathbf{Z} + \frac{\mathbf{Y}}{\mu}\|_F^2 \quad (14)$$

Due to  $\text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) = \frac{1}{2} \text{Tr}(\mathbf{ZQ})$  (the detailed proof is in Appendix A), where  $\mathbf{Q}$  is the distance matrix mentioned at Section 3.1, problem (14) has a closed-form solution as follows:

$$\Rightarrow \mathbf{Z}^* = (2\mathbf{K} + 2\lambda_1 \mathbf{I} + \mu \mathbf{I})^{-1} \left( 2\mathbf{K} + \mu \mathbf{B} + \mathbf{Y} - \frac{\lambda_3}{2} \mathbf{Q}^T \right) \quad (15)$$

**(2) Update  $\mathbf{B}$  as given  $\mathbf{Z}, \mathbf{A}, \mathbf{P}$  and  $\mathbf{K}$ :**

The updating of auxiliary matrix  $\mathbf{B}$  is as follows:

$$\min_{\mathbf{B}} \lambda_2 \text{Tr}(\mathbf{B} \mathbf{L}_A \mathbf{B}^T) + \frac{\mu}{2} \|\mathbf{B} - \mathbf{Z} + \frac{\mathbf{Y}}{\mu}\|_F^2 \quad (16)$$

whose solution is provided by

$$\Rightarrow \mathbf{B}^* = (\mu \mathbf{Z} - \mathbf{Y})(2\lambda_2 \mathbf{L}_A + \mu \mathbf{I})^{-1} \quad (17)$$

**(3) Update  $\mathbf{A}$  as given  $\mathbf{Z}, \mathbf{B}, \mathbf{P}$  and  $\mathbf{K}$ :**

The updating of affinity matrix/graph  $\mathbf{A}$  is as follows:

$$\min_{\mathbf{A}} \lambda_2 \text{Tr}(\mathbf{B} \mathbf{L}_A \mathbf{B}^T) + \lambda_1 \|\mathbf{A}\|_F^2 \quad \text{s.t. } \forall i, \mathbf{A}_i^T \mathbf{1} = 1, \mathbf{A} \succeq 0 \quad (18)$$

This problem can be divide into a series of independent subproblems. For  $\mathbf{A}_i$ , the subproblem is formulated as

$$\mathbf{A}_i = \arg \min \|\mathbf{A}_i + \mathbf{D}_i\|_F^2 \quad (19)$$

where  $\mathbf{D}_{ij} = \frac{\lambda_2 \|\mathbf{B}_i - \mathbf{B}_j\|_F^2}{4\lambda_1}$ . The solution of problem (19) is derived by a top- $k$  probabilistic neighbors selection strategy, given by

$$\Rightarrow \mathbf{A}_i^* = \left[ \frac{1 + \sum_{j=1}^k \mathbf{D}_{ij}^-}{k} \mathbf{1} - \mathbf{D}_i \right]_+ \quad (20)$$

where  $k$  is used to control the number of important edges connected to a vertex in affinity graph  $\mathbf{A}$ , and  $\mathbf{D}_i^-$  is the  $\mathbf{D}_i$  with all elements sorted by descending order. For clarity and completeness, the detailed solution steps of (19) are presented in Appendix B. By doing so, from the matrix perspective, the  $i$ th column  $\mathbf{A}_i$  of affinity matrix  $\mathbf{A}$  has  $k$  number of nonzero entries; from the graph perspective, there are  $k$  number of important vertices connecting to the  $i$ -th vertex in the affinity graph  $\mathbf{A}$ .

**Update  $\mathbf{P}$  as given  $\mathbf{Z}, \mathbf{A}, \mathbf{B}$  and  $\mathbf{K}$ :**

The updating of clustering indicator matrix  $\mathbf{P}$  is as follows:

$$\min_{\mathbf{P} \in \mathbb{R}^{n \times c}} \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) \quad \text{s.t. } \mathbf{P}^T \mathbf{P} = \mathbf{I} \quad (21)$$

whose solution is formed by the  $c$  eigenvectors of  $\mathbf{L}_Z$  corresponding to its  $c$  smallest eigenvalues. **(5) Update  $\mathbf{K}$  as given  $\mathbf{Z}, \mathbf{A}, \mathbf{B}$  and  $\mathbf{P}$ :**

The subproblem for updating consensus kernel  $\mathbf{K}$  is as follows:

$$\min_{\mathbf{K}} \text{Tr}(\mathbf{K} - 2\mathbf{K}\mathbf{Z} + \mathbf{Z}^T\mathbf{K}\mathbf{Z}) + \lambda_4 \sum_i^r w_i \|\mathbf{K}^i - \mathbf{K}\|_F^2 \quad (22)$$

It should be noted that the weight value  $w_i$  is used to control the contribution of the  $i$ -th candidate kernel for generating the consensual  $\mathbf{K}$ .

Taking the partial derivative of (22) with respect to  $\mathbf{K}$  and setting the derivative to zero [5], the closed-form solution of (22) is given by

$$\Rightarrow \mathbf{K} = \frac{-\mathbf{I} - \mathbf{Z}\mathbf{Z}^T + 2\mathbf{Z}^T + 2\lambda_4 \sum_i w_i \mathbf{K}^i}{2\lambda_4 \sum_i w_i} \quad (23)$$

The procedure of solving the problem (11) terminates when  $\|\mathbf{B} - \mathbf{Z}\|_\infty \leq \epsilon$  with  $\epsilon = 1e-5$  or the number of iterators exceeds  $mit = 1e3$ . The pseudo-code is depicted as in Algorithm 1. By using Algorithm 1, both the affinity matrix  $\mathbf{A}$  and self-expressiveness coefficient matrix  $\mathbf{Z}$  can be achieved, and then we fuse the two matrices and calculate the final balanced affinity graph by

$$\mathbf{A} = \frac{|\mathbf{Z} \odot \mathbf{A}| + |(\mathbf{Z} \odot \mathbf{A})^T|}{2} \quad (24)$$

where  $\odot$  denotes the Hadamard element-wise product. After obtaining the undirected graph  $\mathbf{A}$ , we perform spectral clustering [46] onto this fused graph  $\mathbf{A}$  to obtain final clustering results.

---

**Algorithm 1.** Solve the proposed objective model (11) for MKC.

---

Input:  $r$  base kernels  $\{\mathbf{K}^i\}_{i=1}^r$ , parameters  $\lambda_2$  and  $\lambda_3$ , and the number of neighbors  $k$ .

1: **Initialize for ADMM:**  $\rho = 2$ ,  $\mu_{max} = 1e5$ ,  $\epsilon = 1e-5$ ,  $t = 1$ , and  $mit = 1e3$ .

2: **Initialize for model:**  $\mathbf{Z} = \mathbf{I}$ ,  $\mathbf{K} = \frac{1}{r} \sum_{i=1}^r \mathbf{K}^i$ ,  $\{(w_k)^1\}_{k=1}^r = \frac{1}{r}$ ,  $\lambda_1 = 1$ , and  $\lambda_4 = 20$ . Other involved variables are initialized as 0.

3: **while**  $t++ < mit$

4: Update  $(\mathbf{Z})^{t+1}$  as (15).

5: Update  $(\mathbf{B})^{t+1}$  as (17).

6: Update  $(\mathbf{A})^{t+1}$  as (20).

7: Update  $(\mathbf{P})^{t+1}$  as (21).

8: Update  $(\mathbf{K})^{t+1}$  as (23), compute  $\mathbf{w}$  via Proposition 2.

9: Update  $(\mathbf{Y})^{t+1}$  using  $(\mathbf{Y})^t + \mu((\mathbf{B})^{t+1} - (\mathbf{Z})^{t+1})$ .

10: Update  $(\mu)^{t+1}$  using  $\min(\mu_{max}, \rho(\mu)^t)$ .

11: Break if  $\|(\mathbf{B})^{t+1} - (\mathbf{Z})^{t+1}\|_\infty \leq \epsilon$ .

12: **end while**

13: Obtain the optimal  $\mathbf{Z}$  and  $\mathbf{A}$ , and then fuse them via (24).

14: Perform spectral clustering.

Output: The clustering results: ACC, NMI, and purity.

---

### 3.3. Kernel weight strategy

After obtaining the consensus kernel  $\mathbf{K}$  using (23), the following proposition can determine weight  $w_i$  automatically. Proposition 2 The weight  $w_i$  is determined as  $\frac{1}{2\sqrt{\|\mathbf{K} - \mathbf{K}^i\|_F^2 + \xi}}$ , where  $\xi$  is infinitely close to zero.

Proof Define the auxiliary problem with respect to  $\mathbf{K}$  as follows:

$$\min_{\mathbf{K}} \sum_{i=1}^r \sqrt{\|\mathbf{K} - \mathbf{K}^i\|_F^2} \quad (25)$$

By taking the derivative of (25) with respect to  $\mathbf{K}$  and setting to be zero, then we have



$$\sum_{i=1}^r \widehat{w}_i \frac{\partial \|\mathbf{K} - \mathbf{K}^i\|_F^2}{\partial \mathbf{K}} = 0 \quad (26)$$

where  $\widehat{w}_i = \frac{1}{2\sqrt{\|\mathbf{K} - \mathbf{K}^i\|_F^2}}$ . With applying the same operations on problem (10), we can obtain the same result as shown in Eq. (26). Thus,  $\widehat{w}_i$  can be considered as  $w_i$ , i.e.,  $w_i = \frac{1}{2\sqrt{\|\mathbf{K} - \mathbf{K}^i\|_F^2}}$ . To avoid dividing by zero in theory,  $w_i$  can be transformed into  $w_i = \frac{1}{2\sqrt{\|\mathbf{K} - \mathbf{K}^i\|_F^2 + \xi}}$ , where  $\xi$  is infinitely close to zero. The proof is completed. ■

### 3.4. Complexity analysis

The main computational complexity of Algorithm 1 consists of five parts, i.e., the subproblems for updating  $\mathbf{Z}, \mathbf{A}, \mathbf{B}, \mathbf{P}$  and  $\mathbf{K}$ . Detailedly, the complexities of updating  $\mathbf{Z}$  and  $\mathbf{B}$  are  $\mathcal{O}(n^3)$  due to the matrix inverse operation. In fact, since  $(\mathbf{K} + 2\lambda_4 \mathbf{I})$  and  $(2\lambda_2 \mathbf{L}_A + \mu \mathbf{I})$  are positive semi-definite matrices, they also can be solved within  $\mathcal{O}(n^2)$  [44]. For updating  $\mathbf{A}$ , it can be solved within  $\mathcal{O}(n^2)$ . For updating  $\mathbf{P}$ , it requires a SVD operator of  $\mathbf{L}_Z$  within complexity  $\mathcal{O}(n^3)$ , but a rank constrained SVD can be effectively solved by some packages like PROPACK<sup>1</sup> within complexity  $\mathcal{O}(bn^2)$ , where  $b = n - c$ . Since the optimization of  $\mathbf{K}$  takes closed form solution, its computational complexity is  $\mathcal{O}(n^2)$ . Moreover, the updating of  $\mathbf{w}$  costs  $\mathcal{O}(mn^2)$ . In sum, the overall complexity of the proposed algorithm is approximate to  $\mathcal{O}(tn^3)$ , where  $t$  is the iteration number of our algorithm.

## 4. Experiments

We demonstrate the effectiveness and efficiency of our SLMKC method by conducting several experiments on nine public benchmark datasets. All codes are implemented in Matlab and run on a Mac mini PC with an Intel Core i7 processor at 3.2 GHz, RAM 16-GB, and macOS Mojave operating system.

### 4.1. Datasets

We employ nine widely used benchmark datasets, including six image datasets (i.e., Yale,<sup>2</sup> Jaffe,<sup>3</sup> ORL,<sup>4</sup> AR,<sup>5</sup> COIL20,<sup>6</sup> and BA<sup>7</sup>) and three text corporas<sup>8</sup> (i.e., TR11, TR41, and TR45). Some samples and the statistics descriptions of these used image datasets are shown in Fig. 2 and Table 2, respectively.

### 4.2. Compared methods and evaluation criteria

We compare our SLMKC with the following state-of-the-art MKC competitors, including the following: MKKM [14], RMKKM [9], AASC [13], SCMK [18], LKGr [19], JMKSC [46], SMKL [16], and JMKSC [46]. Amongst them, MKKM and RMKKM are  $k$ -means [49] based methods, while others are spectral clustering based methods. For fair comparison, the involved parameters of these competitors have been carefully tuned as recommended by their respective authors.

As in [46], to quantitatively evaluate the clustering results of the compared methods, three widely used clustering performance metrics are applied here, i.e., the clustering accuracy (ACC), normalized mutual information (NMI), and purity. For these three metrics used, the higher values indicate the better performance.

### 4.3. Multiple kernel construction

In the same way as [9,46], we also build 12 candidate base kernels, including seven radial basis function (RBF) kernels, whose  $(i,j)$ th entry is defined as  $\mathbf{K}_{ij} = \exp\left(-\|\mathbf{X}_i - \mathbf{X}_j\|_2^2 / (2\tau\sigma^2)\right)$ , where  $\tau$  varies in the set of  $\{0.01, 0.05, 0.1, 1, 10, 50, 100\}$  and  $\sigma$  is the maximum distance between any two data points, four polynomial kernels, whose  $(i,j)$ th entry is defined as  $\mathbf{K}_{ij} = (a + \mathbf{X}_i^T \mathbf{X}_j)^b$  with  $a = \{0, 1\}$  and  $b = \{2, 4\}$ , and one cosine kernel, whose  $(i,j)$ th entry is defined as

<sup>1</sup> <http://sun.stanford.edu/~rmunk/PROPACK/>.

<sup>2</sup> <http://www.cvc.yale.edu/projects/Yalefaces/Yalefaces.html>.

<sup>3</sup> <http://www.kasrl.org/Jaffe.html>.

<sup>4</sup> <https://www.cl.cam.ac.uk/research/dtg/attarchive/facedataset.html>.

<sup>5</sup> <http://www2.ece.ohio-state.edu/~aleix/ARdataset.html>.

<sup>6</sup> <http://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php>.

<sup>7</sup> <https://cs.nyu.edu/~roweis/data/binaryalphadigs.mat>.

<sup>8</sup> <http://trec.nist.gov>.

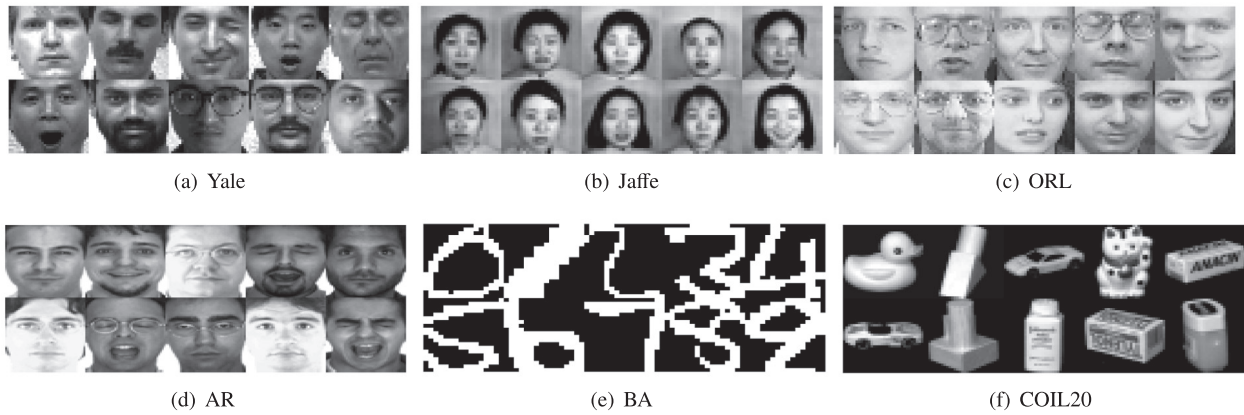


Fig. 2. Sample images of the six image datasets.

Table 2

Statistics of the nine benchmark datasets.

Dataset	# Classes	# Samples	# features	# Note
Yale	15	165	1024	Face image of 15 subjects
Jaffe	10	213	676	Face image of 10 Japanese females
AR	120	840	768	Face image of 70 men and 56 women
ORL	40	400	1024	Face images of 40 subjects
COIL20	20	1440	1024	Object images of 20 different objects
BA	36	1404	320	Digit images of '0' through '9' and 26 letters
TR11	9	414	6429	Text corpora
TR41	10	878	7454	Text corpora
TR45	10	690	8261	Text corpora

$\mathbf{K}_{ij} = (\mathbf{X}_i^T \mathbf{X}_j) / (\|\mathbf{X}_i\| \cdot \|\mathbf{X}_j\|)$ . Finally, all the kernels are normalized to  $[0, 1]$  using  $\mathbf{K}_{ij} = \mathbf{K}_{ij} / \sqrt{\mathbf{K}_{ii} \mathbf{K}_{jj}}$ . Note that the number of candidate base kernels,  $r$ , is equal to 12.

#### 4.4. Performance comparison

The results of all the comparison methods on the nine benchmark datasets are reported in Tables 3–5, and the highest scores are highlighted in boldface, where we repeat each experiment for 20 trials and report the average and standard deviation as final clustering result.

Obviously, the proposed SLMKC method mostly obtains the best performance, followed by JMKSC, SCMK, LKGr, RMKKM, SMKL, MKKM, and AASC. Compared with the second best method, JMKSC, over 4.3%, 3.1%, and 3.2% improvements are achieved in terms of average ACC, NMI and purity, respectively.

To further show the effectiveness of the proposed method, the affinity matrices produced by the comparison methods are also evaluated [6]. Taking the Jaffe dataset, for instance, it consists of ten clusters, and therefore ideally has ten diagonal blocks in affinity matrix. As shown in Fig. 3, the fused affinity matrix  $\mathbf{A} \odot \mathbf{Z}$  produced by our SLMKC has better block diagonal property and inter-cluster separability than that of all other MKC methods.

For better flow of the paper, the detailed discussions are given in Section 4.7.

#### 4.5. Parameter sensitivity

Although there are four trade-off parameters in the final objective function (11), it is very easy to tame these parameters since most of them can be fixed as constants. More specifically, the parameter  $\lambda_1$  is used to balance the term  $\|\mathbf{A}\|_F^2$  to resist trivial solution and avoid scale change; the parameter  $\lambda_4$  is used to control the contribution of the multiple kernel learning term, which can be fixed to  $\lambda_1 = 1$  and  $\lambda_4 = 20$ , respectively, by experience; and the parameters  $\lambda_2$  and  $\lambda_3$  are used to control the term  $\text{Tr}(\mathbf{Z} \mathbf{L}_A \mathbf{Z}^T)$  and the term  $\text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P})$ , respectively. Moreover, due to the top- $k$  probabilistic neighbors selection strategy, the parameter  $k$  is used to limit the number of neighbors/edges when constructing affinity graph, which is also worthy of evaluating. Therefore, take the Yale and ORL datasets for example, we conduct two experiments successively, the former is used to evaluate the sensitivity of the parameters  $\lambda_2$  and  $\lambda_3$ , and the latter is used to evaluate the sensitivity of the parameter  $k$ .

**Table 3**

Clustering performance comparison (average  $\pm$  standard deviation) in term of ACC. Symbol 'Avg' indicates the average performance of one evaluated MKC method on all datasets. The highest scores are in bold.

Data	MKKM	RMKKM	AASC	SCMK	LKGr	SMKL	JMKSC	SLMKC
Yale	0.457 $\pm$ 0.041	0.521 $\pm$ 0.034	0.406 $\pm$ 0.027	0.582 $\pm$ 0.025	0.540 $\pm$ 0.030	0.582 $\pm$ 0.017	0.630 $\pm$ 0.006	<b>0.667 <math>\pm</math> 0.000</b>
ORL	0.475 $\pm$ 0.023	0.556 $\pm$ 0.024	0.272 $\pm$ 0.009	0.656 $\pm$ 0.015	0.616 $\pm$ 0.016	0.573 $\pm$ 0.032	0.725 $\pm$ 0.014	<b>0.735 <math>\pm</math> 0.008</b>
Jaffe	0.746 $\pm$ 0.069	0.871 $\pm$ 0.053	0.304 $\pm$ 0.008	0.869 $\pm$ 0.022	0.861 $\pm$ 0.052	0.967 $\pm$ 0.000	0.967 $\pm$ 0.007	<b>1.000 <math>\pm</math> 0.000</b>
AR	0.286 $\pm$ 0.014	0.344 $\pm$ 0.012	0.332 $\pm$ 0.006	0.544 $\pm$ 0.024	0.314 $\pm$ 0.015	0.263 $\pm$ 0.009	0.609 $\pm$ 0.007	<b>0.630 <math>\pm</math> 0.016</b>
BA	0.405 $\pm$ 0.019	0.434 $\pm$ 0.018	0.271 $\pm$ 0.003	0.384 $\pm$ 0.014	0.444 $\pm$ 0.018	0.246 $\pm$ 0.012	0.484 $\pm$ 0.015	<b>0.511 <math>\pm</math> 0.012</b>
COIL-20	0.548 $\pm$ 0.058	0.667 $\pm$ 0.028	0.349 $\pm$ 0.050	0.591 $\pm$ 0.028	0.618 $\pm$ 0.051	0.487 $\pm$ 0.031	0.696 $\pm$ 0.016	<b>0.884 <math>\pm</math> 0.003</b>
TR11	0.501 $\pm$ 0.048	0.577 $\pm$ 0.094	0.472 $\pm$ 0.008	0.549 $\pm$ 0.015	0.607 $\pm$ 0.043	0.708 $\pm$ 0.033	0.737 $\pm$ 0.002	<b>0.741 <math>\pm</math> 0.004</b>
TR41	0.561 $\pm$ 0.068	0.627 $\pm$ 0.073	0.459 $\pm$ 0.001	0.650 $\pm$ 0.068	0.595 $\pm$ 0.020	0.671 $\pm$ 0.002	0.689 $\pm$ 0.004	<b>0.691 <math>\pm</math> 0.001</b>
TR45	0.585 $\pm$ 0.066	0.640 $\pm$ 0.071	0.526 $\pm$ 0.008	0.634 $\pm$ 0.058	0.663 $\pm$ 0.042	0.671 $\pm$ 0.004	0.687 $\pm$ 0.036	<b>0.751 <math>\pm</math> 0.000</b>
Avg	0.507 $\pm$ 0.045	0.582 $\pm$ 0.045	0.377 $\pm$ 0.013	0.607 $\pm$ 0.030	0.584 $\pm$ 0.032	0.574 $\pm$ 0.016	0.692 $\pm$ 0.012	<b>0.734 <math>\pm</math> 0.005</b>

**Table 4**

Clustering performance comparison (average  $\pm$  standard deviation) in term of NMI. Symbol 'Avg' indicates the average performance of one evaluated MKC method on all datasets. The highest scores are in bold.

Data	MKKM	RMKKM	AASC	SCMK	LKGr	SMKL	JMKSC	SLMKC
Yale	0.501 $\pm$ 0.036	0.556 $\pm$ 0.025	0.468 $\pm$ 0.028	0.576 $\pm$ 0.012	0.566 $\pm$ 0.025	0.614 $\pm$ 0.015	0.631 $\pm$ 0.002	<b>0.651 <math>\pm</math> 0.004</b>
ORL	0.689 $\pm$ 0.016	0.748 $\pm$ 0.018	0.438 $\pm$ 0.007	0.808 $\pm$ 0.008	0.794 $\pm$ 0.008	0.733 $\pm$ 0.027	0.852 $\pm$ 0.012	<b>0.870 <math>\pm</math> 0.004</b>
Jaffe	0.798 $\pm$ 0.058	0.893 $\pm$ 0.041	0.272 $\pm$ 0.006	0.868 $\pm$ 0.021	0.869 $\pm$ 0.031	0.951 $\pm$ 0.000	0.952 $\pm$ 0.010	<b>1.000 <math>\pm</math> 0.000</b>
AR	0.592 $\pm$ 0.014	0.655 $\pm$ 0.015	0.651 $\pm$ 0.005	0.775 $\pm$ 0.009	0.648 $\pm$ 0.007	0.568 $\pm$ 0.014	0.820 $\pm$ 0.002	<b>0.824 <math>\pm</math> 0.002</b>
BA	0.569 $\pm$ 0.008	0.585 $\pm$ 0.011	0.423 $\pm$ 0.004	0.544 $\pm$ 0.012	0.604 $\pm$ 0.009	0.486 $\pm$ 0.011	0.621 $\pm$ 0.007	<b>0.642 <math>\pm</math> 0.009</b>
COIL-20	0.707 $\pm$ 0.033	0.773 $\pm$ 0.017	0.419 $\pm$ 0.027	0.726 $\pm$ 0.011	0.766 $\pm$ 0.023	0.628 $\pm$ 0.018	0.818 $\pm$ 0.007	<b>0.939 <math>\pm</math> 0.002</b>
TR11	0.446 $\pm$ 0.046	0.561 $\pm$ 0.118	0.394 $\pm$ 0.003	0.371 $\pm$ 0.018	0.597 $\pm$ 0.031	0.557 $\pm$ 0.068	<b>0.673 <math>\pm</math> 0.002</b>	0.664 $\pm$ 0.004
TR41	0.578 $\pm$ 0.042	0.635 $\pm$ 0.092	0.431 $\pm$ 0.000	0.492 $\pm$ 0.017	0.604 $\pm$ 0.023	0.625 $\pm$ 0.004	0.660 $\pm$ 0.003	<b>0.679 <math>\pm</math> 0.002</b>
TR45	0.562 $\pm$ 0.056	0.627 $\pm$ 0.092	0.420 $\pm$ 0.014	0.584 $\pm$ 0.051	0.671 $\pm$ 0.020	0.622 $\pm$ 0.007	0.690 $\pm$ 0.022	<b>0.730 <math>\pm</math> 0.000</b>
Avg	0.605 $\pm$ 0.034	0.670 $\pm$ 0.048	0.435 $\pm$ 0.010	0.638 $\pm$ 0.018	0.680 $\pm$ 0.020	0.643 $\pm$ 0.018	0.746 $\pm$ 0.007	<b>0.778 <math>\pm</math> 0.003</b>

In the first experiment, we fix  $k = 10$  and  $k = 5$  for the Yale and ORL datasets, respectively, and then use a grid search strategy to tune  $\lambda_2$  and  $\lambda_3$ . Specifically, the searching regions of  $\lambda_2$  and  $\lambda_3$  are both varied from the candidate set  $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ . The results in terms of ACC, NMI and purity are illustrated in Fig. 4. Obviously, the proposed SLMKC method works well for a wide range of the parameters  $\lambda_2$  and  $\lambda_3$ , and can achieve better performance on both the Yale and ORL datasets when setting  $\lambda_2 \in [10^{-2}, 10^{-1}]$  and  $\lambda_3 \in [10^{-2}, 10^{-1}]$ . Therefore, we fix  $\lambda_2 = 10^{-2}$  and  $\lambda_3 = 10^{-2}$  in all experiments for simplicity.

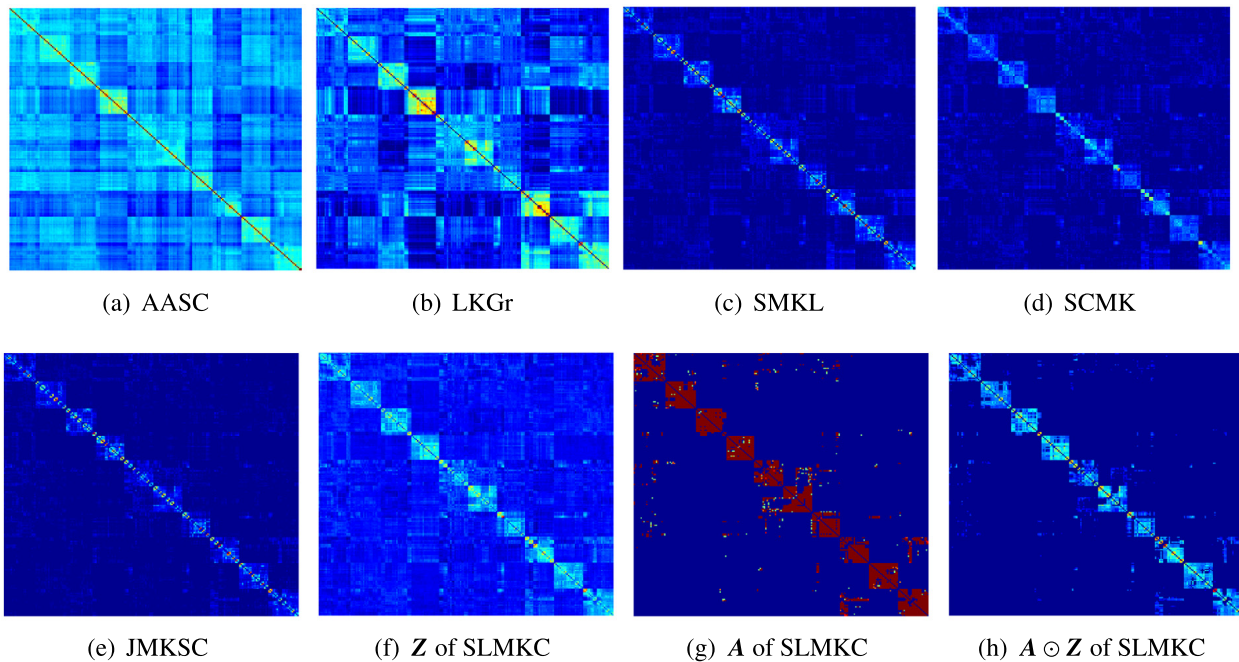
In the second experiment, we fix  $\lambda_2 = 10^{-2}$  and  $\lambda_3 = 10^{-2}$ , and then tune  $k$  from the range of  $[1, 2, \dots, 40]$  to report the influence of different parameter  $k$  to the clustering performance in terms of ACC, NMI, and purity. The results are reported in Fig. 5. We can observe that SLMKC is easily tamed and there exists a large parameter space of  $(\lambda_2, \lambda_3, k)$ . More specifically, for the Yale dataset, the more promising results are achieved when  $k$  ranges from the range of  $[8, 9, \dots, 15]$ ; for the ORL dataset, it is obvious that the clustering performance is not very sensitive to parameter  $k$  when setting  $k \geq 5$ , especially the promising performance is achieved when  $k$  ranges from the range of  $[5, 6, \dots, 10]$ .

In summary, the proposed SLMKC method is very easy to tame, i.e., it works well for a wide range of the parameters

**Table 5**

Clustering performance comparison (average  $\pm$  standard deviation) in term of Purity. Symbol 'Avg' indicates the average performance of one evaluated MKC method on all datasets. The highest scores are in bold.

Data	MKKM	RMKKM	AASC	SCMK	LKGr	SMKL	JMKSC	SLMKC
Yale	0.475 $\pm$ 0.037	0.536 $\pm$ 0.031	0.423 $\pm$ 0.026	0.610 $\pm$ 0.014	0.554 $\pm$ 0.029	0.667 $\pm$ 0.014	0.673 $\pm$ 0.007	<b>0.706</b> $\pm$ <b>0.005</b>
ORL	0.514 $\pm$ 0.021	0.602 $\pm$ 0.024	0.316 $\pm$ 0.007	0.699 $\pm$ 0.015	0.658 $\pm$ 0.017	0.648 $\pm$ 0.017	0.753 $\pm$ 0.012	<b>0.800</b> $\pm$ <b>0.004</b>
Jaffe	0.768 $\pm$ 0.062	0.889 $\pm$ 0.045	0.331 $\pm$ 0.008	0.882 $\pm$ 0.023	0.859 $\pm$ 0.038	0.967 $\pm$ 0.000	0.967 $\pm$ 0.007	<b>1.000</b> $\pm$ <b>0.000</b>
AR	0.305 $\pm$ 0.012	0.368 $\pm$ 0.010	0.350 $\pm$ 0.006	0.642 $\pm$ 0.014	0.330 $\pm$ 0.014	0.530 $\pm$ 0.014	0.656 $\pm$ 0.010	<b>0.680</b> $\pm$ <b>0.004</b>
BA	0.435 $\pm$ 0.014	0.463 $\pm$ 0.015	0.303 $\pm$ 0.004	0.606 $\pm$ 0.009	0.479 $\pm$ 0.017	0.623 $\pm$ 0.011	0.563 $\pm$ 0.018	<b>0.642</b> $\pm$ <b>0.009</b>
COIL-20	0.590 $\pm$ 0.053	0.699 $\pm$ 0.022	0.391 $\pm$ 0.044	0.635 $\pm$ 0.013	0.650 $\pm$ 0.039	0.683 $\pm$ 0.004	0.806 $\pm$ 0.010	<b>0.944</b> $\pm$ <b>0.003</b>
TR11	0.655 $\pm$ 0.044	0.729 $\pm$ 0.096	0.547 $\pm$ 0.000	0.783 $\pm$ 0.011	0.776 $\pm$ 0.030	<b>0.835</b> $\pm$ <b>0.048</b>	0.819 $\pm$ 0.001	0.752 $\pm$ 0.004
TR41	0.728 $\pm$ 0.042	0.776 $\pm$ 0.065	0.621 $\pm$ 0.001	0.758 $\pm$ 0.034	0.759 $\pm$ 0.031	0.761 $\pm$ 0.003	<b>0.799</b> $\pm$ <b>0.003</b>	0.780 $\pm$ 0.001
TR45	0.691 $\pm$ 0.058	0.752 $\pm$ 0.074	0.575 $\pm$ 0.011	0.728 $\pm$ 0.048	0.800 $\pm$ 0.026	0.816 $\pm$ 0.004	0.822 $\pm$ 0.031	<b>0.844</b> $\pm$ <b>0.000</b>
Avg	0.573 $\pm$ 0.038	0.646 $\pm$ 0.042	0.429 $\pm$ 0.012	0.705 $\pm$ 0.020	0.650 $\pm$ 0.027	0.726 $\pm$ 0.013	0.762 $\pm$ 0.011	<b>0.794</b> $\pm$ <b>0.003</b>



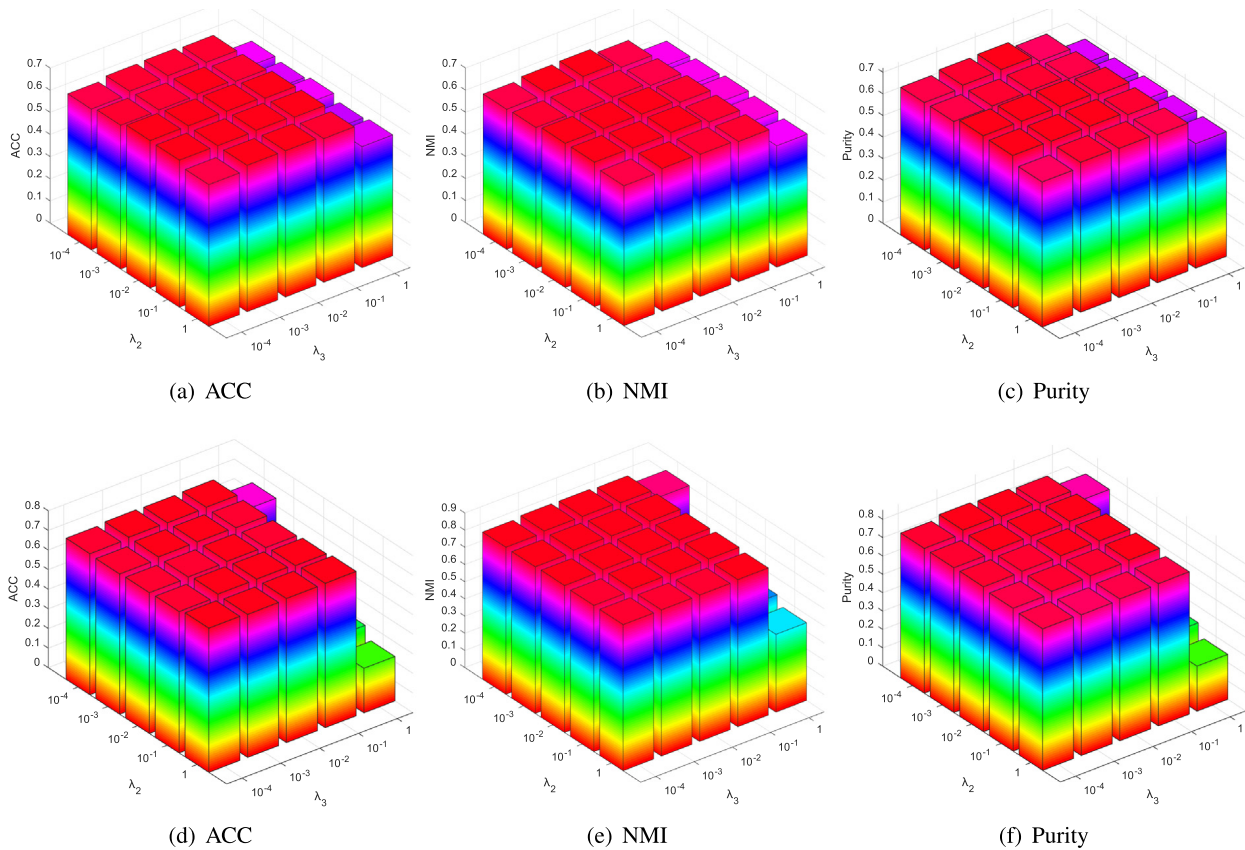
**Fig. 3.** Illustrations of the resulting affinity matrices produced by the competitors and our SLMKC on the Jaffe dataset. The Jaffe dataset consists of ten clusters. Our SLMKC simultaneously learns the self-expressiveness coefficient matrix  $\mathbf{Z}$  and the affinity matrix  $\mathbf{A}$ , and then fuses the two matrices to obtain an optimal affinity matrix  $\mathbf{A} \odot \mathbf{Z}$  for clustering purpose. The last sub-figure suggests that SLMKC can yield an intra-class dense and inter-class sparse affinity matrix having better block diagonal property.

$\lambda_1, \lambda_2, \lambda_3$ , and  $\lambda_4$ . For simplicity, we fix  $\lambda_1 = 1, \lambda_4 = 20$ , and  $\lambda_2 = \lambda_3 = 10^{-2}$ .

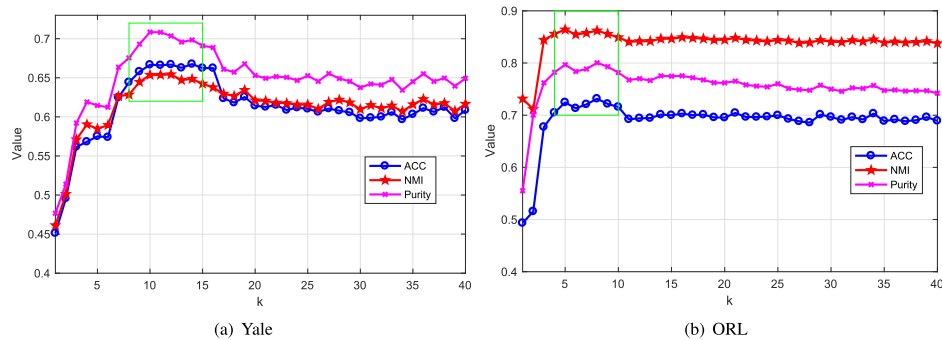
#### 4.6. Convergence study and computational time

We further investigate the convergence behavior of the proposed SLMKC method on four representative datasets (i.e., Yale, ORL, TR11, and TR45). Fig. 6 shows the trend of the convergence curve and difference curve with increasing iteration, where difference curve is defined as the different values of  $\mathbf{Z}$  between two consecutive iterations. As observed, the proposed method converges quickly within about ten iterations and becomes steadily with more iterations; meanwhile, the difference





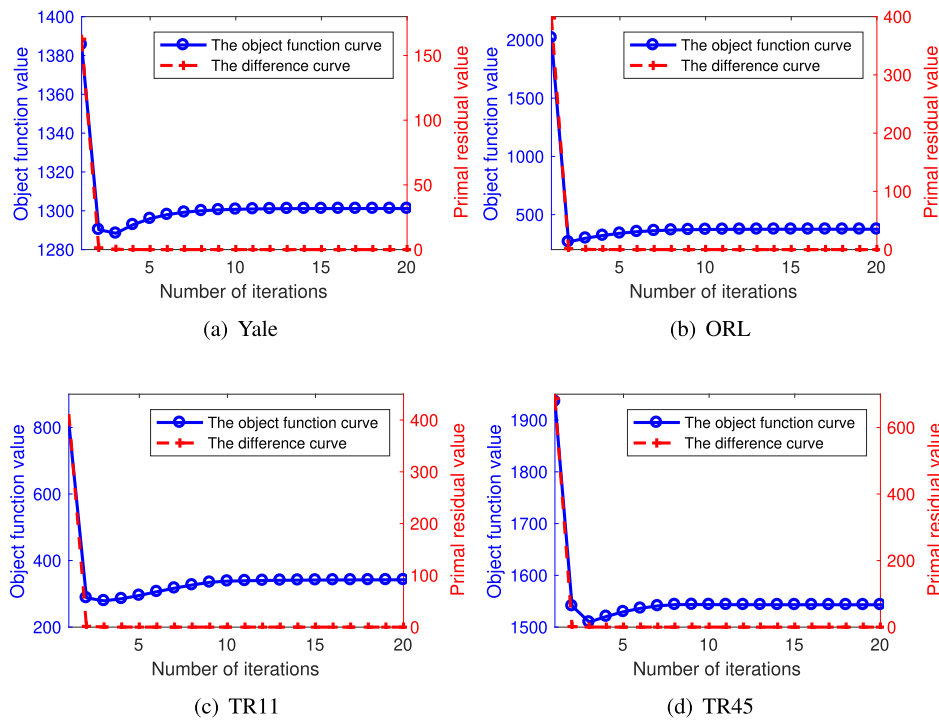
**Fig. 4.** Clustering performance of the proposed SLMKC method with respect to the two parameters  $\lambda_2$  and  $\lambda_3$  on the Yale (the first row) and ORL (the second row) datasets. Specially, the two parameters  $\lambda_2$  and  $\lambda_3$  are varied from the candidate set  $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ ; the two parameters  $\lambda_1$  and  $\lambda_4$  are fixed to  $\lambda_1 = 1$  and  $\lambda_4 = 20$  for all experiments; the parameter  $k$  is fixed to  $k = 10$  and  $k = 5$  for the Yale and ORL datasets, respectively.



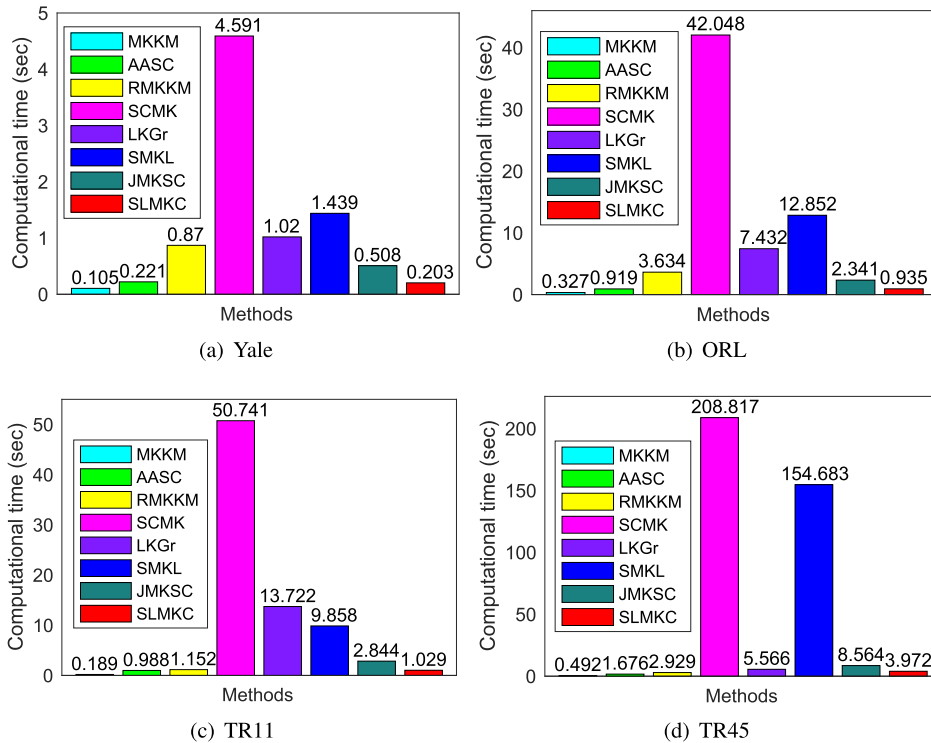
**Fig. 5.** Clustering performance of the proposed SLMKC method with respect to the number of neighbors  $k$  on the Yale and ORL datasets. We fix  $\lambda_1 = 1$ ,  $\lambda_2 = \lambda_3 = 0.01$ , and  $\lambda_4 = 20$  for the two datasets.

curves quickly reduce close to zero. This indicates that our SLMKC can converge efficiently. Note that although it is difficult to present a theoretical proof for the algorithm convergence, these curves demonstrate the convergence and effectiveness of our SLMKC empirically.

Furthermore, the computation time is also the key criterion to evaluate a clustering method. Therefore, we compute the computational time (in seconds) of all the MKC competitors on the Yale, ORL, TR11, and TR45 datasets. In order to enhance the comparability among the competitors, we set the same convergence conditions for each algorithm during the experiments. Seen from Fig. 7, the computational cost of our SLMKC is lower than RMKKM, SCMK, SMKL, LKGr, and JMKSC, whose clustering performance is worse than ours. In addition, such as MKKM and AASC, although their computational costs are less than that of the proposed method, their clustering performance is poor.



**Fig. 6.** Convergence curves and difference curves of our SLMKC on the Yale, ORL, TR11, and TR41 datasets. In each subfigure, the x-axis, the y-left-axis, and the y-right-axis denote the iteration number, the corresponding objective function value, and the corresponding primal residual value, respectively.



**Fig. 7.** Computational time (in seconds) of the proposed SLMKC method and the competitors on the Yale, ORL, TR11, and TR41 datasets.



Therefore, considering both the convergence and computational cost, the proposed SLMKC method is a fast MKC method.

#### 4.7. Discussion

From these experimental results, one could observe that: (1) According to the experimental results in Tables 3–5, we have

- Overall, in most cases, our SLMKC method consistently achieves the highest clustering performance in terms of ACC, NMI, and purity when compared with some preexisting state-of-the-art MKC methods. This demonstrates that SLMKC can be qualified for clustering tasks.
- In most cases, the  $k$ -means based MKC methods (such as MKKM and RMKKM) usually achieve worse clustering performance than these spectral clustering based MKC methods (such as JMKSC, SCMK, SMKL, LKGr, and SLMKC). This suggests that the spectral clustering based methods are generally superior to the standard  $k$ -means based ones [16].
- Our SLMKC obtains the smallest standard deviation in most cases, this phenomenon indicates that SLMKC has good statistical coherence. One main concern is that SLMKC has zero standard deviations on the Jaffe and TR45 datasets.
- Compared with these competitors, our SLMKC simultaneously learns the optimal self-expressiveness coefficient matrix (i.e.,  $\mathbf{Z}$ ) and the affinity matrix (i.e.,  $\mathbf{A}$ ), but the competitors solely learn the self-expressiveness coefficient matrix rather than both. Especially SMKL, the major difference between SLMKC and SMKL is whether to learn both a coefficient matrix and affinity matrix simultaneously. Obviously, the fused affinity graph  $\mathbf{A} \odot \mathbf{Z}$  proposed in this study is more effective for spectral clustering purpose than these competitors.
- Although we do not evaluate the clustering performance of single kernel clustering (SKC) methods, which are designed based on SKL, [9,46] have confirmed that the performance of SKC methods is largely depended on the choice of kernel function, so MKC methods are generally believed to perform better than SKC methods.

(2) According to the results in Fig. 3, we have

- It is obvious that our SLMKC method can yield an intra-class dense and inter-class sparse fused affinity matrix (i.e.,  $\mathbf{A} \odot \mathbf{Z}$ ) with proper block diagonal structure. On the one hand, the rank constraint of Laplacian matrix  $\mathbf{L}_Z$ ,  $\text{rank}(\mathbf{L}_Z) = n - c$ , urges to learn an affinity graph with exactly  $c$  connected components. On the other hand, the top- $k$  affinity graph construction strategy can effectively remove the redundant edges. It is noticeable, although the learned self-expressiveness coefficient matrix  $\mathbf{Z}$  has proper block diagonal structure, it contains too many redundant entries (also known as edges), i.e., the entries that are not on the diagonal blocks are not zeros (dark blue). Moreover, the learned affinity matrix  $\mathbf{A}$  also has proper block diagonal structure, but it cannot reflect the local similarity information of between samples, i.e., the diagonal blocks of  $\mathbf{A}$  only contain flat brown color. Apparently, the fused affinity matrix  $\mathbf{A} \odot \mathbf{Z}$  can preserve the advantages of  $\mathbf{Z}$  and  $\mathbf{A}$ , as well as eliminating the disadvantages of both.

(3) According to the results in Figs. 4 and 5, we have

- From Fig. 4, we can see that our SLMKC method is insensitive to the involved parameters, thus it is easy to be tamed. For simplicity, we can fix  $\lambda_1 = 1$ ,  $\lambda_2 = \lambda_3 = 10^{-2}$  and  $\lambda_4 = 20$ , and then tune  $k$  in all experiments.
- Additionally, from Fig. 5, it is well understood that the best performance is obtained when using a tunable  $k$  rather than a fixed  $k$  or a “fully connected” complete graph. In other words, the top- $k$  affinity graph construction strategy can effectively improve clustering performance; meanwhile, it can reduce the computational cost since we only need to update the first  $k$  important edges in each iteration.

(4) According to the results in Figs. 6 and 7, we have

- Obviously, the proposed SLMKC method converges very fast almost within ten iterations. The fast convergence property ensures the speed of the whole algorithm.

## 5. Conclusion

In this paper, a novel SLMKC method is proposed for multiple kernel clustering. Basically, SLMKC can not only effectively handle non-linear data, but also can avoid selecting a pre-defined kernel and tuning its kernel parameters. Significantly, SLMKC completely differs from existing methods, which first learn a self-expressiveness coefficient matrix and then require an additional postprocessing step to construct the final affinity matrix for clustering, SLMKC simultaneously learns the coefficient matrix and affinity matrix so that the learning of coefficient matrix and the learning of affinity matrix can help each other in each iteration. Comprehensive experimental results on nine widely used datasets testify that SLMKC consistently outperforms the state-of-the-art competitors in terms of clustering performance and computational cost.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## CRediT authorship contribution statement

**Zhenwen Ren:** Conceptualization, Methodology, Software, Visualization, Writing - original draft. **Haoyun Lei:** Data curation, Writing - original draft. **Quansen Sun:** Writing - review & editing, Supervision. **Chao Yang:** Supervision, Validation, Investigation.

## Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant Nos. 61673220), the Intelligent Manufacturing and Robot Special Project of Major Science and Technology Project in Sichuan (Grant No. 2020ZDZX0014), the Science and Technology Special Project of Major Scientific Instruments and Equipment Project in Sichuan (Grant No. 19ZDZX0119), and the Undergraduate Innovation and Entrepreneurship Project of Sichuan Province of China (Grant No. 19xcy099).

## Appendix A. Appendix A

**Theorem 2.** Denote the Laplacian matrix and the degree matrix of the symmetric affinity matrix  $\mathbf{Z} \in \mathbb{R}^{n \times n}$  as  $\mathbf{L}_Z \in \mathbb{R}^{n \times n}$  and  $\mathbf{D} \in \mathbb{R}^{n \times n}$ , respectively, the following equality is held:

$$\text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) = \frac{1}{2} \text{Tr}(\mathbf{Z} \mathbf{Q}) \quad (27)$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a symmetric matrix, whose  $(i, j)$ th element is computed by  $q_{ij} = \|p^i - p^j\|_2^2$ , and  $p^i$  and  $p^j$  are the  $i$ th and  $j$ th rows of matrix  $\mathbf{P}$ , respectively.

**Proof.** Mathematically, we have the following derivation:

$$\begin{aligned} \Rightarrow \text{Tr}(\mathbf{P}^T \mathbf{L}_Z \mathbf{P}) &= \text{Tr}(\mathbf{P}^T \mathbf{D} \mathbf{P}) - \text{Tr}(\mathbf{P}^T \mathbf{Z} \mathbf{P}) \\ &= \frac{1}{2} \left( \sum_{i=1}^n d_{ii} p^i (p^i)^T - 2 \sum_{i,j=1}^n z_{ij} p^i (p^j)^T + \sum_{j=1}^n d_{jj} p^j (p^j)^T \right) \\ &= \frac{1}{2} \sum_{i,j=1}^n z_{ij} \|p^i - p^j\|_2^2 = \frac{1}{2} \mathbf{1}^T (\mathbf{Q} \odot \mathbf{Z}) \mathbf{1} = \frac{1}{2} \text{Tr}(\mathbf{Z} \mathbf{Q}) \end{aligned}$$

Thus, Eq. (27) is achieved. The proof is completed. ■

## Appendix B. Appendix B

**Problem (19)** can be simplified to a vector form, i.e.,

$$\min_{\mathbf{a}} \|\mathbf{a} + \mathbf{e}\|_2^2 \quad \text{s.t. } \mathbf{a}^T \mathbf{1} = 1, \mathbf{a} \geq 0 \quad (28)$$

To obtain a better affinity graph and reduce computing overhead, we only update the first  $k$  important edges and drop the remaining  $n - k$  irrelevant edges, thus a sparse affinity graph can be obtained. Theoretically, the learned  $\mathbf{a}$  has  $k$  non-zero entries and problem (28) can be efficiently solved by

$$\mathbf{a} = \left[ \frac{1 + \sum_{j=1}^k e_j^-}{k} \mathbf{1} - \mathbf{e} \right]_+ \quad (29)$$

where  $[\ast]_+$  indicates the nonnegative part of a vector, and the entries of  $\mathbf{e}^\rightarrow$  are the same as that of  $\mathbf{e}$  but sorted in ascending order. Proof Considering the constraints  $\mathbf{a}^T \mathbf{1} = 1$  and  $\mathbf{a} \geq 0$ , we then remove the two constraints, thus the Lagrangian function of problem (28) is specified as

$$\mathcal{L}(\mathbf{a}, \pi, \psi) = \frac{1}{2} \|\mathbf{a} + \mathbf{e}\|_F^2 - \pi(\mathbf{a}^T \mathbf{1} - 1) - \psi^T \mathbf{a} \quad (30)$$

where  $\pi$  and  $\psi \geq 0$  are two Lagrangian multipliers. According to KKT condition, we then have

$$\begin{cases} \mathbf{a} + \mathbf{e} - \pi \mathbf{1} - \psi = 0 \\ \mathbf{a}^T \mathbf{1} - 1 = 0 \\ \psi^T \mathbf{a} = 0 \end{cases} \quad (31)$$

Note that the equation  $\psi^T \mathbf{a} = 0$  is held when the condition that if  $\mathbf{a} \geq 0$  then  $\psi_j = 0$  is satisfied. Therefore, we obtain

$$\mathbf{a} = [\pi \mathbf{1} - \mathbf{e}]_+ \quad (32)$$

In our top- $k$  probabilistic neighbors selection strategy, the first important  $k$  edges are need to be updated, and other edges are removed. Namely, the first  $k$  entries of  $\mathbf{a}$  should be positive, i.e.,

$$\pi - \mathbf{e}_k^\rightarrow > 0 \text{ and } \pi - \mathbf{e}_{k+1}^\rightarrow = 0 \quad (33)$$

Based on Eqs. (32) and (33), and the constraint,  $\mathbf{a}^T \mathbf{1} = 1$ , there is

$$\sum_{j=1}^k (\pi - \mathbf{e}_j^\rightarrow) = 1 \quad (34)$$

Thus, we have

$$\pi = \frac{1 + \sum_{j=1}^k \mathbf{e}_j^\rightarrow}{k} \quad (35)$$

Then, we substitute  $\pi$  in Eq. (32) by  $\frac{1 + \sum_{j=1}^k \mathbf{e}_j^\rightarrow}{k}$  in Eq. (35), the solution of problem (28) (i.e., Eq. (29)) is obtained. The proof is completed. ■

## References

- [1] Ali Akgül, A novel method for a fractional derivative with non-local and non-singular kernel, *Chaos Solitons Fract.* 114 (2018) 478–482.
- [2] Ali Akgül, Reproducing kernel hilbert space method based on reproducing kernel functions for investigating boundary layer flow of a powell–eyring non-newtonian fluid, *J. Taibah Univ. Sci.* 13 (1) (2019) 858–863.
- [3] Ali Akgül, Mahmut Modanli, Crank–nicholson difference method and reproducing kernel function for third order fractional differential equations in the sense of atangana–baleanu caputo derivative, *Chaos Solitons Fract.* 127 (2019) 10–16.
- [4] Esra Karatas Akgül, Solutions of the linear and nonlinear differential equations within the generalized fractional derivatives. *Chaos Interdisc. J. Nonlinear Sci.* 29 (2) (2019) 023108.
- [5] Dumitru Baleanu, Arran Fernandez, Ali Akgül, On a fractional operator combining proportional and classical differintegrals, *Mathematics* 8 (3) (2020) 360.
- [6] James C. Bezdek, Richard J. Hathaway, Vat: a tool for visual assessment of (cluster) tendency, in: *Proceedings of the 2002 International Joint Conference on Neural Networks. IJCNN'02 (Cat. No. 02CH37290)*, vol. 3, IEEE, 2002, pp. 2225–2230.
- [7] Serhat S. Bucak, Rong Jin, Anil K. Jain, Multiple kernel learning for visual object recognition: a review, *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (7) (2013) 1354–1369.
- [8] Tingquan Deng, Dongsheng Ye, Rong Ma, Hamido Fujita, Lvnan Xiong, Low-rank local tangent space embedding for subspace clustering, *Inf. Sci.* 508 (2020) 1–21.
- [9] Liang Du, Peng Zhou, Lei Shi, Hanmo Wang, Mingyu Fan, Wenjian Wang, Yi-Dong Shen, Robust multiple kernel k-means using l21-norm, in: *IJCAI*, 2015, pp. 3476–3482.
- [10] Ehsan Elhamifar, Rene Vidal, Sparse subspace clustering: algorithm, theory, and applications, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (11) (2013) 2765–2781.
- [11] Reza Ghaemi, Md. Nasir Sulaiman, Hamidah Ibrahim, Norwati Mustapha, et al., A survey: clustering ensembles techniques. *World Acad. Sci. Eng. Technol.* 50 (2009) 636–645.
- [12] Jeffrey Ho, Ming-Hsuan Yang, Jongwoo Lim, Kuang-Chih Lee, David Kriegman, Clustering appearances of objects under varying illumination conditions, in: *Proceedings of the 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2003, pp. 11–18.
- [13] Hsin-Chien Huang, Yung-Yu Chuang, Chu-Song Chen, Affinity aggregation for spectral clustering, in: *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2012, pp. 773–780.
- [14] Hsin-Chien Huang, Yung-Yu Chuang, Chu-Song Chen, Multiple kernel fuzzy clustering, *IEEE Trans. Fuzzy Syst.* 20 (1) (2012) 120–134.
- [15] Zhao Kang, Xiao Lu, Yiwei Lu, Chong Peng, Wenyu Chen, Zenglin Xu, Structure learning with similarity preserving, *Neural Networks* page 2020, 10.1016/j.neunet.2020.05.030.
- [16] Zhao Kang, Xiao Lu, Jinfeng Yi, Zenglin Xu, Self-weighted multiple kernel learning for graph-based clustering and semi-supervised classification, *IJCAI* (2018) 2312–2318.
- [17] Zhao Kang, Haiqi Pan, Steven CH Hoi, Zenglin Xu, Robust graph learning from noisy data, *IEEE Trans. Cybern.* (2019).

- [18] Zhao Kang, Chong Peng, Qiang Cheng, Zenglin Xu, Unified spectral clustering with optimal graph, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018, pp. 3366–3373.
- [19] Zhao Kang, Liangjian Wen, Wenyu Chen, Xu. Zenglin, Low-rank kernel learning for graph-based clustering, *Knowl.-Based Syst.* 163 (2019) 510–517.
- [20] Chun-Guang Li, Rene Vidal, Structured sparse subspace clustering: a unified optimization framework, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 277–286.
- [21] Chun-Guang Li, Chong You, René Vidal, Structured sparse subspace clustering: a joint affinity learning and subspace clustering framework, *IEEE Trans. Image Process.* 26 (6) (2017) 2988–3001.
- [22] Miaomiao Li, Xinwang Liu, Lei Wang, Yong Dou, Jianping Yin, En Zhu, Multiple kernel clustering with local kernel alignment maximization, in: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, AAAI Press, 2016, pp. 1704–1710. .
- [23] Yong Li, Jing Liu, Hanqing Lu, Songde Ma, Learning robust face representation with classwise block-diagonal structure, *IEEE Trans. Inf. Forensics Secur.* 9 (12) (2014) 2051–2062.
- [24] Guangcan Liu, Zhouchen Lin, Shuicheng Yan Ju, Yong Yu Sun, Yi Ma, Robust recovery of subspace structures by low-rank representation, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (1) (2012) 171–184.
- [25] Guangcan Liu, Zhouchen Lin, Yong Yu, Robust subspace segmentation by low-rank representation, in: Proceedings of the 27th International Conference on Machine Learning (ICML-10), 2010, pp. 663–670. .
- [26] Hua Feng Liu, Li Ping Jing, Y.U. Jian, Survey of matrix factorization based recommendation methods by integrating social information, *J. Software* (2018).
- [27] Xinwang Liu, Yong Dou, Jianping Yin, Lei Wang, En Zhu, Multiple kernel k-means clustering with matrix-induced regularization, in: Thirtieth AAAI Conference on Artificial Intelligence, 2016, pp. 1888–1894.
- [28] Can-Yi Lu, Hai Min, Jie Gui, Lin Zhu, Ying-Ke Lei, Face recognition via weighted sparse representation, *J. Visual Commun. Image Represent.* 24 (2) (2013) 111–116.
- [29] Canyi Lu, Jiashi Feng, Zhouchen Lin, Tao Mei, Shuicheng Yan, Subspace clustering by block diagonal representation, *IEEE Trans. Pattern Anal. Mach. Intell.* 41 (2) (2019) 487–501.
- [30] Canyi Lu, Jinhui Tang, Min Lin, Liang Lin, Shuicheng Yan, Zhouchen Lin, Correntropy induced l2 graph for robust subspace clustering, in: Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 1801–1808.
- [31] Saeid Niazmardi, Begüm Demir, Lorenzo Bruzzone, Abdolreza Safari, Saeid Homayouni, Multiple kernel learning for remote sensing image classification, *IEEE Trans. Geosci. Remote Sens.* 56 (3) (2017) 1425–1443.
- [32] Feiping Nie, Guohao Cai, Jing Li, Xuelong Li, Auto-weighted multi-view learning for image clustering and semi-supervised classification, *IEEE Trans. Image Process.* 27 (3) (2017) 1501–1511.
- [33] Feiping Nie, Xiaoqian Wang, Heng Huang, Clustering and projected clustering with adaptive neighbors, in: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 977–986. ACM, 2014.
- [34] Feiping Nie, Xiaoqian Wang, Michael I. Jordan, Heng Huang, The constrained laplacian rank algorithm for graph-based clustering, in: Thirtieth AAAI Conference on Artificial Intelligence, 2016, pp. 1969–1976.
- [35] Zhenwen Ren, Haoran Li, Chao Yang, Quansen Sun, Multiple kernel subspace clustering with local structural graph and low-rank consensus kernel learning, *Knowl.-Based Syst.* (2019) 105040. .
- [36] Zhenwen Ren, Quansen Sun, Simultaneous global and local graph structure preserving for multiple kernel clustering, *IEEE Trans. Neural Networks Learn. Syst.* (2020). 10.1109/TNNLS.2020.2991366. .
- [37] Zhenwen Ren, Quansen Sun, Bin Wu, Xiaoqian Zhang, Wenzhu Yan, Learning latent low-rank and sparse embedding for robust image feature extraction, *IEEE Trans. Image Process.* 29 (1) (2019) 2094–2107.
- [38] Zhenwen Yang Ren, X. Simon, Quansen Sun, Tao Wang, Consensus affinity graph learning for multiple kernel clustering, *IEEE Trans. Cybern.* (2020) 10.1109/TCYB.2020.3000947. .
- [39] Jingjing Tang, Yingjie Tian, Dalian Liu, Gang Kou, Coupling privileged kernel method for multi-view learning, *Inf. Sci.* 481 (2019) 110–127.
- [40] Chuanli Wang, En Zhu, Xinwang Liu, Long Gao, Jianping Yin, Ning Hu, Multiple kernel clustering with global and local structure alignment, *IEEE Access* 6 (2018) 77911–77920.
- [41] Hao Wang, Yan Yang, Bing Liu, Hamido Fujita, A study of graph-based system for multi-view clustering, *Knowl.-Based Syst.* 163 (2019) 1009–1019.
- [42] Tinghua Wang, Jie Lu, Guangquan Zhang, Two-stage fuzzy multiple kernel learning based on hilbert–schmidt independence criterion, *IEEE Trans. Fuzzy Syst.* 26 (6) (2018) 3703–3714.
- [43] Jie Wen, Bob Zhang, Yong Xu, Jian Yang, Na. Han, Adaptive weighted nonnegative low-rank representation, *Pattern Recogn.* 81 (2018) 326–340.
- [44] Shijie Xiao, Minghui Tan, Dong Xu, Weighted block-sparse low rank representation for face clustering in videos, in: European Conference on Computer Vision, Springer, 2014, pp. 123–138.
- [45] Zhe Xue, Junping Du, Dawei Du, Siwei Lyu, Deep low-rank subspace ensemble for multi-view clustering, *Inf. Sci.* 482 (2019) 210–227.
- [46] Chao Yang, Zhenwen Ren, Quansen Sun, Mingna Wu, Maowei Yin, Yuan Sun, Joint correntropy metric weighting and block diagonal regularizer for robust multiple kernel subspace clustering, *Inf. Sci.* 500 (2019) 48–66.
- [47] Shuangyan Yi, Yingyi Liang, Zhenyu He, Yi Li, Yiu-Ming Cheung, Dual pursuit for subspace learning, *IEEE Trans. Multimedia* 21 (6) (2018) 1399–1411.
- [48] Sihang Zhou, Xinwang Liu, Miaomiao Li, En Zhu, Li Liu, Changwang Zhang, Jianping Yin, Multiple kernel clustering with neighbor-kernel subspace segmentation, *IEEE Trans. Neural Networks Learn. Syst.* (2020) 1351–1362. .
- [49] Jihua Zhu, Zutao Jiang, Georgios D. Evangelidis, Changqing Zhang, Shanmin Pang, Zhongyu Li, Efficient registration of multi-view point sets by k-means clustering, *Inf. Sci.* 488 (2019) 205–218. .