

Structured Graph Learning for Clustering and Semi-supervised Classification

Zhao Kang^a, Chong Peng^b, Qiang Cheng^c, Xinwang Liu^d, Xi Peng^e, Zenglin Xu^f, Ling Tian^a

^a*School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu, 611731, China.*

^b*College of Computer Science and Technology, Qingdao University, Qingdao, 266071, China.*

^c*Institute of Biomedical Informatics and Department of Computer Science, University of Kentucky, Lexington, KY, 40506, USA.*

^d*School of Computer Science, National University of Defense Technology, Changsha, 410073, China.*

^e*College of Computer Science, Sichuan University, Chengdu, 610064, China.*

^f*Department of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, 518055, China.*

Abstract

Graphs have become increasingly popular in modeling structures and interactions in a wide variety of problems during the last decade. Graph-based clustering and semi-supervised classification techniques have shown impressive performance. This paper proposes a graph learning framework to preserve both the local and global structure of data. Specifically, our method uses the self-expressiveness of samples to capture the global structure and adaptive neighbor approach to respect the local structure. Furthermore, most existing graph-based methods conduct clustering and semi-supervised classification on the graph learned from the original data matrix, which doesn't have explicit cluster structure, thus they might not achieve the optimal performance. By considering rank constraint, the achieved graph will have exactly c connected components if there are c clusters or classes. As a byproduct of this, graph learning and label inference are jointly and iteratively implemented in a principled way. Theoretically, we show that our model is equivalent to a combination of kernel k-means and k-means methods under certain condition. Extensive experiments on clustering and semi-supervised classification demonstrate that the proposed method

outperforms other state-of-the-art methods.

Keywords: Similarity graph, Rank constraint, Clustering, Semi-supervised classification, Local and global structure, Kernel method.

1. Introduction

As a natural way to represent structure or connections in data, graphs have broad applications including world wide web, social networks, information retrieval, bioinformatics, computer vision, natural language processing, and many others. Some special cases of graph algorithms, such as graph-based clustering [1, 2], graph embedding [3], graph-based semi-supervised classification [4], signal processing [5], have attracted increasing attention in the recent years.

Clustering refers to the task of finding subsets of similar samples and grouping them together, such that samples in the same cluster would share high similarity to each other, whereas samples in different groups are dissimilar [6, 7]. By leveraging a small set of labeled data, semi-supervised classification aims at determining the labels of a large collection of unlabeled samples based on relationships among the samples [8]. In essence, both clustering and semi-supervised classification algorithms are trying to predict labels for samples [9]. As fundamental techniques in machine learning and pattern recognition, they have been facilitating various research fields and have been extensively studied.

Among numerous clustering and semi-supervised classification methods developed in the past decades, graph based techniques often provide impressive performance. In general, these methods consist of two key steps. First, an affinity graph is constructed from all data points to represent the similarity among the samples. Second, spectral clustering [10] algorithm or label propagation [11] method is utilized to obtain the final labels. Therefore, the start step of building graph might heavily impact the subsequent step and finally lead to suboptimal performance. Since underlying structures of data are often unknown in advance, this pose a major challenge for graph construction. Consequently, the final result might be far from optimal. Unfortunately, constructing a good graph that

best captures the essential data structure is still known to be fundamentally challenging [12].

The existing strategies to define adjacency graph can be roughly divided into three categories: a) the metric based approaches, which use some functions to measure the similarity among data points [13], such as Cosine, Euclidean distance, Gaussian function; b) the local structure approaches, which induce the similarity by representing each datum as a linear combination of local neighbors [14] or learning a probability value for two points as neighbors [15]; c) the global self-expressiveness property based approaches, which encode each datum as a weighted combination of all other samples, i.e., its direct neighbors and reachable indirect neighbors [16, 17]. The traditional metric based approaches and the local neighbor based methods depend upon the selection of metric or the local neighborhood parameter, which heavily influence final accuracy. Hence, they are not reliable in practice [18].

On the other hand, adaptive neighbor [15] and self-expressiveness approaches [19, 20] automatically learn graph from data. As a matter of fact, they share a similar spirit as locality preserve projection (LPP) and locally linear embedding (LLE), respectively. Different from LPP and LLE, they don't specify the neighborhood size and predefine the similarity graph. In realistic applications, they enjoy several benefits. First, automatically determining the most informative neighbors for each data point will avoid the inconsistent drawback in widely used k -nearest-neighborhood and ϵ -nearest-neighborhood graph construction techniques, which provide unstable performance with respect to different k or ϵ values [21]. Second, they are independent of measure metric, while traditional methods are often data-dependent and sensitive to noise and outliers [22]. Third, they can tackle data with structures at different scales of size and density [23]. Therefore, they are preferred in practice. For example, [24] performs dimension reduction and graph learning based on adaptive neighbor in a unified framework.

Nevertheless, they emphasize different aspects of data structure information, i.e., local and global, respectively. As demonstrated in many problems, such as

dimension reduction [25], feature selection [26], semi-supervised classification [27], clustering [14], local and global structure information are both important to algorithm performance since they can provide complementary information to each other and thus enhance the performance. In the paper, we combine them into a unified framework for graph learning task.

Moreover, most existing graph-based methods conduct clustering and semi-supervised classification on the graph learned from the original data matrix, which doesn't have explicit cluster structure, thus they might not achieve the optimal performance. For example, the seminal work [20] assumes a low-rank structure of graph, whose solution might not be optimal due to the bias of nuclear norm [28]. Ideally, the achieved graph should have exactly c connected components if there are c clusters or classes. Most existing methods fail to take this information into account. In this paper, we consider rank constraint to meet this requirement. As an extension to our previous work [22], we establish the theoretical connection of our clustering model to kernel k-means and k-means and consider semi-supervised classification application. As an added bonus, graph learning and label inference are seamlessly integrated into a unified objective function. This is quite different from traditional ways, where graph learning and label inference are performed in two separate steps, which easily lead to suboptimal results. To overcome the limitation of single kernel method, we further extend our model to accommodate multiple kernels.

Though there are many other lines of research on graph. For instance, [29] discusses the transformation issue; [30] introduces a fitness metric to learn the adjacency matrix; [31] focuses on the graph that is sampled from a graphon. Different from them, this work aims to learn a graph that has explicit cluster structure. In particular, the number of clusters/classes is employed as a prior knowledge to enhance the quality of graph, which leads to improved performance of clustering and semi-supervised classification. Additionally, graph neural networks (GNN) has gained increasing popularity recently [32]. The main difference between GNN and our method is that GNN targets to process a graph that is already available in existing data, while our method is designed to learn

a good graph from feature data for further processing. Hence, our method and GNN focus on different types of data. In practice, feature data is more common than graph data. From this point of view, our method could be useful for GNN applications when the graph is not available or the graph has low quality. As a matter of fact, how to refine the graph used in GNN is a promising research direction.

To sum up, the main contributions of this paper are:

1. The similarity graph and labels are adaptively learned from the data by preserving both global and local structure information. By leveraging the interactions among them, they are mutually reinforced towards an overall optimal solution.
2. Theoretical analysis shows the connections of our model to kernel k-means, k-means, and spectral clustering methods. Our framework is more general than k-means and kernel k-means. At the same time, it solves the graph construction challenge of spectral clustering.
3. Based on our method with a single kernel, we further extend our model into an integrated framework which can simultaneously learn the similarity graph, labels, and the optimal combination of multiple kernels. Each subtask can be iteratively boosted by using the results of the others.
4. Extensive experiments on real-world data sets are conducted to testify the effectiveness and advantages of our framework over other state-of-the-art clustering and semi-supervised classification algorithms.

The rest of the paper is organized as follows. Section 2 introduces the proposed clustering method based on a single kernel. In Section 3, we show the theoretical analysis of our model. An extended model with multiple kernel learning ability is provided in Section 4. Clustering and semi-supervised classification experimental results and analysis are presented in Section 5 and 6, respectively. Section 7 draws conclusions.

Notations. Given a data set $X \in \mathcal{R}^{n \times m}$ with m features and n instances, its i -th sample and (i, j) -th element are denoted by $x_i \in \mathcal{R}^{m \times 1}$ and x_{ij} , respectively.

The ℓ_2 -norm of x_i is denoted as $\|x_i\| = \sqrt{x_i^T \cdot x_i}$, where T means transpose. The definition of squared Frobenius norm is $\|X\|_F^2 = \sum_{ij} x_{ij}^2$. I represents the identity matrix and $\mathbf{1}$ denotes a column vector with all the elements as one. $Tr(\cdot)$ is the trace operator. $0 \leq Z \leq 1$ indicates that elements of Z are in the range of $[0, 1]$.

2. Structured Graph Learning with Single Kernel

In this section, we first review local and global structure learning, then describe our model and its optimization.

2.1. Local Structure Learning

It is reasonable to assume that the similarity z_{ij} between the i -th sample x_i and the j -th sample x_j is big if distance $\|x_i - x_j\|^2$ is small. Intuitively, we follow the adaptive neighbor approach [15] to have the following objective function:

$$\min_{z_i} \sum_{j=1}^n (\|x_i - x_j\|^2 z_{ij} + \alpha z_{ij}^2) \quad s.t. \quad z_i^T \mathbf{1} = 1, \quad 0 \leq z_{ij} \leq 1, \quad (1)$$

where α is a tuning parameter and it can be computed in advance as we show later. By solving above problem, we obtain a graph matrix $Z \in \mathcal{R}^{n \times n}$, which characterizes the pairwise relationships among samples.

Define $d_{ij}^x = \|x_i - x_j\|^2 = x_i^T x_i + x_j^T x_j - 2x_i^T x_j$, then its corresponding matrix is $D^x = \text{Diag}(XX^T)\mathbf{1}\mathbf{1}^T + \mathbf{1}\mathbf{1}^T \text{Diag}(XX^T) - 2XX^T$, where $\text{Diag}(XX^T)$ is a diagonal matrix with the diagonal elements of XX^T . Thus (1) can be reformulated in matrix format as:

$$\min_Z Tr(Z^T D^x) + \alpha \|Z\|_F^2 \quad s.t. \quad Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1. \quad (2)$$

The achieved graph Z from (2) will capture the local structure information. Since choosing local neighbors may lead to disjoint components and incorrect neighbors, we advocate to preserve global neighborhoods.

2.2. Global Structure Learning

Self-expressive property has been applied to many applications and demonstrates its capability in capturing the global structure of data [33, 34]. In particular, subspace clustering is built on this property to learn an adjacency matrix [35, 23, 19]. It assumes that each data point can be linearly reconstructed from weighted combinations of all other data points, i.e., its direct neighbors and reachable indirect neighbors. The weight coefficient matrix Z also behaves like similarity matrix, since the weight z_{ij} should be big if x_i and x_j are similar. In mathematical language, this problem is written as:

$$\min_Z \|X^T - X^T Z\|_F^2 + \alpha f(Z) \text{ s.t. } Z^T \mathbf{1} = \mathbf{1}, 0 \leq Z \leq 1. \quad (3)$$

where $f(Z)$ is a regularizer on Z . For simplicity, squared Frobenius norm of Z is adopted in this paper. As a result, (3) will learn the graph matrix by following the distribution of the data points, which will reflect the global relationships.

It is easy to see that (3) is a linear model and assumes that data points are drawn from a union of subspaces. Hence, it may not work well when data points reside in a union of manifolds. As we know, nonlinear data display linearity if mapped to an implicit, higher-dimensional space [36, 37]. Therefore, we extend (3) to kernel representation through transformation ϕ , then $K_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$. It yields:

$$\begin{aligned} \min_Z \text{Tr}(K - 2KZ + Z^T K Z) + \alpha f(Z) \\ \text{s.t. } Z^T \mathbf{1} = \mathbf{1}, 0 \leq Z \leq 1. \end{aligned} \quad (4)$$

(4) will recover the nonlinear relationships in the raw space.

To make use of possible complementary information provided by the local structure and the global structure of the samples, we combine (2) and (4) into a single unified objective function:

$$\begin{aligned} \min_Z \text{Tr}(K - 2KZ + Z^T K Z) + \text{Tr}(Z^T D^x) + \alpha \|Z\|_F^2 \\ \text{s.t. } Z^T \mathbf{1} = \mathbf{1}, 0 \leq Z \leq 1. \end{aligned} \quad (5)$$

As a consequence, (5) will provide a graph matrix Z that respects the global and local structure hidden in the data. However, Z doesn't display an explicit cluster structure, thus it may not produce the optimal performance. Specifically, we expect that the connections among data samples from different classes are as weak as possible; whereas the connections among data points within the same class are as strong as possible. Ideally, the achieved graph should have exactly c connected components if there are c clusters or classes, i.e., Z is block diagonal (with proper permutations) in which each block is connected and corresponds to data samples from the same class.

2.3. Structured Graph Learning

To achieve the desired structure of graph matrix Z , we impose constraint on the rank of its Laplacian, which is defined as $L = D - \frac{Z+Z^T}{2}$, where $D \in \mathcal{R}^{n \times n}$ is the diagonal degree matrix with $d_{ii} = \sum_j \frac{z_{ij}+z_{ji}}{2}$. Concretely, we are based on the following important theorem [16].

Theorem 1. *The multiplicity c of the eigenvalue 0 of the Laplacian matrix L is equal to the number of connected components in the graph associated with Z .*

Theorem 1 indicates that $\text{rank}(L) = n - c$ if Z contains exactly c connected components. Thus our proposed **Structured Graph** learning framework with **Single Kernel** (SGSK) is:

$$\begin{aligned} \min_Z & \text{Tr}(K - 2KZ + Z^T K Z) + \text{Tr}(Z^T D^x) + \alpha \|Z\|_F^2 \\ \text{s.t.} & \quad Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1, \quad \text{rank}(L) = n - c. \end{aligned} \quad (6)$$

The problem (6) seems very difficult to solve since L also depends on Z . In the next subsection, we will design a novel algorithm to solve this problem.

2.4. Optimization

Let $\sigma_i(L)$ denotes the i -th smallest eigenvalue of L . Since L is positive semi-definite, we have $\sigma_i(L) \geq 0$. Then $\text{rank}(L) = n - c$ means $\sum_{i=1}^c \sigma_i(L) = 0$. The

problem (6) is equivalent to the following problem for a large enough γ :

$$\begin{aligned} \min_Z Tr(K - 2KZ + Z^T KZ) + Tr(Z^T D^x) + \alpha \|Z\|_F^2 \\ + \gamma \sum_{i=1}^c \sigma_i(L) \quad s.t. \quad Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1. \end{aligned} \quad (7)$$

According to the Ky Fan's Theorem [16], we have:

$$\sum_{i=1}^c \sigma_i(L) = \min_{P^T P = I} Tr(P^T L P), \quad (8)$$

where $P \in \mathcal{R}^{n \times c}$ is the cluster/label matrix. Therefore, the problem (7) can be reformulated as:

$$\begin{aligned} \min_{Z, P} Tr(K - 2KZ + Z^T KZ) + Tr(Z^T D^x) + \alpha \|Z\|_F^2 + \\ \gamma Tr(P^T L P) \quad s.t. \quad Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1, \quad P^T P = I. \end{aligned} \quad (9)$$

Then we can solve problem (9) using an alternating optimization strategy.

When Z is fixed, the problem (9) becomes:

$$\min_{P^T P = I} Tr(P^T L P). \quad (10)$$

The optimal solution P is formed by the c eigenvectors of L corresponding to the c smallest eigenvalues.

When P is fixed, the problem (9) becomes:

$$\begin{aligned} \min_Z Tr(K - 2KZ + Z^T KZ) + Tr(Z^T D^x) + \alpha \|Z\|_F^2 + \\ \gamma Tr(P^T L P) \quad s.t. \quad Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1. \end{aligned} \quad (11)$$

According to the property of Laplacian matrix, we have the following equation:

$$\sum_{i,j} \frac{1}{2} \|P_{i,:} - P_{j,:}\|^2 z_{ij} = Tr(P^T L P) \quad (12)$$

Based on it, the problem (11) can be rewritten in the vector form as:

$$\begin{aligned} \min_{z_i} z_i^T (\alpha I + K) z_i + [(d_i^x + \frac{\gamma}{2} d_i^p)^T - 2K_{i,:}] z_i \\ s.t. \quad z_i^T \mathbf{1} = 1, \quad 0 \leq z_{ij} \leq 1. \end{aligned} \quad (13)$$

Algorithm 1 The algorithm of SGSK

Input: Kernel matrix K , parameter $\gamma > 0$, α .**Initialize:** Random matrix Z .**REPEAT**

- 1: Calculate P as the c smallest eigenvectors of $L = D - \frac{Z+Z^T}{2}$.
- 2: For each i , update the i -th column of Z according to (13).

UNTIL stopping criterion is met.

where we denote $d_i^p \in \mathcal{R}^{n \times 1}$ as a vector with the j -th element $d_{ij}^p = \|P_{i,:} - P_{j,:}\|^2$. Note that the nearest neighbors to any data point x_i are not steady and they change in each iteration. Thus the neighbors are learned adaptively here, which is quite different from traditional approaches. Problem (13) can be solved in parallel by various quadratic programming packages.

We can observe that when graph Z is given, our algorithm solves a spectral clustering problem; when P is known, our algorithm learns graph to well respect the local and global structure of the data under the guidance of the cluster structure. For clarity, the complete procedure is outlined in Algorithm 1.

2.5. Convergence Analysis

SGSK is solved in an alternative way, the optimization procedure will monotonically decrease the objective function value of the problem in (9) in each iteration [38]. Since the objective function has a lower bound, such as zero, the above iteration converges.

2.6. Determination of Parameter α

In our proposed model, parameter α controls the balance between the trivial solution ($\alpha = 0$) and the uniform distribution ($\alpha = \infty$). To alleviate computational burden, a sparse z_i , i.e., only x_i 's k nearest neighbors are connected to x_i , is expected for local structure learning. Motivated by this, we introduce a practical way to set α value.

For subproblem (1), its corresponding Lagrangian function is

$$(d_i^x)^T z_i + \alpha_i z_i^T z_i - \beta(z_i^T \mathbf{1} - 1) - \rho_i^T z_i, \quad (14)$$

where β and ρ_i are the Lagrangian multipliers. For each i , we introduce a parameter α_i . By Karush-Kuhn-Tucker (KKT) condition, we have

$$z_{ij} = \left(\frac{\beta - d_{ij}^x}{2\alpha_i} \right)_+ \quad (15)$$

Considering the constraint $z_i^T \mathbf{1} = 1$, we have

$$\sum_{j=1}^k \left(\frac{\beta - d_{ij}^x}{2\alpha_i} \right) = 1 \Rightarrow \beta = \frac{2\alpha_i + \sum_{j=1}^k d_{ij}^x}{k} \quad (16)$$

To keep k nonzero components, we can have $z_{ik} > 0$ and $z_{i,k+1} = 0$ if we sort each row of D^x in ascending order denoted by $d_{i1}^x, d_{i2}^x, \dots, d_{in}^x$. Then the following inequalities hold

$$\begin{cases} \frac{\beta - d_{ik'}^x}{2\alpha_i} > 0 & \text{for } k' = 1, \dots, k \\ \frac{\beta - d_{i,k''}^x}{2\alpha_i} \leq 0 & \text{for } k'' = k+1, \dots, n. \end{cases} \quad (17)$$

Inserting the β value, we have the following inequality for α_i :

$$\frac{k}{2} d_{ik}^x - \frac{1}{2} \sum_{j=1}^k d_{ij}^x < \alpha_i \leq \frac{k}{2} d_{i,k+1}^x - \frac{1}{2} \sum_{j=1}^k d_{ij}^x \quad (18)$$

This range of α_i values will make sure z_i has exactly k nonzero elements. For convenience, we set $\alpha_i = \frac{k}{2} d_{i,k+1}^x - \frac{1}{2} \sum_{j=1}^k d_{ij}^x$. Then, the average number of nonzero elements in each row of Z is close to k if we set α to be the mean value of $\alpha_1, \alpha_2, \dots, \alpha_n$. That is,

$$\alpha = \frac{1}{n} \sum_{i=1}^n \left(\frac{k}{2} d_{i,k+1}^x - \frac{1}{2} \sum_{j=1}^k d_{ij}^x \right). \quad (19)$$

In this way, we can avoid tuning α blindly and instead we search the neighborhood size $k \in (0, n]$.

3. Theoretical Connection

3.1. Connection to Kernel K-means and K-means Clustering

Theorem 2. *When $\alpha \rightarrow \infty$, the proposed SGSK model is equivalent to a combination of kernel k-means and k-means problems.*

Proof. As aforementioned, the constraint $\text{rank}(L) = n - c$ in (6) will make Z block diagonal. Suppose $Z_i \in \mathcal{R}^{n_i \times n_i}$ is the similarity graph matrix of the i -th component, where n_i is the number of data samples in this component. Then problem (6) can be written for each i :

$$\begin{aligned} \min_{Z_i} & \|\phi(X_i) - \phi(X_i)Z_i\|_F^2 + \text{Tr}(Z_i^T D_i^x) + \alpha \|Z_i\|_F^2 \\ \text{s.t.} & \quad Z_i^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z_i \leq 1, \end{aligned} \quad (20)$$

where X_i consists of the points in the Z_i . When $\alpha \rightarrow \infty$, the above problem becomes:

$$\min_{Z_i} \|Z_i\|_F^2 \quad \text{s.t.} \quad Z_i^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z_i \leq 1. \quad (21)$$

The solution is all elements in Z_i are with the same value $\frac{1}{n_i}$.

Therefore, when $\alpha \rightarrow \infty$, the solution to problem (6) is:

$$z_{ij} = \begin{cases} \frac{1}{n_k}, & \text{if } x_i \text{ and } x_j \text{ are in the same } k\text{-th component} \\ 0, & \text{otherwise} \end{cases} \quad (22)$$

Denote the solution set of this form as \mathcal{C} . We can see that $\|Z\|_F^2 = c$ and $Z\mathbf{1} = \mathbf{1}^T Z = \mathbf{1}$. Thus (6) can be written as:

$$\min_{Z \in \mathcal{C}} \sum_i \|\phi(x_i) - \phi(X)z_i\|^2 + \text{Tr}(Z^T D^x) \quad (23)$$

For the first term, it is easy to deduce that $\phi(X)z_i$ is the mean of cluster c_i in the kernel space. Therefore, the first term in (23) is exactly the kernel k-means.

For the second term in (23), we first introduce the centering matrix, i.e., $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$. It is obvious that $H\mathbf{1} = \mathbf{0}$ and also $\mathbf{1}^T H = 0$. It can be shown that $HD^x H = -2HXX^T H$. Moreover, $\text{Tr}(Z^T D^x) = \text{Tr}(D^x Z) = \text{Tr}(HD^x HZ) +$

$\frac{1}{n}\mathbf{1}^T D^x \mathbf{1}$. Therefore, we have

$$\begin{aligned}
& \min_{Z \in \mathcal{C}} \text{Tr}(Z^T D^x) \iff \min_{Z \in \mathcal{C}} \text{Tr}(H D^x H Z) \\
& \iff \max_{Z \in \mathcal{C}} \text{Tr}(H X X^T H Z) \iff \max_{Z \in \mathcal{C}} \text{Tr}(X^T H Z H X) \\
& \iff \min_{Z \in \mathcal{C}} \text{Tr}(X^T H (I - Z) H X) \iff \min_{Z \in \mathcal{C}} \text{Tr}(S_w)
\end{aligned} \tag{24}$$

which is exactly the problem of k-means. Here, S_w is the so-called within-class scatter matrix.

Therefore, our proposed model is to solve a combination of kernel k-means and k-means clustering problems when $\alpha \rightarrow \infty$. When α is not very large, our model becomes a generalization of kernel k-means and k-means, so it can partition data of an arbitrary shape. \square

3.2. Connection to Spectral Clustering

With a prespecified graph Z , spectral clustering solves the following problem:

$$\min_{P^T P = I} \text{Tr}(P^T L P). \tag{25}$$

In general, Z does not have exactly c connected components and P may not be optimal. Unlike existing spectral clustering method, Z is not predefined in (10). Also, Z is achieved by incorporating cluster/class structure. Z and P are learned simultaneously in a coupled way, so that collaboratively improve each of them. This results in overall optimal solutions, which are confirmed by our experiments.

4. Structured Graph Learning with Multiple Kernel

The only input for our proposed model (9) is kernel K . It is well known that the performance of kernel method is strongly dependent on the selection of kernel. It is also time consuming and impractical to exhaustively search the optimal kernel. Multiple kernel learning [39] which lets an algorithm do the

picking or combination from a set of candidate kernels is an effective way to tackle this issue. Here we present an approach to identify a suitable kernel or construct a consensus kernel from a pool of predefined kernels.

Transforming and concatenating r kernel spaces with different weights $\sqrt{w_i}$ ($w_i \geq 0$), we have $\tilde{\phi}(x) = [\sqrt{w_1}\phi_1(x), \sqrt{w_2}\phi_2(x), \dots, \sqrt{w_r}\phi_r(x)]^T$. Then the combined kernel K_w becomes

$$K_w(x, y) = \langle \tilde{\phi}_w(x), \tilde{\phi}_w(y) \rangle = \sum_{i=1}^r w_i K^i(x, y). \quad (26)$$

Replacing single kernel with combined kernel, we obtain our proposed **Structured Graph learning framework with Multiple Kernel (SGSK)** as:

$$\begin{aligned} \min_{Z, P, w} & Tr(K_w - 2K_w Z + Z^T K_w Z) + Tr(Z^T D^x) + \alpha \|Z\|_F^2 \\ & + \gamma Tr(P^T L P), \\ \text{s.t.} & Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1, \quad P^T P = I, \\ & K_w = \sum_{i=1}^r w_i K^i, \quad \sum_{i=1}^r \sqrt{w_i} = 1, \quad w_i \geq 0. \end{aligned} \quad (27)$$

4.1. Optimization

We can iteratively solve Z, P , and w , so that each of them will be adaptively refined by the results of the other two.

When w is fixed, we can directly calculate K_w , and the optimization problem goes back to (9). We can update Z and P by following Algorithm 1 with K_w as the input kernel.

When Z and P are known, solving (27) with respect to w can be rewritten as:

$$\min_w \sum_{i=1}^r w_i h_i \quad \text{s.t.} \quad \sum_{i=1}^r \sqrt{w_i} = 1, \quad w_i \geq 0, \quad (28)$$

where

$$h_i = Tr(K^i - 2K^i Z + Z^T K^i Z). \quad (29)$$

The Lagrange function corresponding to (28) is

$$\mathcal{J}(w) = w^T h + g(1 - \sum_{i=1}^r \sqrt{w_i}). \quad (30)$$

According to the KKT condition, we require $\frac{\partial \mathcal{J}(w)}{\partial w_i} = 0$. Then, w has the following expression:

$$w_i = \left(h_i \sum_{j=1}^r \frac{1}{h_j} \right)^{-2}. \quad (31)$$

In summary, our algorithm for solving (27) is provided in Algorithm 2.

Algorithm 2 The algorithm of SGMK

Input: Kernel matrices $\{K^i\}_{i=1}^r$, parameter $\gamma > 0$, α .

Initialize: Random matrix Z , $w_i = 1/r$.

REPEAT

- 1: Compute K_w by (26).
- 2: Calculate P as the c smallest eigenvectors of $L = D - \frac{Z+Z^T}{2}$.
- 3: For each i , update the i -th column of Z according to (13).
- 4: Compute h by (29).
- 5: Calculate w by (31).

UNTIL stopping criterion is met.

4.2. Extend to Semi-supervised Classification

Model (6) also lends itself to semi-supervised classification. Graph construction and label inference are two fundamental stages in semi-supervised learning (SSL). Solving two separate problems only once is suboptimal since label information is not exploited when learning the graph. SGMK unifies these two fundamental components into a unified framework. Then the given labels and estimated labels will be utilized to build the graph and to predict the unknown labels.

Based on a similar approach, we can reformulate SGMK for semi-supervised

classification as:

$$\begin{aligned}
& \min_{Z, P, w} Tr(K_w - 2K_w Z + Z^T K_w Z) + Tr(Z^T D^x) + \alpha \|Z\|_F^2 \\
& \quad + \gamma Tr(P^T L P) \\
& \text{s.t. } Z^T \mathbf{1} = \mathbf{1}, \quad 0 \leq Z \leq 1, \quad P_l = Y_l, \\
& \quad K_w = \sum_{i=1}^r w_i K^i, \quad \sum_{i=1}^r \sqrt{w_i} = 1, \quad w_i \geq 0,
\end{aligned} \tag{32}$$

where $Y_l = [y_1, \dots, y_l]^T$ denote the label matrix. $y_i \in \mathcal{R}^{c \times 1}$ and l is the number of labeled points. y_i is one-hot and $y_{ij} = 1$ indicates that the i -th sample belongs to the j -th class. (32) can be solved in the same procedure as (27), the only difference is updating P .

For convenience, we rearrange all the points and put the unlabeled u points in the back, e.g., $P = [Y_l; P_u]$. To solve P , we take the derivative of (32) with respect to P , we have $LP = 0$, i.e.,

$$\begin{bmatrix} L_{ll} & L_{lu} \\ L_{ul} & L_{uu} \end{bmatrix} \begin{bmatrix} Y_l \\ P_u \end{bmatrix} = 0.$$

Then $P_u = -L_{uu}^{-1} L_{ul} Y_l$. Finally, the class label for unlabeled points could be assigned according to following decision rule:

$$y_i = \underset{j}{\operatorname{argmax}} P_{ij}. \tag{33}$$

5. Clustering Experiments

In this section, we demonstrate the effectiveness of our proposed method on clustering application.

5.1. Data Sets

We implement experiments on eight publicly available data sets. The statistics information of these data sets is summarized in Table 1. Specifically, the first

Table 1: Description of the data sets

	# instances	# features	# classes
YALE	165	1024	15
JAFFE	213	676	10
ORL	400	1024	40
AR	840	768	120
BA	1404	320	36
TR11	414	6429	9
TR41	878	7454	10
TR45	690	8261	10

five data sets include four face databases (ORL¹, YALE², AR³, and JAFFE⁴) and a binary alpha digits data set BA⁵. Tr11, Tr41, and Tr45 are derived from NIST TREC Document Database⁶.

Following the setting in [40], we design 12 kernels. They are: seven Gaussian kernels of the form $K(x, y) = \exp(-\|x - y\|_2^2 / (td_{max}^2))$, where d_{max} is the maximal distance between samples and t varies over the set $\{0.01, 0.0, 0.1, 1, 10, 50, 100\}$; a linear kernel $K(x, y) = x^\top y$; four polynomial kernels $K(x, y) = (a + x^\top y)^b$ with $a \in \{0, 1\}$ and $b \in \{2, 4\}$. Besides, all kernels are rescaled to $[0, 1]$ by dividing each element by the largest pairwise squared distance.

5.2. Comparison Methods

To fully investigate the performance of our method on clustering, we choose a good set of methods to compare.

¹<http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>

²<http://vision.ucsd.edu/content/yale-face-database>

³<http://www2.ece.ohio-state.edu/~aleix/ARdatabase.html>

⁴<http://www.kasrl.org/jaffe.html>

⁵<http://www.cs.nyu.edu/~roweis/data.html>

⁶<http://www-users.cs.umn.edu/~han/data/tmdata.tar.gz>

- **Spectral Clustering (SC)** [10]: SC is a widely used clustering technique. It enjoys the advantage of exploring the intrinsic data structures. However, how to construct a good similarity graph is an open issue. Here, we directly use kernel matrix as its input.
- **Robust Kernel K-means (RKKM)**[40]: As an extension to classical k-means clustering method, RKKM has the capability of dealing with non-linear structure, noise, and outliers in the data, since ℓ_{21} -norm is adopted to measure the loss of k-means. RKKM shows promising results on a number of real-world data sets.
- **Low Rank Representation (LRR)** [20]: Based on self-expressive property, a low-rank graph is obtained.
- **Simplex Sparse Representation (SSR)** [41]: Based on self-expressive property, a sparse graph is obtained. SSR achieves satisfying performance in numerous data sets.
- **Local structure learning approach (Local)** [42]: By using adaptive neighbor idea, this method considers local structure (1) and the rank constraint.
- **Global structure learning approach (Global)** [16]: Based on self-expressive property, this method incorporates global structure (3) and the rank constraint.
- Our proposed **SGSK** and **SGMK** methods: Our method combines both local and global structure information. The code for our method is publicly available ⁷.
- **Multiple Kernel K-means (MKKM)** [43]: It is an extension of k-means in a multiple-kernel setting. Besides, a different way of kernel weight learning is used.

⁷<https://github.com/sckangz/ICDE>

- **Affinity Aggregation for Spectral Clustering (AASC)** [44]: It is a version of spectral clustering where multiple affinity graphs exist.
- **Robust Multiple Kernel K-means (RMKKM)** [40]: It extends RKKM to the situation of multiple kernels.

5.3. Clustering Results

To quantitatively assess the performance of our proposed method, we adopt the commonly used metrics, accuracy (Acc), normalized mutual information (NMI), and Purity [45]. We present the experimental results of different methods in Table 2. We can see that our proposed methods obtain promising results. More precisely, we have the following observations.

- Compared to traditional spectral clustering and recently proposed robust kernel k-means techniques, our method can enhance the performance considerably. For instance, in terms of the best acc, SGSK improves over SC and RMMK by 42.95%, 36.01% on average, respectively.
- Adaptive neighbor and self-expressiveness based approaches outperform spectral clustering and k-means based methods. Specifically, LRR, SSR, Local, Global, SGSK perform much better than SC and RKKM on YALE, JAFFE, ORL, AR datasets. Among them, SGSK, which combines the complementary information carried by local and global structure, works the best in most cases. This confirms the equally importance of local and global structure information.
- For multiple kernel learning based methods, our proposed SGMK achieves much better results than MKKM, AASC, RMKKM. Furthermore, the performance of multiple kernel methods are close to or better than their corresponding single kernel methods.

5.4. Ablation Study

The Local and Global results in Table 2 have demonstrate the importance of local and global strcture learning. Here we further investigate their importance

in the multiple kernel setting. In particular, we show the results of SGMK and SGMK without local structure part in Table 3.

Once again, we can observe that our global and local structure unified model generally outperforms the model only with global structure learning. This strongly verifies the benefit of incorporating both global and local structure in graph learning. Furthermore, it can be seen that global part can obtain better performance than SGMK in several cases. This could be caused by the fact that we treat the global and local structure terms equally important in our model (27). In real-world applications, global structure might be more important than local structure in some data sets. In such cases, it would be more practical to introduce a parameter to balance the first two terms in Eq. (27).

5.5. Parameter Sensitivity

There are two parameters in our model: α and γ . As we discussed in subsection 2.6, the search for α can be better handled by searching for a proper neighborhood size k . Therefore, we perform grid search for the γ and k that produce the best performance. Taking YALE and JAFFE data sets as examples, we demonstrate the sensitivity of our model SGMK to γ and k in Figure 1 and 2. They illustrate that our method works well γ and k over wide ranges of values. For k , we can increase its value when there are more samples in the data set.

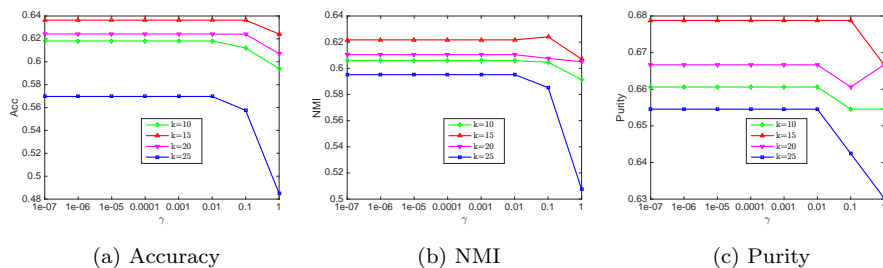


Figure 1: Parameter influence on YALE data set.

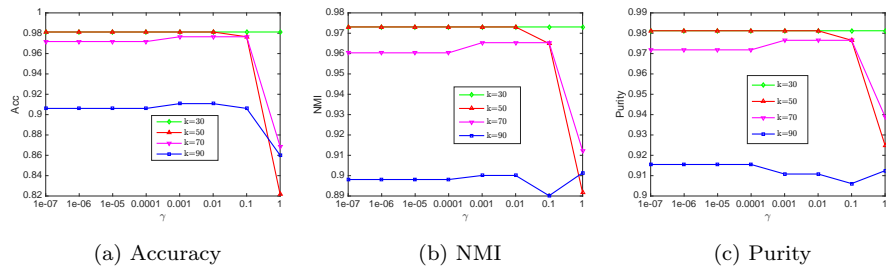


Figure 2: Parameter influence on JAFFE data set.

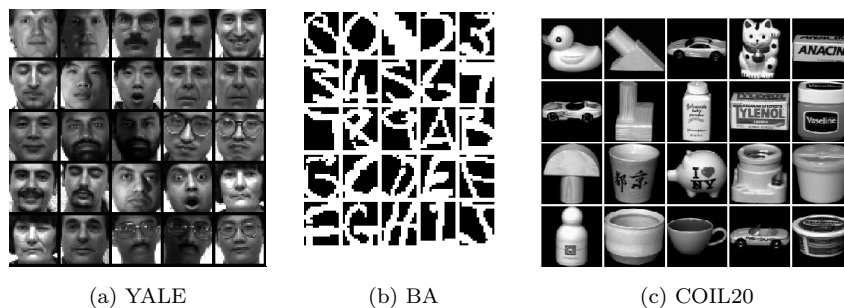


Figure 3: Sample images of YALE, BA, and COIL20.

6. Semi-supervised Classification Experiments

In this section, we assess the effectiveness of SGMK on semi-supervised learning (SSL) task.

6.1. Data Sets

1) **Evaluation on Face Recognition:** We examine the effectiveness of our graph learning for face recognition on two frequently used face databases: YALE and JEFFE. The YALE face data set contains 15 individuals, and each person has 11 near frontal images taken under different illuminations. Each image is resized to 32×32 pixels. Some sample images are shown in Figure 3a. The JAFFE face database consists of 10 individuals, and each subject has 7 different facial expressions (6 basic facial expressions +1 neutral). The images are resized to 26×26 pixels.

2) **Evaluation on Digit/Letter Recognition:** In this experiment, we address the digit/letter recognition problem on the BA database. The data set consists of digits of “0” through “9” and letters of capital “A” to “Z”. Therefore, there are 39 classes and each class has 39 samples. Figure 3b shows some sample images from BA database.

3) **Evaluation on Visual Object Recognition:** We conduct visual object recognition experiment on the COIL20 database. The database consists of 20 objects and 72 images for each object. For each object, the images were taken 5 degrees apart as the object is rotating on a turntable. The size of each image is 32×32 pixels. Some sample images are shown in Figure 3c.

Similar to clustering experiment, we construct 7 kernels for each data set. They include: four Gaussian kernels with t varies over $\{0.1, 1, 10, 100\}$; a linear kernel $K(x, y) = x^\top y$; two polynomial kernels $K(x, y) = (a + x^\top y)^2$ with $a \in \{0, 1\}$.

6.2. Comparison Methods

We compare our method with several other state-of-the-art algorithms.

- **Local and Global Consistency (LGC)** [27]: LGC is a popular label propagation method. For this method, kernel matrix is used to compute L .
- **Gaussian Field and Harmonic function (GFHF)** [13]: Different from LGC, GFHF is another mechanics to infer those unknown labels as a process of propagating labels through the pairwise similarity.
- **Semi-supervised Classification with Adaptive Neighbours (SCAN)** [15]: Based on adaptive neighbors method, SCAN adds the rank constraint to ensure that Z has exact c connected components. As a result, the similarity matrix and class indicator matrix F are learned simultaneously. It shows much better performance than many other techniques.
- **A Unified Optimization Framework for Semisupervised Learning** [8]: Li et al. propose a unified framework based on self-expressiveness approach. Similar to SCAN, the similarity matrix and class indicator matrix F are updated alternatively. By using low-rank and sparse regularizer, they have S^2LRR and S^3R method, respectively.
- Our Proposed **SGMK**: SGMK integrates both local and global structure information, with a rank constraint to improve the quality of graph.

6.3. Classification Results

We randomly choose some portions of samples as labeled data and repeat 20 times. In our experiment, 10%, 30%, 50% of samples in each class are randomly selected and labeled. Then, classification accuracy and deviation are shown in Table 4. For GFHF and LGC, the aforementioned seven kernels are tested and the best performance is reported. For these two methods, more importantly, the label information is only used in the label propagation stage. For SCAN, S^2LRR , S^3R , and SGMK, the label prediction and graph learning are conducted in a unified framework, which often leads to better performance.

As expected, the classification accuracy for all methods monotonically increase with the increase of the percentage of labeled samples. As can be observed, our SGMK method outperforms other state-of-the-art methods in general. This confirms the effectiveness of our proposed method on SSL task. Remember that S^2LRR and S^3R are using self-expressive property to capture the global information, while SCAN is developed to reveal local structure information, so the advantages of our SGMK method over them verify the necessity of incorporating both global and local structure information.

7. Conclusion

In this paper, we propose a new graph learning framework by iteratively learning the graph matrix and the labels. Specifically, both local and global structure information is incorporated in our model. We also consider rank constraint on the graph Laplacian, to yield an optimal graph for clustering and classification tasks, so the achieved graph is more informative and discriminative. This turns out to be a unified model for both graph and label learning, both are improved collaboratively. A multiple kernel learning method is also developed to avoid extensive search for the most suitable kernel. Extensive experiments show the high potential of our method on real-world applications.

Though impressive performance is achieved, the proposed approach has a high time complexity. In the future, we plan to improve its computation efficiency. This can be addressed by borrowing the idea of anchor point. Specifically, we only need to learn a graph between the whole data points and some landmarks. Considering the crucial role of graph in many algorithms, researchers from many other communities could benefit from this line of research.

Acknowledgment

This paper was in part supported by Grants from the National Key R&D Program of China (No. 2018YFC0807500), the Natural Science Foundation of China (Nos. 61806045, U19A2059), the Sichuan Science and Technology Program

under Project 2020YFS0057, the Ministry of Science and Technology of Sichuan Province Program (Nos. 2018GZDZX0048, 20ZDYF0343), the Fundamental Research Fund for the Central Universities under Project ZYGX2019Z015.

8. References

References

- [1] Z. Li, J. Liu, J. Tang, H. Lu, Robust structured subspace learning for data representation, *IEEE transactions on pattern analysis and machine intelligence* 37 (10) (2015) 2085–2098.
- [2] S. Huang, Z. Kang, I. W. Tsang, Z. Xu, Auto-weighted multi-view clustering via kernelized graph learning, *Pattern Recognition* 88 (2019) 174–184.
- [3] S. Yan, D. Xu, B. Zhang, H.-J. Zhang, Q. Yang, S. Lin, Graph embedding and extensions: A general framework for dimensionality reduction, *IEEE transactions on pattern analysis and machine intelligence* 29 (1) (2007) 40–51.
- [4] Z. Zhang, M. Zhao, T. W. Chow, Graph based constrained semi-supervised learning framework via label propagation over adaptive neighborhood, *IEEE Transactions on Knowledge and Data Engineering* 27 (9) (2013) 2362–2376.
- [5] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, P. Vandergheynst, The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains, *IEEE Signal Processing Magazine* 30 (3) (2013) 83–98.
- [6] X. Shen, W. Liu, I. Tsang, F. Shen, Q.-S. Sun, Compressed k-means for large-scale clustering, in: *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [7] S. Huang, Z. Kang, Z. Xu, Auto-weighted multi-view clustering via deep matrix decomposition, *Pattern Recognition* 97 (2020) 107015.

- [8] C.-G. Li, Z. Lin, H. Zhang, J. Guo, Learning semi-supervised representation towards a unified optimization framework for semi-supervised learning, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 2767–2775.
- [9] Z. Kang, H. Pan, S. C. Hoi, Z. Xu, Robust graph learning from noisy data, IEEE Transactions on Cybernetics 50 (5) (2020) 1833–1843.
- [10] A. Y. Ng, M. I. Jordan, Y. Weiss, et al., On spectral clustering: Analysis and an algorithm, Advances in neural information processing systems 2 (2002) 849–856.
- [11] Z. Zhang, F. Li, L. Jia, J. Qin, L. Zhang, S. Yan, Robust adaptive embedded label propagation with weight learning for inductive classification, IEEE transactions on neural networks and learning systems 29 (8) (2017) 3388–3403.
- [12] X. Zhu, Y. Zhu, W. Zheng, Spectral rotation for deep one-step clustering, Pattern Recognition (2019) 107175.
- [13] X. Zhu, Z. Ghahramani, J. D. Lafferty, Semi-supervised learning using gaussian fields and harmonic functions, in: Proceedings of the 20th International conference on Machine learning (ICML-03), 2003, pp. 912–919.
- [14] F. Wang, C. Zhang, T. Li, Clustering with local and global regularization, IEEE Transactions on Knowledge and Data Engineering 21 (12) (2009) 1665–1678.
- [15] F. Nie, G. Cai, X. Li, Multi-view clustering and semi-supervised classification with adaptive neighbours., in: AAAI, 2017, pp. 2408–2414.
- [16] Z. Kang, C. Peng, Q. Cheng, Twin learning for similarity and clustering: A unified kernel approach., in: AAAI, 2017, pp. 2080–2086.
- [17] L. Zhuang, H. Gao, Z. Lin, Y. Ma, X. Zhang, N. Yu, Non-negative low rank and sparse graph for semi-supervised learning, in: Computer Vision

- and Pattern Recognition (CVPR), 2012 IEEE Conference on, IEEE, 2012, pp. 2328–2335.
- [18] C. A. R. de Sousa, S. O. Rezende, G. E. Batista, Influence of graph construction on semi-supervised learning, in: Joint European Conference on Machine Learning and Knowledge Discovery in Databases, Springer, 2013, pp. 160–175.
- [19] X. Peng, J. Feng, J. T. Zhou, Y. Lei, S. Yan, Deep subspace clustering, IEEE Transactions on Neural Networks and Learning Systems.
- [20] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, Y. Ma, Robust recovery of subspace structures by low-rank representation, IEEE Transactions on Pattern Analysis and Machine Intelligence 35 (1) (2013) 171–184.
- [21] M. Maier, U. V. Luxburg, M. Hein, Influence of graph construction on graph-based clustering measures, in: Advances in neural information processing systems, 2009, pp. 1025–1032.
- [22] Z. Kang, C. Peng, Q. Cheng, Clustering with adaptive manifold structure learning, in: Data Engineering (ICDE), 2017 IEEE 33rd International Conference on, IEEE, 2017, pp. 79–82.
- [23] Z. Kang, X. Zhao, Shi, C. Peng, H. Zhu, J. T. Zhou, X. Peng, W. Chen, Z. Xu, Partition level multiview subspace clustering, Neural Networks 122 (2020) 279–288.
- [24] L. Zhang, L. Qiao, S. Chen, Graph-optimized locality preserving projections, Pattern Recognition 43 (6) (2010) 1993–2002.
- [25] C. Hou, C. Zhang, Y. Wu, Y. Jiao, Stable local dimensionality reduction approaches, Pattern Recognition 42 (9) (2009) 2054–2066.
- [26] X. Zhu, X. Li, S. Zhang, C. Ju, X. Wu, Robust joint graph sparse coding for unsupervised spectral feature selection, IEEE transactions on neural networks and learning systems 28 (6) (2017) 1263–1275.

- [27] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, B. Schölkopf, Learning with local and global consistency, in: *Advances in neural information processing systems*, 2004, pp. 321–328.
- [28] Z. Kang, C. Peng, Q. Cheng, Robust pca via nonconvex rank approximation, in: *2015 IEEE International Conference on Data Mining, IEEE*, 2015, pp. 211–220.
- [29] R. Kondor, K. M. Borgwardt, The skew spectrum of graphs, in: *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 496–503.
- [30] S. I. Daitch, J. A. Kelner, D. A. Spielman, Fitting a graph to vector data, in: *Proceedings of the 26th Annual International Conference on Machine Learning*, 2009, pp. 201–208.
- [31] J. Eldridge, M. Belkin, Y. Wang, Graphons, mergeons, and so on!, in: *Advances in Neural Information Processing Systems*, 2016, pp. 2307–2315.
- [32] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, *International Conference on Learning Representations (ICLR)*.
- [33] C. Tang, X. Liu, X. Zhu, J. Xiong, M. Li, J. Xia, X. Wang, L. Wang, Feature selective projection with low-rank embedding and dual laplacian regularization, *IEEE Transactions on Knowledge and Data Engineering*.
- [34] K. Zhan, F. Nie, J. Wang, Y. Yang, Multiview consensus graph clustering, *IEEE Transactions on Image Processing* 28 (3) (2018) 1261–1270.
- [35] C. Zhang, H. Fu, Q. Hu, X. Cao, Y. Xie, D. Tao, D. Xu, Generalized latent multi-view subspace clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (1) (2020) 86–99.
- [36] X. Liu, X. Zhu, M. Li, L. Wang, C. Tang, J. Yin, D. Shen, H. Wang, W. Gao, Late fusion incomplete multi-view clustering, *IEEE transactions on pattern analysis and machine intelligence* 41 (10) (2018) 2410–2423.

- [37] Z. Kang, X. Lu, J. Liang, K. Bai, Z. Xu, Relation-guided representation learning, *Neural Networks* 131 (2020) 93–102.
- [38] J. C. Bezdek, R. J. Hathaway, Convergence of alternating optimization, *Neural, Parallel & Scientific Computations* 11 (4) (2003) 351–368.
- [39] X. Liu, X. Zhu, M. Li, L. Wang, E. Zhu, T. Liu, M. Kloft, D. Shen, J. Yin, W. Gao, Multiple kernel k-means with incomplete kernels, *IEEE transactions on pattern analysis and machine intelligence*.
- [40] L. Du, P. Zhou, L. Shi, H. Wang, M. Fan, W. Wang, Y.-D. Shen, Robust multiple kernel k-means using ℓ_2 , 1-norm, in: *Proceedings of the 24th International Conference on Artificial Intelligence*, AAAI Press, 2015, pp. 3476–3482.
- [41] J. Huang, F. Nie, H. Huang, A new simplex sparse learning model to measure data similarity for clustering, in: *Proceedings of the 24th International Conference on Artificial Intelligence*, AAAI Press, 2015, pp. 3569–3575.
- [42] F. Nie, X. Wang, H. Huang, Clustering and projected clustering with adaptive neighbors, in: *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2014, pp. 977–986.
- [43] H.-C. Huang, Y.-Y. Chuang, C.-S. Chen, Multiple kernel fuzzy clustering, *IEEE Transactions on Fuzzy Systems* 20 (1) (2012) 120–134.
- [44] H. Huang, Y. Chuang, C. Chen, Affinity aggregation for spectral clustering, in: *Computer Vision and Pattern Recognition (CVPR)*, 2012 IEEE Conference on, IEEE, 2012, pp. 773–780.
- [45] C. Peng, Z. Kang, S. Cai, Q. Cheng, Integrate and conquer: Double-sided two-dimensional k-means via integrating of projection and manifold construction, *ACM Transactions on Intelligent Systems and Technology (TIST)* 9 (5) (2018) 1–25.

Table 2: Clustering results of various methods. The average performance of those 12 kernels are put in parenthesis. Single and multiple kernel methods are separated by double lines. The best performance of single and multiple kernel methods are highlighted in boldface. ‘-’ denotes the results are unavailable due to numerical error (text data is sparse).

(a) Accuracy(%)

Data	SC	RKKM	LRR	SSR	Local	Global	SGSK	MKKM	AASC	RMKKM	SGMK
YALE	49.42(40.52)	48.09(39.71)	53.94	54.55	58.79	55.85(45.35)	62.75 (62.05)	45.70	40.64	52.18	63.62
JAFFE	74.88(54.03)	75.61(67.98)	70.89	87.32	98.12	99.83 (86.64)	99.53(98.12)	74.55	30.35	87.07	99.53
ORL	57.96(46.65)	54.96(46.88)	71.50	69.00	61.50	62.35(50.50)	70.05(62.10)	47.51	27.20	55.60	70.02
AR	28.83(22.22)	33.43 (31.20)	32.02	65.00	42.26	56.79(41.35)	62.59(48.21)	28.61	33.23	34.37	63.45
BA	31.07(26.25)	42.17(34.35)	25.93	23.97	36.82	47.72(39.50)	48.32 (37.59)	40.52	27.07	43.42	49.37
TR11	50.98(43.32)	53.03(45.04)	-	41.06	38.89	71.26(54.88)	71.74 (54.92)	50.13	47.15	57.71	74.40
TR41	63.52(44.80)	56.76(46.80)	-	63.78	62.89	67.43(53.13)	72.67 (69.15)	56.10	45.90	62.65	79.38
TR45	57.39(45.96)	58.13(45.69)	-	71.45	56.96	74.02(53.38)	77.54 (75.33)	58.46	52.64	64.00	77.54

(b) NMI(%)

Data	SC	RKKM	LRR	SSR	Local	Global	SGSK	MKKM	AASC	RMKKM	SGMK
YALE	52.92(44.79)	52.29(42.87)	59.39	57.26	57.67	56.50(45.07)	61.58 (60.47)	50.06	46.83	55.58	62.04
JAFFE	82.08(59.35)	83.47(74.01)	75.73	92.93	97.31	99.35 (84.67)	99.18(97.62)	79.79	27.22	89.37	99.18
ORL	75.16(66.74)	74.23(63.91)	85.40	84.23	76.59	78.96(63.55)	82.65(75.93)	68.86	43.77	74.83	81.94
AR	58.37(56.05)	65.44 (60.81)	67.23	84.16	65.73	76.02(59.70)	82.61(67.63)	59.17	65.06	65.49	83.51
BA	50.76(40.09)	57.82(46.91)	40.74	30.29	49.32	63.04 (52.17)	61.94(52.71)	56.88	42.34	58.47	62.25
TR11	43.11(31.39)	49.69(33.48)	-	27.60	19.17	58.60(37.58)	62.07 (38.98)	44.56	39.39	56.08	64.18
TR41	61.33(36.60)	60.77(40.86)	-	59.56	51.13	65.50(43.18)	70.59 (63.67)	57.75	43.05	63.47	69.85
TR45	48.03(33.22)	57.86(38.96)	-	67.82	49.31	74.24 (44.36)	70.7(69.70)	56.17	41.94	62.73	70.92

(c) Purity(%)

Data	SC	RKKM	LRR	SSR	Local	Global	SGSK	MKKM	AASC	RMKKM	SGMK
YALE	51.61(43.06)	49.79(41.74)	55.15	58.18	59.39	57.27(55.79)	66.77 (66.19)	47.52	42.33	53.64	67.79
JAFFE	76.83(56.56)	79.58(71.82)	74.18	96.24	98.12	99.85 (96.53)	99.53(98.17)	76.83	33.08	88.90	99.53
ORL	61.45(51.20)	59.60(51.46)	75.25	76.50	76.59	74(70.37)	75.35(71.62)	52.85	31.56	60.23	77.00
AR	33.24(25.99)	35.87 (33.88)	33.33	69.52	44.64	63.45(62.37)	80.60 (62.54)	30.46	34.98	36.78	83.57
BA	34.50(29.07)	45.28(36.86)	28.70	40.85	39.67	52.36(49.79)	57.36 (55.74)	43.47	30.29	46.27	58.27
TR11	58.79(50.23)	67.93(56.40)	-	85.02	44.20	82.85(80.76)	81.40(80.07)	65.48	54.67	72.93	82.37
TR41	73.68(56.45)	74.99(60.21)	-	75.40	67.54	73.23(71.21)	78.36 (77.19)	72.83	62.05	77.57	87.13
TR45	61.25(50.02)	68.18(53.75)	-	83.62	60.87	78.26(77.76)	78.70(78.06)	69.14	57.49	75.20	78.70

Table 3: Global and local structure effect in multiple kernel learning setting.

Metric	Method	YALE	JAFFE	ORL	AR	BA	TR11	TR41	TR45
Acc	SGMK without Local	56.97	100	65.25	62.38	47.34	73.43	67.31	74.35
	SGMK	63.62	99.53	70.02	63.45	49.37	74.40	79.38	77.54
NMI	SGMK without Local	56.52	100	80.04	81.51	62.94	60.15	65.11	74.97
	SGMK	62.04	99.18	81.94	83.51	62.25	64.18	69.85	70.92
Purity	SGMK without Local	60.00	100	77.00	82.62	52.12	87.44	73.69	78.26
	SGMK	67.79	99.53	77.00	83.57	58.27	82.37	87.13	78.70

Table 4: Classification accuracy (%) on benchmark data sets (mean±standard deviation).

The best results are in bold font.

Data	Labeled Percentage(%)	GFHF	LGC	S ³ R	S ² LRR	SCAN	SGMK
YALE	10	38.00±11.91	47.33±13.96	38.83±8.60	28.77±9.59	45.07±1.30	52.40 ±0.19
	30	54.13±9.47	63.08±2.20	58.25±4.25	42.58±5.93	60.92±4.03	75.58 ±0.04
	50	60.28±5.16	69.56±5.42	69.00±6.57	51.22±6.78	68.94±4.57	82.11 ±0.05
JAFFE	10	92.85±7.76	96.68±2.76	97.33±1.51	94.38±6.23	96.92±1.68	99.57 ±0.02
	30	98.50±1.01	98.86±1.14	99.25±0.81	98.82±1.05	98.20±1.22	99.90 ±0.01
	50	98.94±1.11	99.29±0.94	99.82±0.60	99.47±0.59	99.25±5.79	100 ±0.00
BA	10	45.09±3.09	48.37±1.98	25.32±1.14	20.10±2.51	55.05±1.67	58.77 ±0.83
	30	62.74±0.92	63.31±1.03	44.16±1.03	43.84±1.54	68.84±1.09	89.88 ±0.27
	50	68.30±1.31	68.45±1.32	54.10±1.55	52.49±1.27	72.20±1.44	90.60 ±0.13
COIL20	10	87.74±2.26	85.43±1.40	93.57 ±1.59	81.10±1.69	90.09±1.15	90.74±0.64
	30	95.48±1.40	87.82±1.03	96.52±0.68	87.69±1.39	95.27±0.93	96.85 ±0.32
	50	96.27±0.71	88.47±0.45	97.87±0.10	90.92±1.19	97.53±0.82	98.74 ±0.08