

# Multi-view Low-rank Preserving Embedding: A Novel Method for Multi-view Representation

Xiangzhu Meng, Lin Feng\*, Huibing Wang

**Abstract**—In recent years, we have witnessed a surge of interest in multi-view representation learning, which is concerned with the problem of learning representations of multi-view data. When facing multiple views that are highly related but slightly different from each other, most of existing multi-view methods might fail to fully integrate multi-view information. Besides, correlations between features from multiple views always vary seriously, which makes multi-view representation challenging. Therefore, how to learn appropriate embedding from multi-view information is still an open problem but challenging. To handle this issue, this paper proposes a novel multi-view learning method, named Multi-view Low-rank Preserving Embedding (MvLPE). It integrates different views into one centroid view by minimizing the disagreement term, based on distance or similarity matrix among instances, between the centroid view and each view meanwhile maintaining low-rank reconstruction relations among samples for each view, which could make more full use of compatible and complementary information from multi-view features. Unlike existing methods with additive parameters, the proposed method could automatically allocate a suitable weight for each view in multi-view information fusion. However, MvLPE couldn't be directly solved, which makes the proposed MvLPE difficult to obtain an analytic solution. To this end, we approximate this solution based on stationary hypothesis and normalization post-processing to efficiently obtain the optimal solution. Furthermore, an iterative alternating strategy is provided to solve this multi-view representation problem. The experiments on six benchmark datasets demonstrate that the proposed method outperforms its counterparts while achieving very competitive performance.

**Index Terms**—Multi-view learning, Low-rank preserving, Dimension reduction

## I. INTRODUCTION

In general, one object could be characterized by different kinds of views [1–3], because data is often collected from diverse domains or obtained from different feature extractors. For examples, web pages could be usually presented by the page-text and hyperlink information; Color, text or shape information could be used as different kinds of features, in image and video processing, such as HSV, Local Binary Pattern (LBP) [4], Gist [5], Histogram of Gradients (HoG) [6], Edge Direction Histogram (EDH) [7]. Since different views describe distinct properties of the instance, multiple views contain more complete information than just one view. Generally, it could achieve better performance in many real-world applications [8–13] by taking the complementary information from multiple views into consideration. As we all

know, the performance of machine learning methods is heavily dependent on the expressive power of feature representation. Consequently, multi-view representation learning has received great research efforts, such as multi-view information fusion, multi-view representation alignment, and so on. Those multi-view methods focused on exploiting the diverse information and complementary information among multiple views to achieve a comprehensive representation of the instance.

Nowadays, multi-view representation methods [14–25] have been well studied in many applications. Multi-view information fusion methods [14–20] aimed to fuse multi-view features into single compact representation. Multiview Spectral Embedding (MSE) [14] was an extension of Laplacian Eigenmaps (LE) [26] and incorporated it with multi-view data to find a common low-dimensional subspace, which exploited low-dimensional representations based on the graph. Tang et al. [15] fused the information from multiple graphs with linked matrix factorization, where each graph was approximated by the graph-specific factor and the common factor. Tzortzis et al. [16] expressed each view as a given kernel matrix and learned a weighted combination of those kernels in parallel. Multi-view sparse coding [17, 18] associated the shared latent representation for the multi-view data by a set of linear mappings that are defined as dictionaries. Ngiam et al. [19] proposed a novel method to extract shared representations via training deep auto-encoder [27], which utilized the concatenation of the final hidden coding of audio and video modalities as inputs and mapped these inputs to a shared representation layer. Inspired by the great success of Convolutional Neural Networks (CNN) [28], Su et al. [20] introduced a multi-view CNN for 3D object recognition, which integrated information from multiple 2D views of an object into a single and compact representation. However, those multi-view fusion methods might ignore the consistent correlation information among multiple views so that compatible and complementary information couldn't be made full advantage. To capture the relationships among different views, multi-view representation alignment methods [21–25] were proposed to explore consistent correlation information by feature alignment. In particular, Canonical Correlation Analysis (CCA) [21] and its kernel extension [29] were representative features alignment methods, which could project two views into the common subspace by maximizing the cross correlation between two views. Furthermore, CCA was further generalized for a multi-view scenario termed as multi-view canonical correlation analysis (MCCA) [22]. Kan et al. [23] proposed multi-View Discriminant Analysis to extend Linear Discriminant Analysis (LDA) [30, 31] based on CCA into a multi-view setting, which projected multi-view features to one

X. Meng and L. Feng (corresponding author, denoted by \*) are with the School of Computer Science and Technology, Dalian University of Technology, Dalian, China (xiangzhu\_meng@mail.dlut.edu.cn; fenglin@dlut.edu.cn). H. Wang is with the College of Information Science and Technology, Dalian Maritime University, Dalian, China (huibing.wang@dmlu.edu.cn).

discriminative common subspace. Inspired by the success of deep neural networks [27, 32], Andrew et al. [24] proposed the method of deep CAA to capture the high-level association between multi-view data by coupling the joint representation among multiple views at the higher level. Zhang et al. [25] proposed a Generalized Latent Multi-View Subspace Clustering, which jointly learns the latent representation and multi-view subspace representation within the unified framework. Nevertheless, these alignment methods mainly employed the linear projection to model the cross correlation for forcible alignment of pairwise views so that algorithm performance would be not enough robust when facing such multiple views that were highly related but sometimes different from each other.

It's also attracted wide attention to achieve the multi-view clustering agreement [33–37] to yield a substantial superior clustering performance over the single view paradigm. For example, Kumar et al. [33] proposed a co-regularized multi-view spectral clustering framework that captured complementary information among different viewpoints by co-regularizing the clustering hypotheses. Besides, other works in [38–43] could also obtain promising performance in the multi-view learning environment. Even though the above multi-view methods have achieved promising performance in many applications, most of them could not make use of compatible and complementary information among multiple views or introduce additional learnable parameters in fusing multi-view information. Moreover, the limitations of their generalization and scalability exist all the time.

### A. Contributions

In this paper, we first propose a novel single-view representation method called Low-rank Preserving Embedding (LPE), which provides with three different manners, including direct embedding, linear projection, and kernel method, to maintain the low-rank reconstruction relationships among samples. In this way, we could flexibly choose the embedding manner to preserve the low-rank reconstruction structure under each view when fusing multi-view information. Then we extend LPE into multi-view setting to develop a multi-view method called Multi-view Low-rank Preserving Embedding (MvLPE), which integrates all views into one centroid view by minimizing the disagreement between centroid view and other views and combines it with the low-rank reconstruction structure in each view. Specially, the proposed multi-view method could learn an optimal weight for each view without additive parameters when fusing all views into centroid view, and the obtained the embedding of centroid view could affect the solution of the embeddings of other views in turn. Consequently, both compatible and complementary information from multi-view feature sets and the low-rank reconstruction structure under all instances in each view could be considered at the same time. To obtain the optimal solution, we further design an iterative alternating strategy for the proposed MvLPE and also analyze its convergence. Furthermore, we discuss potential extensions for single-view methods to improve the generalization of our method. Finally, extensive experiments on six

benchmark datasets demonstrate that the proposed MvLPE achieves comparable performance. The major contributions of this paper are summarized as follows:

- We propose a novel single-view representation method providing with three different manners to maintain the low-rank reconstruction relationships among samples, called Low-rank Preserving Embedding (LPE). It could flexibly choose the embedding manner to keep the low-rank reconstruction structure when fusing multi-view information.
- We extend LPE into multi-view setting to propose a novel multi-view method, called Multi-view Low-rank Preserving Embedding (MvLPE), to integrate different information into one centroid view. It considers both compatible and complementary information from multi-view feature sets and the low-rank reconstruction structure under all instances in each view at the same time.
- An effective and robust iterative alternating algorithm is developed to seek an approximate optimal solution for MvLPE. Moreover, we provide with the convergence analysis of this method and its extensions for those single-view methods.
- The experimental results on 6 benchmark datasets demonstrate that the proposed method outperforms its counterparts and achieves comparable performance.

### B. Organization

The rest of the paper is organized as follows: in Section II, we provide briefly some related methods which have attracted extensive attention; in Section III, we describe the construction procedure of MvLPE and optimization algorithm for MvLPE; in Section IV, extensive experiments on text and image datasets demonstrate the effectiveness of our proposed approach; in Section V, we finally conclude this paper.

## II. RELATED WORKS

In this section, we first introduce a low-rank representation method [44], which seeks the low-rank representation among all the candidates that can represent the data samples as linear combinations of the bases in a given dictionary. Then, we review a multi-view learning method called Multi-view Spectral Embedding [14].

### A. Low-Rank Representation

Liu et al. [44] proposed a representative low-rank representation method to handle the subspace recovery problem, which was quite superior in terms of its effectiveness, intuitiveness and robustness to noise corruptions. Assume that we are provided a features set consisting of  $N$  samples, which are extracted from the  $v$ th view. We denote the features set in the  $v$ th view as  $\mathbf{X}^v = [\mathbf{x}_1^v, \mathbf{x}_2^v, \dots, \mathbf{x}_N^v]$ . When we choose the matrix  $\mathbf{X}^v$  itself as a dictionary that linearly spans the data space. We could get the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{Z}^v, \mathbf{E}^v} \text{rank}(\mathbf{Z}^v) + \lambda \|\mathbf{E}^v\|_{2,1} \\ & \text{s.t. } \mathbf{X}^v = \mathbf{X}^v \mathbf{Z}^v + \mathbf{E}^v \end{aligned} \quad (1)$$

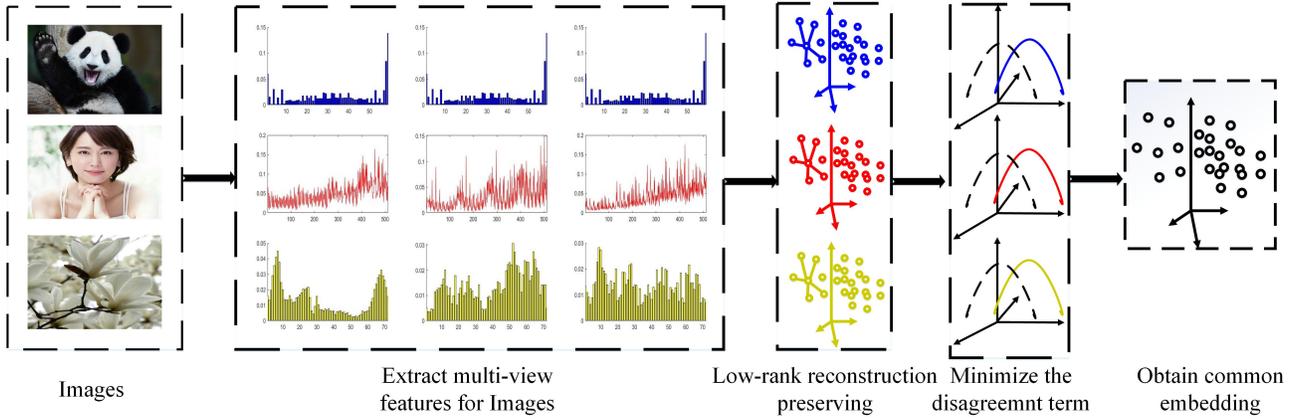


Fig. 1: The working procedure of Multi-view Low-rank Preserving Embedding (MvLPE), which aims to handle multi-view representation problem. Taking the images as an example, we first extract LBP, GIST, and EDH descriptors for images as multi-view information. Then, we apply the representation method called Low-rank Preserving Embedding (LPE) for images to obtain low-rank reconstruction structure and their embeddings. By minimizing the disagreement term between centroid view and all views, we integrate compatible and complementary information from multi-view feature sets to construct common embedding as multi-view representation. Finally, an iterative alternating strategy is adopted to find the optimal solution for MvLPE.

where  $\lambda$  is a hyperparameter and  $Z^v \in \mathbb{R}^{N \times N}$  is the lowest rank representation of data  $X^v$ . After obtaining an optimal solution, we could recover the original data by using  $X^v Z^v$ . Since  $\text{rank}(X^v Z^v) \leq \text{rank}(X^v)$ ,  $X^v Z^v$  is also a low-rank recovery to the original data. By choosing an appropriate dictionary, as we will see, the lowest-rank representation can recover the underlying row space so as to reveal the true segmentation of data. Therefore, LRR could handle well the data drawn from a union of multiple subspaces.

### B. Multi-view Spectral Embedding

Multi-view Spectral Embedding (MSE) [14] is a spectral embedding method that could map different features from multiple views into a common subspace. It aims to find a low-dimensional and physically meaningful embedding. Assume that given data has  $m$  views. Let  $X^v = \{x_1^v, x_2^v, \dots, x_N^v\}$  denote the features set in the  $v$ th view. MSE aims to find a low-dimensional embedding  $U$  as follows:

$$\begin{aligned} & \max_{U, \alpha^v \in \mathbb{R}^{N \times k}} \sum_{v=1}^m \alpha^v \text{tr}(U^T L^v U) \\ & \text{s.t.} \quad U^T U = I, \sum_{v=1}^m \alpha^v = 1, \forall 1 \leq v \leq m \end{aligned} \quad (2)$$

where  $L^v$  denotes the normalized graph Laplacian matrix in the  $v$ th view,  $\alpha = [\alpha^1, \alpha^2, \dots, \alpha^m]$  is a non-negative weight vector. And  $\alpha^v$  reflects the importance which the view  $X^v$  plays in learning to obtain the low-dimensional embedding. Global coordinate alignment is utilized such that low-dimensional embedding in different views could keep consistent with each other globally. And, to solve the above problem, an iterative method could be adopted to update  $\alpha$  and  $U$  respectively.

TABLE I: Important notations used in this paper

Notation	Description
$X^v$	The set of all instances in the $v$ th view
$x_i^v$	The $i$ th instance in the $v$ th view
$d^v$	Dimension of the subspace in the $v$ -th view
$Z^v$	The low-rank reconstruction weights in the $v$ th view
$M^v$	The low-rank reconstructive matrix for the $v$ -th view
$w^v$	The projection direction for the $v$ -th view
$K^v$	The kernel matrix for the $v$ -th view
$\psi(X^v)$	The set of all low-dimensional embedding in the $v$ th view
$w^v$	The weight coefficient for the $v$ -th view
$U^*$	The low-dimensional embedding for the centroid view
$\mathbf{1}_{1/N}$	The $N \times N$ matrix that each element is filled with $1/N$
$I_N$	The $N \times N$ identity matrix
$I_{d^v}$	The $d^v \times d^v$ identity matrix
$m$	The number of views
$N$	The number of samples

## III. MULTI-VIEW LOW-RANK PRESERVING EMBEDDING

In this section, we first propose a new single-view representation method called Low-rank Preserving Embedding (LPE), which provides with three different manners to maintain the low-rank reconstruction relationship among samples. Then, we extend LPE into the multi-view setting to propose a multi-view method called Multi-view Low-rank Preserving Embedding (MvLPE), which fully integrates compatible and complementary information from multi-view features sets to construct common embedding for all views. Then, an iterative alternating strategy is derived to find the optimal solution for our method and the optimization procedure is illustrated in detail. Fig.1 shows the working procedure of MvLPE. Moreover, we provide the convergence discussion of our method in detail and the extensions for those single-view methods according to our proposed MvLPE. For convenience, the important notations in this paper are listed in Table I.

### A. Low-rank Preserving Embedding

Recall that we are provided a features set consisting of  $N$  samples, which is extracted from the  $v$ th view. We denote the features set in the  $v$ th view as  $\mathbf{X}^v = [\mathbf{x}_1^v, \mathbf{x}_2^v, \dots, \mathbf{x}_N^v]$ . As discussed in Section II-A, LRR employs  $\mathbf{X}^v$  itself as a dictionary, which exists such two issues consisting of unsuitable correlation reconstruction and high computational cost caused by the number of instances. Inspired by the fact that the compact combination of samples always lies on the local subspace of the test sample, we replace dictionary  $\mathbf{X}^v$  with the near neighbors set responding to each sample. As a result, we could achieve more ideal space reconstruction than LRR. Meanwhile, the issue of high computational cost could be handled by choosing its near neighbors for individual sample, which could significantly reduce the scale of the dictionary. Therefore, it's feasible and necessary to choose such a dictionary for individual sample by using  $K$  its near neighbors. Combining this with the low-rank hypothesis, we could get the following optimization problem:

$$\begin{aligned} \min_{\mathbf{Z}^v, \mathbf{E}^v} \text{rank}(\mathbf{Z}^v) + \lambda \|\mathbf{E}^v\|_F^2 \\ \text{s.t. } \mathbf{X}_i^v = \widetilde{\mathbf{X}}_i^v \mathbf{Z}_i^v + \mathbf{E}_i^v, \forall 1 \leq i \leq N \end{aligned} \quad (3)$$

where  $\widetilde{\mathbf{X}}_i^v$  is the dictionary of  $i$ th sample consisting of  $K$  its closed neighbors,  $\mathbf{Z}_i^v$  and  $\mathbf{E}_i^v$  denote the  $i$ th column data in the matrix  $\mathbf{Z}^v$  and  $\mathbf{E}^v$  respectively. As a common practice in rank minimization problems, we replace the rank function with the nuclear norm and subject to the constraints the columns of the matrix  $\mathbf{Z}_i^v$  sum to one. By this means, it is deduced to the following optimization problem:

$$\begin{aligned} \min_{\mathbf{Z}^v, \mathbf{E}^v} \|\mathbf{Z}^v\|_* + \lambda \|\mathbf{E}^v\|_F^2 \\ \text{s.t. } \mathbf{X}_i^v = \widetilde{\mathbf{X}}_i^v \mathbf{Z}_i^v + \mathbf{E}_i^v, \mathbf{Z}_i^{vT} \mathbf{1} = 1, \forall 1 \leq i \leq N \end{aligned} \quad (4)$$

And we aim to maintain the low-rank reconstruction relationships among samples, which are obtained by Eq.(4). For the convenience of modeling and solving, a simple trick is used to transform the matrix  $\mathbf{Z}^v \in \mathbb{R}^{K \times N}$  into a matrix  $\mathbf{M}^v \in \mathbb{R}^{N \times N}$ , which fills column elements in the matrix  $\mathbf{M}^v$  according to the low-rank coefficients  $\mathbf{Z}_i^v$  of its neighbors and fills zeros into other elements. Accordingly, we could define the following objective function to seek low dimensional embedding while maintaining the low-rank reconstruction relationships:

$$\begin{aligned} \min_{\mathbf{U}^v} \text{tr} \left( \mathbf{U}^v (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \mathbf{U}^{vT} \right) \\ \text{s.t. } \mathbf{U}^v \mathbf{U}^{vT} = \mathbf{I}_{d^v} \end{aligned} \quad (5)$$

where  $\mathbf{U}^v \in \mathbb{R}^{d^v \times N}$  denotes the embedding in the  $v$ th view,  $d^v$  is the dimension of  $\mathbf{U}^v$ , and  $\text{tr}(\cdot)$  denotes the matrix trace. Furthermore, we further propose two additional variants, which are based on linear transform and kernel trick respectively.

Suppose that  $\mathbf{W}^v$  is a transformation matrix, that is  $\mathbf{U}^v = \mathbf{W}^{vT} \mathbf{X}^v$ , which is a linear approximation. By simple algebra

formulation, the objective function in Eq.(5) can be expressed as follows:

$$\begin{aligned} \min_{\mathbf{W}^v} \text{tr} \left( \mathbf{W}^{vT} \mathbf{X}^v (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \mathbf{X}^{vT} \mathbf{W}^v \right) \\ \text{s.t. } \mathbf{W}^{vT} \mathbf{X}^v \mathbf{X}^{vT} \mathbf{W}^v = \mathbf{I}_{d^v} \end{aligned} \quad (6)$$

Furthermore, suppose that the Euclidean space is mapped to a Hilbert space, that is  $\mathbf{X}_\phi^v = [\phi^v(\mathbf{x}_1^v), \phi^v(\mathbf{x}_2^v), \dots, \phi^v(\mathbf{x}_N^v)]$ , where  $\phi^v(\cdot)$  is a nonlinear map. It has been verified [45] that  $\mathbf{W}_\phi^v$  is that mapped space spanned by  $\phi^v(\mathbf{x}_1^v), \phi^v(\mathbf{x}_2^v), \dots, \phi^v(\mathbf{x}_N^v)$ . Consequently,  $\mathbf{W}_\phi^v$  could be expressed as follows:

$$\mathbf{W}_\phi^v = \sum_{i=1}^N \phi^v(\mathbf{x}_i^v) \beta_i^v = \mathbf{X}_\phi^v \beta^v \quad (7)$$

where  $\beta^v = [\beta_1^v, \beta_2^v, \dots, \beta_N^v]^T \in \mathbb{R}^{N \times d^v}$  consists of the expansion coefficients. Set  $\mathbf{K}_\phi^v = \mathbf{X}_\phi^{vT} \mathbf{X}_\phi^v$ . Combining this with the Eq.(6), we could obtain the following low-rank preserving problem based on kernel:

$$\begin{aligned} \min_{\beta^v} \text{tr} \left( \beta^{vT} \mathbf{K}_\phi^v (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \mathbf{K}_\phi^v \beta^v \right) \\ \text{s.t. } \beta^{vT} \mathbf{K}_\phi^v \mathbf{K}_\phi^v \beta^v = \mathbf{I}_{d^v} \end{aligned} \quad (8)$$

In terms of the discussion above, we provide with three manners to obtain the low-dimensional embedding based on preserving the low-rank reconstruction relation among the samples. Therefore, we could obtain a unified low-rank preserving embedding method, including direct embedding, linear transform, and kernel method, so that we could more flexibly choose embedding manner to fully preserve the low-rank reconstruction relations information among samples when fusing multi-view features based on LPE. For convenience, we utilize  $\psi(\mathbf{X}^v)$  to generally stand for the low dimensional embedding obtained by low-rank preserving embedding method. That is  $\psi(\mathbf{X}^v) = \mathbf{U}^v$ ,  $\mathbf{W}^{vT} \mathbf{X}^v$ , or  $\beta^{vT} \mathbf{K}_\phi^v$ , which is responding to different modes of low-rank preserving embedding respectively.

### B. The construction process of Multi-view Low-rank Preserving Embedding

When facing with multi-view problems, solving the problem for all views separately will fail to integrate multi-view features and make favorable use of the complementary information from multiple views. For solving this problem, we propose a multi-view method called Multi-view Low-rank Preserving Embedding (MvLPE), to fully apply all features from different views into one centroid view and learn common representations, which extends LPE into the multi-view setting. However, the dimension of the features set in each view owns its size, which is different from the other views. Besides, it isn't easy to obtain common embedding directly because of its intrinsic geometric properties in each view. Therefore, integrating different views into one centroid view is still full of challenges.

Inspired by these works [33, 43], we firstly make such hypothesis that similarities among the instances in each view and the centroid view should be consistent under the novel

representations. This hypothesis means that all similarity matrices from the  $v$ th views should be consistent with the similarity of the centroid view by aligning the similarities matrix computed from the centroid view and the  $v$ th view. Noting that,  $\mathbf{U}^*$  in the centroid view and  $\psi(\mathbf{X}^v)$  in the  $v$ th view have different dimensions  $d^*$  and  $d^v$ . To implement this hypotheses and deal with the dimensional difference problem, we utilize the following cost function as a measurement of agreement between the centroid view and the  $v$ th view:

$$\text{Sim}(\mathbf{U}^*, \psi(\mathbf{X}^v)) = -\|\mathbf{K}^* - \mathbf{K}^v\|_F^p \quad (9)$$

where  $\mathbf{K}^*$  and  $\mathbf{K}^v$  stand for the similarity matrix of the centroid view and the  $v$ th view respectively,  $\|\cdot\|_F^p$  denotes the exponential function of Frobenius norm ( $F$ -norm), and  $0 < p \leq 2$  is a scalar. With the change of the value of  $p$ , a series of exponential function could be utilized. In fact, we could choose the general kernel function as our similarity measurement, such as linear kernel, polynomial kernel, Gaussian kernel and so on. For example, when we choose linear kernel as similarity measurement in the centroid view,  $\mathbf{K}^v(\mathbf{U}_i^*, \mathbf{U}_j^*) = \mathbf{U}_i^{*T} \mathbf{U}_j^*$  denotes the similarity between the instance  $\mathbf{U}_i^*$  and the instance  $\mathbf{U}_j^*$ . In this way,  $\text{Sim}(\mathbf{U}^*, \psi(\mathbf{X}^v))$  reflects the consensus measure of the pairwise similarity among all instances under the centroid view and the  $v$ th view.

To further express the consensus term, we expand Eq.(9) as follows:

$$\begin{aligned} \text{Sim}(\mathbf{U}^*, \psi(\mathbf{X}^v)) &= -\|\mathbf{K}^* - \mathbf{K}^v\|_F^p \\ &= (-\text{tr}(\mathbf{K}^{*T} \mathbf{K}^* + \mathbf{K}^{vT} \mathbf{K}^v - \mathbf{K}^{*T} \mathbf{K}^v - \mathbf{K}^{vT} \mathbf{K}^*))^{\frac{p}{2}} \\ &= (2\text{tr}(\mathbf{K}^* \mathbf{K}^v) - \text{tr}(\mathbf{K}^* \mathbf{K}^*) - \text{tr}(\mathbf{K}^v \mathbf{K}^v))^{\frac{p}{2}} \end{aligned} \quad (10)$$

In the Eq.(10), the second term  $\text{tr}(\mathbf{K}^* \mathbf{K}^*)$  and the third term  $\text{tr}(\mathbf{K}^v \mathbf{K}^v)$  in the above equation just depend on individual view, so these two terms couldn't work in integrating two different views. Consequently, we could approximate the consensus term as follows:

$$\text{Sim}(\mathbf{U}^*, \psi(\mathbf{X}^v)) = (\text{tr}(\mathbf{K}^* \mathbf{K}^v))^{\frac{p}{2}} \quad (11)$$

Even though the consensus term in Eq.(11) could work in integrating multi-view information, it's full of challenge to choose suitable kernel function for centroid view and each view at the same time under the consideration for the solving process and effectiveness. Besides, how to directly solve the consensus term based on such complicated kernel function is not easy. Therefore, constructing a meaningful and feasible term that reflects the consistent information between centroid view and each view is very necessary. Inspired by the above hypothesis minimizing the gap between the similarities computed in the centroid view and the similarities in the  $v$ th view, we expect that the distance between two instances in the centroid view is expected to be smaller if the similarity between two instances in the  $v$ th view is larger. In this way, we could formulate the following disagreement term by utilizing the square of

Euclidean distance to substitute the matrix similarity among all instances in the centroid view:

$$\begin{aligned} \text{Dis}(\mathbf{U}^*, \psi(\mathbf{X}^v)) &= \left( \sum_{i,j=1}^n \|\mathbf{U}_i^* - \mathbf{U}_j^*\|_2^2 \mathbf{K}_{ij}^v \right)^{\frac{p}{2}} \\ &= \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}} \end{aligned} \quad (12)$$

where  $\mathbf{D}^v = \text{diag}(d_{11}^v, d_{22}^v, \dots, d_{NN}^v)$  is a diagonal matrix,  $d_{NN}^v = \sum_{i=1}^N \mathbf{K}_{iN}^v$ . Accordingly, we just consider the choice of kernel function for each view but centroid view, which is more convenient to solve.

To integrate rich information among different features, we could obtain the following optimization problem by adding up cost function in Eq.(12) among all views:

$$\begin{aligned} \min_{\mathbf{U}^*} \sum_{v=1}^m \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}} \\ \text{s.t. } \mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*} \end{aligned} \quad (13)$$

The Lagrange function of Eq.(13) could be written as follows:

$$\sum_{v=1}^m \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}} + \mathcal{G}(\mathcal{A}, \mathbf{U}^*) \quad (14)$$

where  $\mathcal{A}$  is the Lagrange multiplier,  $\mathcal{G}(\mathcal{A}, \mathbf{U}^*)$  is the formalized term derived from constraints. Taking the derivative of Eq.(14) w.r.t  $\mathbf{U}^*$  and setting the derivative to zero, we have

$$\sum_{v=1}^m \mathbf{w}^v \frac{\partial \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)}{\mathbf{U}^*} + \frac{\partial \mathcal{G}(\mathcal{A}, \mathbf{U}^*)}{\mathbf{U}^*} = 0 \quad (15)$$

where

$$\mathbf{w}^v = \frac{p}{2} \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}-1} \quad (16)$$

It's easy to find that  $\mathbf{w}^v > 0$  is depended on the target variable  $\mathbf{U}^*$ , so Eq.(15) couldn't be directly solved. If  $\mathbf{w}^v$  is set to be stationary, Eq.(13) could be considered as the solution of the following equation:

$$\begin{aligned} \min_{\mathbf{U}^*} \sum_{v=1}^m \mathbf{w}^v \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right) \\ \text{s.t. } \mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*} \end{aligned} \quad (17)$$

To further analyze the  $\mathbf{w}^v$ , we add the normalization on  $\mathbf{w}^v$  in Eq.(17) after calculating  $\mathbf{w}^v$  by Eq.(16), i.e.  $\sum_{v=1}^m \mathbf{w}^v = 1$ . If the  $v$ th view is close to the centroid view, then  $\text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)$  should be small, thus the learned weight  $\mathbf{w}^v$  for the  $v$ th view is large. Accordingly, such view that isn't close to the centroid view will be assigned a small weight. Therefore, our method optimizes the weight  $\mathbf{w}$  meaningfully. Accordingly,  $\mathbf{w}^v$  could be realized as the weights of different views, which play different contribution in obtaining the common embedding  $\mathbf{U}^*$ .

To further utilize low-rank reconstruction structure information in each view, we expect that the low-dimensional embedding in each view could also be adjusted by minimizing the disagreement measurement against the centroid view rather than only obtained by its low-rank structure. As a result, not only the low-rank structure in this view could be considered but complementary information from other views and centroid view would be utilized when solving the low-dimensional embedding in each view. Therefore, combining the loss function in Eq.(17) with the LPE objectives across all views, we can get the joint loss function for MvLPE as follows:

$$\begin{aligned} & \min_{\mathbf{U}^*, \psi(\mathbf{X}^1), \psi(\mathbf{X}^2), \dots, \psi(\mathbf{X}^m)} \gamma \sum_{v=1}^m \mathbf{w}^v \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right) + \\ & \sum_{v=1}^m \text{tr} \left( \psi(\mathbf{X}^v) (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \psi(\mathbf{X}^v)^T \right) \\ & \text{s.t. } \mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*}, \psi(\mathbf{X}^v) \psi(\mathbf{X}^v)^T = \mathbf{I}_{d^v}, \forall 1 \leq v \leq m \end{aligned} \quad (18)$$

where  $\gamma$  is a hyperparameter that controls the trade-off between the two terms of equation (18). The first term is the agreement between the centroid and all views to follow the multi-view subspace hypotheses. The second term is the sum of LPE loss function for all views. From Eq. (18), we could find that different embedding  $\psi(\mathbf{X}^v)$  inflects each other for the centroid representations. Differing from those fusion methods with additional parameters, our proposed method could automatically assign an optimal weight for each view according to theoretical explanations. Besides, the disagreement term based on distance or similarity matrix encourages to keep consistency between centroid view and other views, which is more robust and scalable than those multi-view representation methods based on features alignment. Therefore, the process of minimizing Eq. (18) aims to find the common embedding which could integrate features from multiple views and preserve low-rank structure among instances.

### C. Optimization Process for MvLPE

In this section, we provide the optimization process for MvLPE in detail. In order to find the optimal solution of Eq.(18), we develop an algorithm based on alternative strategy, which separates the problem into several sub-problems such that each sub-problem is tractable. That is, we alternatively update each variable when fixing others. And we summarized the optimization process in Algorithm 1.

**Updating  $\mathbf{U}^*$ :** By fixing all variables but  $\mathbf{U}^*$ , Eq.(18) will reduce to the following equation without considering constant additive and scaling term:

$$\begin{aligned} & \min_{\mathbf{U}^*} \sum_{v=1}^m \mathbf{w}^v \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right) \\ & \text{s.t. } \mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*} \end{aligned} \quad (19)$$

which has a feasible solution. According to the operational rules of matrix trace, the above equation could be transformed

as follows:

$$\begin{aligned} & \min_{\mathbf{U}^*} \text{tr} \left( \mathbf{U}^* \left( \sum_{v=1}^m \mathbf{w}^v (\mathbf{D}^v - \mathbf{K}^v) \right) \mathbf{U}^{*T} \right) \\ & \text{s.t. } \mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*} \end{aligned} \quad (20)$$

Set  $L^* = \sum_{v=1}^m \mathbf{w}^v (\mathbf{D}^v - \mathbf{K}^v)$ . Therefore, with the constraint  $\mathbf{U}^* \mathbf{U}^{*T} = \mathbf{I}_{d^*}$ , the optimal  $\mathbf{U}^*$  could be solved by eigen-decomposition.  $\mathbf{U}^*$  consists of eigenvectors corresponding to the smallest  $d^*$  eigenvalues.

**Updating  $\psi(\mathbf{X}^v)$ :** By fixing all variables but  $\mathbf{U}^*$ , Eq.(18) will reduce to the following equation :

$$\begin{aligned} & \min_{\psi(\mathbf{X}^v)} \text{tr} \left( \psi(\mathbf{X}^v) (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \psi(\mathbf{X}^v)^T \right) \\ & + \gamma \mathbf{w}^v \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right) \\ & \text{s.t. } \psi(\mathbf{X}^v) \psi(\mathbf{X}^v)^T = \mathbf{I}_{d^v} \end{aligned} \quad (21)$$

Noting that the above equation isn't easy to be directly solved, because the expression of  $K^v$  isn't readily certain and the disagreement term in Eq.(12) is unsymmetric, that is  $Dis(\mathbf{U}^*, \psi(\mathbf{X}^v)) \neq Dis(\psi(\mathbf{X}^v), \mathbf{U}^*)$ . Inspired by co-training methods, which limit the search for the compatible hypothesis that predict the same labels for co-occurring in each view, we utilize the  $Dis(\psi(\mathbf{X}^v), \mathbf{U}^*)$  as the disagreement measurement between the  $v$ th view and the centroid view rather than  $Dis(\mathbf{U}^*, \psi(\mathbf{X}^v))$  when fixing the centroid view  $\mathbf{U}^*$ . Based on the above assumption that  $\mathbf{w}^v$  is set to be stationary, the above equation could be further transformed as follows:

$$\begin{aligned} & \min_{\psi(\mathbf{X}^v)} \text{tr} \left( \psi(\mathbf{X}^v) (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) \psi(\mathbf{X}^v)^T \right) \\ & + \gamma \mathbf{w}^v \text{tr} \left( \psi(\mathbf{X}^v) (\mathbf{D}^* - \mathbf{K}^*) \psi(\mathbf{X}^v)^T \right) \\ & \text{s.t. } \psi(\mathbf{X}^v) \psi(\mathbf{X}^v)^T = \mathbf{I}_{d^v} \end{aligned} \quad (22)$$

where  $\mathbf{K}^*$  stands for the similarity matrix of the centroid view, and  $\mathbf{D}^* = \text{diag}(\mathbf{d}_{11}^*, \mathbf{d}_{22}^*, \dots, \mathbf{d}_{NN}^*)$  is a diagonal matrix,  $\mathbf{d}_{NN}^* = \sum_{i=1}^N \mathbf{K}_{iN}^*$ . Set  $L^v = (\mathbf{I}_N - \mathbf{M}^v)^T (\mathbf{I}_N - \mathbf{M}^v) + \gamma \mathbf{w}^v (\mathbf{D}^* - \mathbf{K}^*)$ . Therefore, with the constraint  $\psi(\mathbf{X}^v) \psi(\mathbf{X}^v)^T = \mathbf{I}_{d^v}$ , the optimal  $\psi(\mathbf{X}^v)$  could be solved by eigen-decomposition.  $\psi(\mathbf{X}^v)$  consists of eigenvectors corresponding to the smallest  $d^v$  eigenvalues.

**Updating  $\mathbf{w}$ :** By fixing all variables but  $\mathbf{w}^v$ , we could calculating  $\mathbf{w}^v$  by Eq.(16) and normalization for each view as follows:

$$\mathbf{w}^v = \frac{\text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}-1}}{\sum_{v=1}^m \text{tr} \left( \mathbf{U}^* (\mathbf{D}^v - \mathbf{K}^v) \mathbf{U}^{*T} \right)^{\frac{p}{2}-1}} \quad (23)$$

It's notable that the value of  $p$  could directly influence the weighting factor  $\mathbf{w}$ . When  $p \rightarrow 0$ ,  $\mathbf{w}^v$  is proportional to the reciprocal of the disagreement term between the  $v$ th view and centroid view. Conversely, when  $p \rightarrow 2$ , all elements in  $\mathbf{w}$  tend to be equal to  $1/m$ .

**Algorithm 1** The optimization procedure of MVLPE**Require:**

1. A multi-view features set with N training samples having m views  $\mathbf{X}^v = [\mathbf{x}_1^v, \mathbf{x}_2^v, \dots, \mathbf{x}_N^v] \in \mathbb{R}^{D_v \times N}$ ,  $v = 1, 2, \dots, m$ .
2. Set the parameters  $\gamma$  and  $p$  in Eq.(18).

**The Main Procedure:****for**  $v=1:m$  **do**

3. Initialize  $\mathbf{w}^v = 1/m$ .
4. Specialize the  $M^v$  in Eq.(4).
5. Initialize  $\psi(\mathbf{X}^v)$  according to  $M^v$

**end for****repeat**

6. Update  $\mathbf{U}^*$  by solving Eq.(21).
- for**  $v=1:m$  **do**
  7. Update  $\psi(\mathbf{X}^v)$  for the  $v$ th view by solving Eq.(22).
- end for**
8. Update  $\mathbf{w}$  by solving Eq.(16).

**until**  $[\mathbf{U}^*, \psi(\mathbf{X}^1), \psi(\mathbf{X}^2), \dots, \psi(\mathbf{X}^m)]$  converges.

**return**  $\mathbf{U}^*$ .

*D. Convergence Analysis*

Because our proposed MvLPE is solved by alternating optimization strategy, it's essential to analyze its convergence. We first need to utilize the following lemma introduced by [46].

**Lemma 1.** For any positive number a and b, the following inequality holds:

$$a^{\frac{p}{2}} - \frac{p}{2} \frac{a}{b^{1-\frac{p}{2}}} \leq b^{\frac{p}{2}} - \frac{p}{2} \frac{b}{a^{1-\frac{p}{2}}} \quad (24)$$

**Theorem 1.** Each updated  $\mathbf{U}^*$  in **Algorithm 1** will monotonically decrease the objective in Eq.(13) in each iteration.

**Proof:** We use  $\tilde{\mathbf{U}}^*$  to denote the updated  $\mathbf{U}^*$  in each iteration. According to the optimization to  $\mathbf{U}^*$  in **Algorithm 1**, we know that  $\tilde{\mathbf{U}}^*$  makes the objective of Eq.(21) have the smaller than  $\mathbf{U}^*$ . Combining  $\mathbf{w}$  computed in **Algorithm 1**, we could drive:

$$\begin{aligned} & \sum_{v=1}^m \frac{p}{2} \frac{\text{tr}(\tilde{\mathbf{U}}^*(\mathbf{D}^v - \mathbf{K}^v)\tilde{\mathbf{U}}^{*T})}{\text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{1-\frac{p}{2}}} \\ & \leq \sum_{v=1}^m \frac{p}{2} \frac{\text{tr}(\mathbf{U}^{*T}(\mathbf{D}^v - \mathbf{K}^v))}{\text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{1-\frac{p}{2}}} \end{aligned} \quad (25)$$

According to **Lemma 1**, we have:

$$\begin{aligned} & \sum_{v=1}^m \text{tr}(\tilde{\mathbf{U}}^*(\mathbf{D}^v - \mathbf{K}^v)\tilde{\mathbf{U}}^{*T})^{\frac{p}{2}} - \sum_{v=1}^m \frac{p}{2} \frac{\text{tr}(\tilde{\mathbf{U}}^*(\mathbf{D}^v - \mathbf{K}^v)\tilde{\mathbf{U}}^{*T})}{\text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{1-\frac{p}{2}}} \\ & \leq \sum_{v=1}^m \text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{\frac{p}{2}} - \sum_{v=1}^m \frac{p}{2} \frac{\text{tr}(\mathbf{U}^{*T}(\mathbf{D}^v - \mathbf{K}^v))}{\text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{1-\frac{p}{2}}} \end{aligned} \quad (26)$$

Sum over Eq.(25) and Eq.(26) in the two sides, we could derive:

$$\sum_{v=1}^m \text{tr}(\tilde{\mathbf{U}}^*(\mathbf{D}^v - \mathbf{K}^v)\tilde{\mathbf{U}}^{*T})^{\frac{p}{2}} \leq \sum_{v=1}^m \text{tr}(\mathbf{U}^*(\mathbf{D}^v - \mathbf{K}^v)\mathbf{U}^{*T})^{\frac{p}{2}} \quad (27)$$

Thus the alternating optimization will monotonically decrease the objective in Eq.(13).

**Theorem 2.** The objective function in Eq.(18) is bounded. The proposed optimization algorithm monotonically decreases the loss value in each step, which makes the solution converge to a local optimum.

**Proof:** It's easy to find that there must exist one view which can make  $e_{min} = \text{tr}(\mathbf{U}^v(\mathbf{I}_N - \mathbf{M}^v)^T(\mathbf{I}_N - \mathbf{M}^v)\mathbf{U}^{vT}) > 0$  to be smallest among all views. Similarly, we also find such a view that is closest to the centroid view, that is  $d_{min} > 0$ . Because the hyperparameter  $\gamma > 0$ , it is provable that the objective value in Eq.(13) is greater than  $m(e_{min} + d_{min})$ . Therefore, The objective function in Eq.(18) has a lower bound.

For **Algorithm 1**, it's obvious to see that  $\{\psi(\mathbf{X}^1), \psi(\mathbf{X}^2), \dots, \psi(\mathbf{X}^m)\}$  generated via solving Eq.(22) are the exact minimum points of Eq.(22) respectively. As a result, the value of the objective function on  $\{\psi(\mathbf{X}^1), \psi(\mathbf{X}^2), \dots, \psi(\mathbf{X}^m)\}$  in Eq.(18) is decreasing in each iteration of **Algorithm 1**. Combining this with **Theorem 1**, thus the alternating optimization will monotonically decrease the objective in Eq.(18). Therefore, according to the bounded monotone convergence theorem [47] that asserts the convergence of every bounded monotone sequence, the proposed optimization algorithm converges.

*E. Extensions*

Differing from these multi-views methods existing limitations of generalization and scalability, we could extend those single view-based methods, which could be cast as a special form of the quadratically constrained quadratic program (QCQP), into multi-view setting referring to MvLPE. Specially, such methods [13, 26, 30, 48–51] could be solved by the following equation:

$$\begin{aligned} & \min_{\mathbf{U}^v} \text{tr}(\mathbf{U}^v \mathbf{Q}^v \mathbf{U}^{vT}) \\ & \text{s.t. } \mathbf{U}^v \mathbf{C}^v \mathbf{U}^{vT} = \mathbf{I} \end{aligned} \quad (28)$$

where  $\mathbf{M}^v \in \mathbb{R}^{N \times N}$  reflects the intrinsic structure for the  $v$ th view and  $\mathbf{C}^v$  stands for the different constraint term according to different methods. For the example of LPE, we could utilize  $\mathbf{I}_N$  to reformulate  $\mathbf{C}^v$ . Taking LDA [30], NPE [49], and LPP [50] as examples, we could express  $\mathbf{M}^v$  and  $\mathbf{C}^v$  as follows:

- **LDA:**  $\mathbf{Q}_{i,j}^v = -1/N_v^c$  if  $\mathbf{X}_i^v$  and  $\mathbf{X}_j^v$  belong to the class  $c$ , 0 otherwise, where  $N_v^c$  is the number of samples for class  $c$  in the  $v$ th view. And  $\mathbf{C}^v = \mathbf{M}^v - \mathbf{I}_N$ , where  $\mathbf{I}_N \in \mathbb{R}^{N \times N}$  is an identity matrix.
- **NPE:**  $\mathbf{Q}^v = (\mathbf{I}_N - \mathbf{S}^v)^T(\mathbf{I}_N - \mathbf{S}^v)$ , where  $\mathbf{S}^v \in \mathbb{R}^{N \times N}$  is the reconstruction coefficient matrix in the  $v$ th view. And  $\mathbf{C}^v = \mathbf{I}_N$ , where  $\mathbf{I}_N \in \mathbb{R}^{N \times N}$  is an identity matrix.
- **LPP:**  $\mathbf{Q}^v$  is the Laplacian matrix in the  $v$ th view and  $\mathbf{C}^v$  is a diagonal matrix, in which  $\mathbf{C}_{ii}^v$  is the sum of all elements in the  $i$ th row of  $\mathbf{Q}^v$ .

To further improve the performance of those single-view representation methods, we could also provide with three manners, including direct embedding, linear projection, and

kernel tricks, to keep their intrinsic information as much as possible. Then, we could extend such QCQP-specific single-view methods into multi-view setting to integrate the information from multiple views according to the construction process of the proposed MvLPE. As a result, we could take full advantage of these works based on single view meanwhile integrating rich information among different views.

#### IV. EXPERIMENTS

In this section, we evaluate the performance of MvLPE compared to several classical single-view and multi-view learning methods in the multi-view datasets of texts and images. We first introduce the details of the utilized datasets and comparing methods in IV-A. Then we show the experiments in IV-B and IV-C. These experiment results verify the excellent performance of MvLPE. Finally, we empirically validate the convergence of MvLPE in IV-D according to the curve of objectives.

##### A. Datasets and Competitors

Texts and images are usually represented by multi-view features, and the feature in each view is represented in high-dimensional space. Therefore, we conduct our experiments on six datasets in the form of texts and images. Three text datasets adopted in the experiments are widely used in works, including WebKB, 3Source, Cora. Three images datasets adopted in the experiments are widely used in works, including: ORL, Yale, Caltech101. We extract features for images using three different image descriptors including LBP, Gist, and EDH. To be more specific, those datasets are summarized as follows:

**WebKB** contains 4 subsets of documents over six labels and each subset consists of three views, including the text, the anchor text, and the title.

**3Sources** consists of 3 well-known online news sources: BBC, Reuters, and the Guardian, and each source is treated as one view. We select the 169 stories which are reported in all these 3 sources.

**Cora** consists of 2708 scientific publications which come from 7 classes. Because the document is represented by content and citation views, Cora could be considered as a two-views dataset.

**ORL** and **Yale** are two face image datasets that have been widely used in face recognition, where ORL consists of 400 faces corresponding to 40 peoples and Yale consists of 165 faces from 15 peoples.

**Caltech101** is a benchmark image dataset that contains 9144 images corresponding to 102 objects and it's a benchmark dataset for image classification.

More specifically, all views information of these utilized datasets is summarized in Table II. The effectiveness of MvLPE is evaluated by comparing it with the following

TABLE II: The detail information of the multi-view datasets

Datasets	Samples	Classes	Views
WebKB	226	4	Text, Anchor Text, and Title
3Sources	169	6	BBC, Reuters, and Guardian
Cora	2708	7	Content, and Cites
ORL	400	40	LBP, Gist, and EDH
Yale	165	15	LBP, Gist, and EDH
Caltech101	9144	102	LBP, Gist, and EDH

algorithms, including the best performance of the single view based LE(BLE), the feature concatenation based LE(CLE), MSE, Auto-weighted [52], Co-regularized, Co-training [53], MvCCA. Besides, we also compared the single view low-dimensional embedding in our framework with original low-dimensional embedding using MvLPE, and additional experiments on the single feature in multi-view framework by correcting and complemented by ones from other views are to verify the fact that our method could make use of complementary information among different views by correcting and complementing ones from other views.

##### B. Experiments on textual datasets

To show the superior performance of MvLPE, the experiments on three multi-view textual datasets (WebKB, 3Source, and Cora) are conducted in this section. And 1NN classifier is adopted here to classify all testing samples to verify the performances of all methods when we have obtained the low-dimensional embedding using all methods. And the mean(MEAN) and max(MAX) classification accuracy on multi-view datasets are employed as the evaluation index.

For WebKB dataset, we randomly select 50% of the samples for each subset as training samples every time. The embedding dimensionality of all the methods is set as 30. We run all methods 20 times with different random training samples and testing samples. Table III shows the MEAN and MAX value on WebKB dataset.

For 3Source dataset, we randomly select 50% of the samples as training samples and remaining samples as testing samples every time. The dimension of the embedding obtained by all methods all maintains 20 and 30 dimensions. We run all methods 20 times with different random training samples and testing samples. Table IV shows the MEAN and MAX value on 3Source dataset.

For Cora dataset, we randomly select 50% of the samples as training samples and remaining samples as testing samples every time. The dimension of the embedding obtained by all methods all maintains 20 and 30 dimensions. We run all methods 20 times with different random training samples and testing samples. Table V shows the MEAN and MAX value on Cora dataset.

##### C. Experiments on images datasets

To show the superior performance of our framework, the experiments on three multi-view images datasets (Yale, ORL, Caltech101) are conducted in this section. And 1NN classifier is adopted here to classify all testing samples to verify the performances of all methods when we have obtained the embedding using all methods.

<sup>1</sup> <http://www.webkb.org/>

<sup>2</sup> <http://mlg.ucd.ie/datasets/3sources.html>

<sup>3</sup> <http://lig-membres.imag.fr/grimal/data.html>

<sup>4</sup> <http://www.uk.research.att.com/facedatabase.html>

<sup>5</sup> <http://cvc.yale.edu/projects/yalefaces/yalefaces.html>

<sup>6</sup> <http://www.vision.caltech.edu/ImageDatasets/Caltech101/>

TABLE III: The classification accuracy on 3Source dataset

Methods	WebKB-1		WebKB-2		WebKB-3		WebKB-4	
	MEAN(%)	MAX(%)	MEAN(%)	MAX(%)	MEAN(%)	MAX(%)	MEAN(%)	MAX(%)
BLE	78.67	84.07	71.56	77.77	66.16	74.40	74.52	79.22
CLE	67.48	75.39	71.79	78.90	70.44	75.78	76.40	82.46
MSE	81.06	86.72	<b>82.64</b>	<b>87.10</b>	82.10	89.50	80.46	85.15
Auto-weighted	82.18	84.11	78.91	84.12	81.43	87.50	72.79	82.64
Co-regularized	81.05	88.64	73.12	80.95	80.30	85.71	80.12	84.41
Co-training	81.94	90.17	73.01	78.89	77.39	82.03	79.66	84.24
MvCCA	82.17	89.64	78.71	81.58	77.78	79.12	72.73	79.85
MvLPE	<b>83.17</b>	<b>91.84</b>	78.86	84.21	<b>85.28</b>	<b>92.96</b>	<b>79.96</b>	<b>85.63</b>

TABLE IV: The classification accuracy on 3Source dataset

Methods	Dims=20		Dims=30	
	MEAN(%)	MAX(%)	MEAN(%)	MAX(%)
BLE	66.47	74.11	59.72	69.41
CLE	66.50	74.71	62.78	72.94
MSE	50.47	57.64	46.86	60.00
Auto-weighted	49.92	57.64	48.15	56.47
Co-regularized	81.25	87.05	78.50	85.88
Co-training	80.80	88.23	<b>80.37</b>	90.58
MvCCA	53.88	76.45	54.37	73.56
MvLPE	<b>82.64</b>	<b>89.41</b>	79.70	<b>90.9</b>

TABLE V: The classification accuracy on Cora dataset

Methods	Dims=20		Dims=30	
	MEAN(%)	MAX(%)	MEAN(%)	MAX(%)
BLE	58.98	60.85	61.05	63.44
CLE	51.00	53.61	52.86	55.31
MSE	64.65	66.24	67.72	69.64
Auto-weighted	63.71	65.73	66.90	69.57
Co-regularized	55.73	57.45	57.19	59.01
Co-training	70.53	72.23	72.11	73.54
MvCCA	71.11	72.35	71.52	72.05
MvLPE	<b>73.7</b>	<b>75.23</b>	<b>73.45</b>	<b>75.84</b>

For Yale dataset, we extract gray-scale intensity, local binary patterns, and edge direction histogram as 3 views. The dimension of embedding obtained by all methods all maintains from 5 to 30 dimensions. we randomly select 50% samples as training ones while the other samples are assigned as the testing ones every time and run all methods 30 times with different random training samples and testing samples. Fig. 2 shows the accuracy values on Yale dataset.

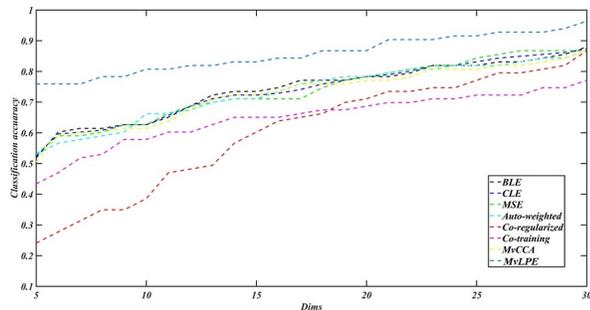


Fig. 2: The classification accuracy on Yale dataset

For ORL dataset, we also randomly select 50% samples as training ones while the other samples are assigned as the testing ones every time and run all methods 30 times with different random training samples and testing samples. And gray-scale intensity, local binary patterns, and edge direction

histogram are utilized as 3 views. The dimension of embedding obtained by all methods all maintains from 5 to 30 dimensions. Fig. 3 shows the accuracy values on ORL dataset.

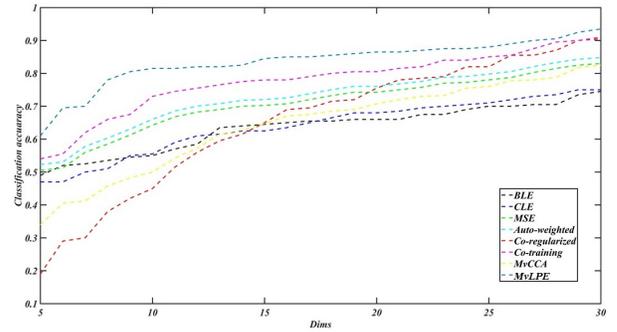


Fig. 3: The classification accuracy on ORL dataset

For Caltech101 dataset, the first 20 classes are utilized in our experiments. Meanwhile, we extract EDH, LBP, and Gist features for an image as 3 views. The dimension of embedding obtained by all methods maintains 20 and 30 dimensions. We randomly select 50% of the samples for Caltech101 dataset as training samples every time and run all methods 30 times with different random training samples and testing samples. Fig. 4 shows the mean accuracy values on Caltech101 dataset.

#### D. Convergence

Because our framework adopts an iterative procedure to obtain the optimal solution, it is essential to discuss the convergence in detail. In this section, we summarize the objective values of MvLPE on Cora and Caltech101 datasets according to the above experiments. All the training parameters (such as training numbers, dimensions) can be found above Fig.5, which summarizes the objective values of Cora and Caltech101 datasets.

## V. CONCLUSION

In this paper, we propose a novel multi-view learning method for multi-view representation, named Multi-view Low-rank Preserving Embedding (MvLPE). MvLPE deals with multi-view problems by integrating different views into one centroid view, which fully integrates compatible and complementary information from multi-view features set meanwhile maintaining low-rank reconstruction relations among samples for each view. Then, an iterative alternating strategy is adopted

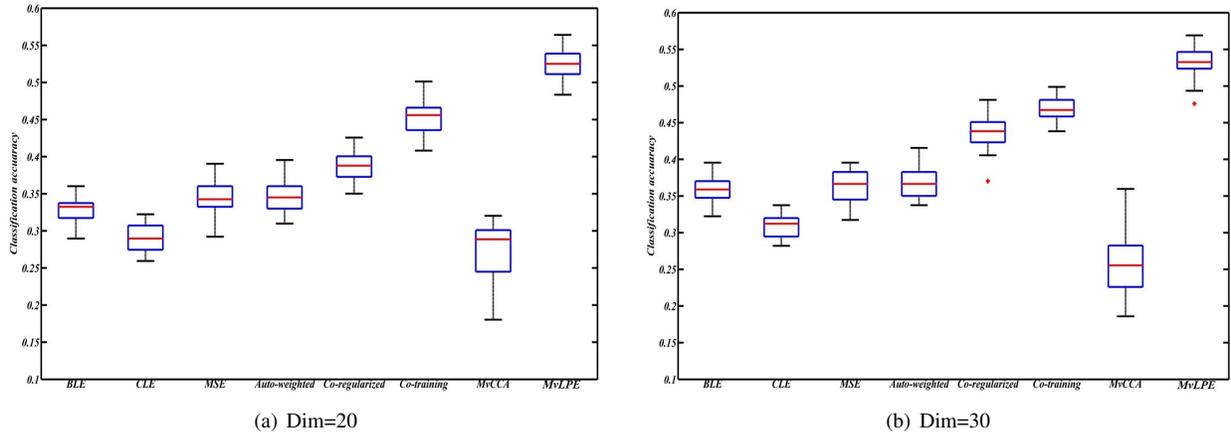


Fig. 4: Classification results on Caltech101 dataset in different dimension

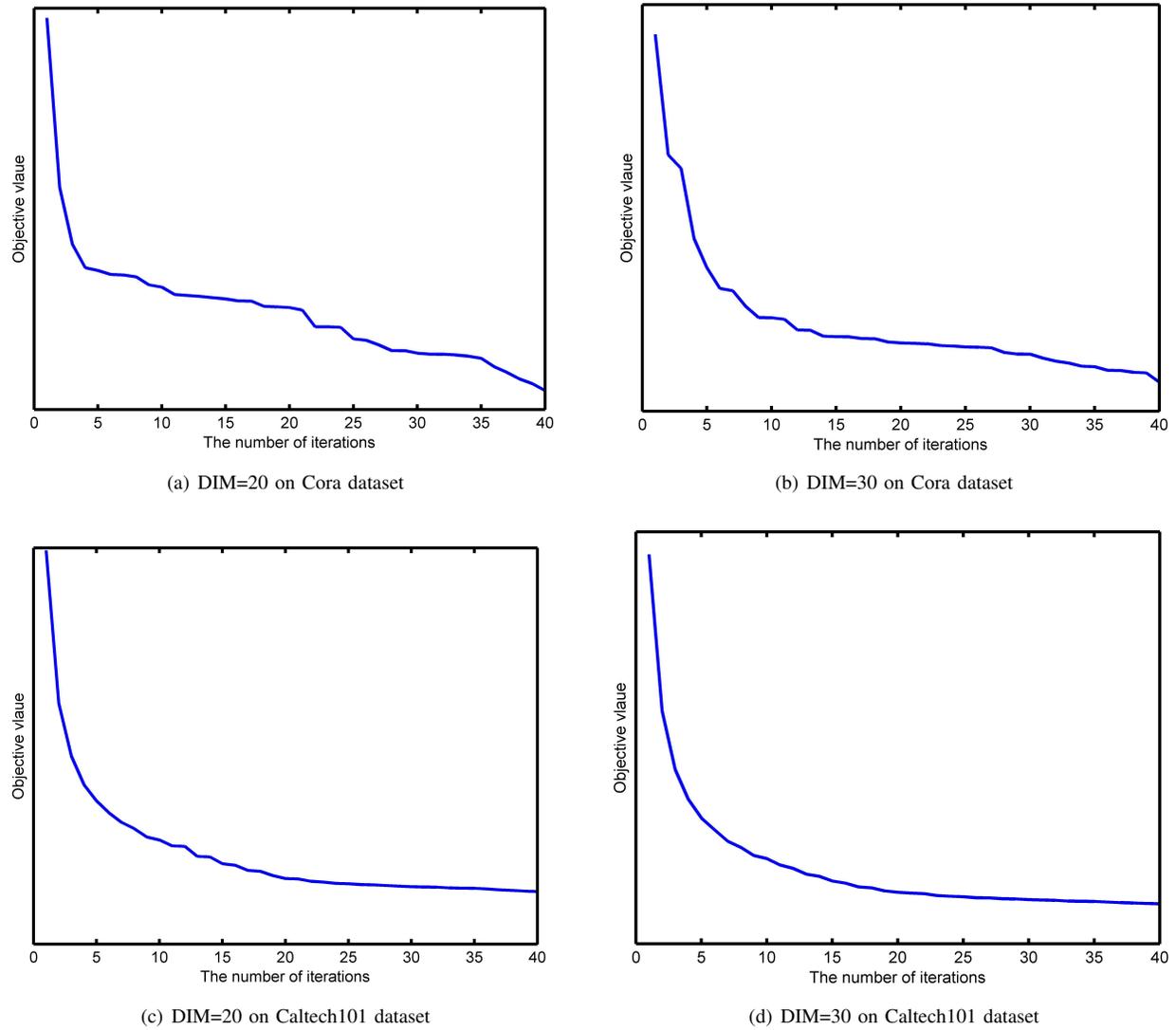


Fig. 5: Objective values of CMLLE on Cora and Caltech101

to find the optimal solution for our method and the optimization procedure is illustrated in detail. Moreover, we provide the convergence discussion of this method and its extensions for those single-view methods. The experiments have verified that the proposed MvLPE could effectively explore the underlying complementary information among multi-view data and achieve the superiority than other multi-view methods used in the experiments.

#### ACKNOWLEDGMENT

This work was supported by National Natural Science Foundation of PR China(61672130, 61972064) and LiaoNing Revitalization Talents Program(XLYC1806006).

#### REFERENCES

- [1] Y. Li, M. Yang, and Z. M. Zhang, "A survey of multi-view representation learning," *IEEE Transactions on Knowledge and Data Engineering*, 2018.
- [2] L. Feng, X. Meng, and H. Wang, "Multi-view locality low-rank embedding for dimension reduction," *Knowledge-Based Systems*, 11 2019.
- [3] C. Xu, D. Tao, and C. Xu, "A survey on multi-view learning," *arXiv preprint arXiv:1304.5634*, 2013.
- [4] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 7, pp. 971–987, 2002.
- [5] M. Douze, H. Jégou, H. Sandhwalia, L. Amsaleg, and C. Schmid, "Evaluation of gist descriptors for web-scale image search," in *Proceedings of the ACM International Conference on Image and Video Retrieval*. ACM, 2009, p. 19.
- [6] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *international Conference on computer vision and Pattern Recognition (CVPR'05)*, vol. 1. IEEE Computer Society, 2005, pp. 886–893.
- [7] X. Gao, B. Xiao, D. Tao, and X. Li, "Image categorization: Graph edit distance+ edge direction histogram," *Pattern Recognition*, vol. 41, no. 10, pp. 3179–3191, 2008.
- [8] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 12, pp. 1349–1380, 2000.
- [9] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Image retrieval: Ideas, influences, and trends of the new age," *ACM Computing Surveys (Csur)*, vol. 40, no. 2, p. 5, 2008.
- [10] J.-Y. Jiang, R.-J. Liou, and S.-J. Lee, "A fuzzy self-constructing feature clustering algorithm for text classification," *IEEE transactions on Knowledge and Data Engineering*, vol. 23, no. 3, pp. 335–349, 2011.
- [11] W. Wei and C. Guo, "A text semantic topic discovery method based on the conditional co-occurrence degree," *Neurocomputing*, vol. 368, pp. 11–24, 2019.
- [12] D. Tao and L. Jin, "Discriminative information preservation for face recognition," *Neurocomputing*, vol. 91, pp. 11–20, 2012.
- [13] L. Qiao, S. Chen, and X. Tan, "Sparsity preserving projections with applications to face recognition," *Pattern Recognition*, vol. 43, no. 1, pp. 331–341, 2010.
- [14] T. Xia, D. Tao, T. Mei, and Y. Zhang, "Multiview spectral embedding," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 40, no. 6, pp. 1438–1446, 2010.
- [15] W. Tang, Z. Lu, and I. S. Dhillon, "Clustering with multiple graphs," in *2009 Ninth IEEE International Conference on Data Mining*. IEEE, 2009, pp. 1016–1021.
- [16] G. Tzortzis and A. Likas, "Kernel-based weighted multi-view clustering," in *2012 IEEE 12th International Conference on Data Mining*. IEEE, 2012, pp. 675–684.
- [17] T. Cao, V. Jojic, S. Modla, D. Powell, K. Czymmek, and M. Niethammer, "Robust multimodal dictionary learning," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2013, pp. 259–266.
- [18] W. Liu, D. Tao, J. Cheng, and Y. Tang, "Multiview hessian discriminative sparse coding for image annotation," *Computer Vision and Image Understanding*, vol. 118, pp. 50–60, 2014.
- [19] J. Ngiam, A. Khosla, M. Kim, J. Nam, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th International Conference on Machine Learning*, 2011, pp. 689–696.
- [20] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 945–953.
- [21] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural computation*, vol. 16, no. 12, pp. 2639–2664, 2004.
- [22] J. Rupnik and J. Shawe-Taylor, "Multi-view canonical correlation analysis," in *Conference on Data Mining and Data Warehouses (SiKDD 2010)*, 2010, pp. 1–4.
- [23] M. Kan, S. Shan, H. Zhang, S. Lao, and X. Chen, "Multi-view discriminant analysis," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 1, pp. 188–194, 2016.
- [24] G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *International Conference on Machine Learning*, 2013, pp. 1247–1255.
- [25] C. Zhang, H. Fu, Q. Hu, X. Cao, Y. Xie, D. Tao, and D. Xu, "Generalized latent multi-view subspace clustering," *IEEE transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [26] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *Advances in Neural Information Processing Systems*, 2002, pp. 585–591.
- [27] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks,"

- in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [29] F. R. Bach and M. I. Jordan, “Kernel independent component analysis,” *Journal of Machine Learning Research*, vol. 3, no. Jul, pp. 1–48, 2002.
- [30] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, and K.-R. Mullers, “Fisher discriminant analysis with kernels,” in *Neural networks for signal processing IX: Proceedings of the 1999 IEEE signal processing society workshop (cat. no. 98th8468)*. Ieee, 1999, pp. 41–48.
- [31] S. Yu, Z. Cao, and X. Jiang, “Robust linear discriminant analysis with a laplacian assumption on projection distribution,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 2567–2571.
- [32] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, “Extracting and composing robust features with denoising autoencoders,” in *Proceedings of the 25th International Conference on Machine Learning*, 2008, pp. 1096–1103.
- [33] A. Kumar, P. Rai, and H. Daume, “Co-regularized multi-view spectral clustering,” in *Advances in Neural Information Processing Systems*, 2011, pp. 1413–1421.
- [34] Y. Wang, X. Lin, L. Wu, W. Zhang, Q. Zhang, and X. Huang, “Robust subspace clustering for multi-view data by exploiting correlation consensus,” *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 3939–3949, 2015.
- [35] Y. Wang, W. Zhang, L. Wu, X. Lin, M. Fang, and S. Pan, “Iterative views agreement: An iterative low-rank based structured optimization method to multi-view spectral clustering,” *arXiv preprint arXiv:1608.05560*, 2016.
- [36] C. Zhang, H. Fu, Q. Hu, P. Zhu, and X. Cao, “Flexible multi-view dimensionality co-reduction,” *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 648–659, 2016.
- [37] Y. Wang, L. Wu, X. Lin, and J. Gao, “Multiview spectral clustering via structured low-rank matrix factorization,” *IEEE transactions on Neural Networks and Learning Systems*, no. 99, pp. 1–11, 2018.
- [38] Y. Wang, W. Zhang, L. Wu, X. Lin, and X. Zhao, “Unsupervised metric fusion over multiview data by graph random walk-based cross-view diffusion,” *IEEE transactions on neural networks and learning systems*, vol. 28, no. 1, pp. 57–70, 2017.
- [39] Y. Wang, X. Lin, L. Wu, and W. Zhang, “Effective multi-query expansions: Collaborative deep networks for robust landmark retrieval,” *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1393–1404, 2017.
- [40] L. Wu, Y. Wang, and L. Shao, “Cycle-consistent deep generative hashing for cross-modal retrieval,” *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1602–1612, 2019.
- [41] Y. Wang and L. Wu, “Beyond low-rank representations: Orthogonal clustering basis reconstruction with optimized graph structure for multi-view spectral clustering,” *Neural Networks*, vol. 103, pp. 1–8, 2018.
- [42] L. Feng, X. Meng, and H. Wang, “Multi-view locality low-rank embedding for dimension reduction,” *Knowledge-Based Systems*, vol. 191, p. 105172, 2020.
- [43] X. Meng, H. Wang, and L. Feng, “The similarity-consensus regularized multi-view learning for dimension reduction,” *Knowledge-Based Systems*, p. 105835, 2020.
- [44] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, “Robust recovery of subspace structures by low-rank representation,” *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171–184, 2013.
- [45] B. Schölkopf, A. Smola, and K.-R. Müller, “Kernel principal component analysis,” in *International Conference on Artificial Neural Networks*. Springer, 1997, pp. 583–588.
- [46] F. Nie, H. Huang, and C. Ding, “Low-rank matrix recovery via efficient Schatten p-norm minimization,” in *Twenty-sixth AAAI conference on Artificial Intelligence*, 2012.
- [47] W. Rudin *et al.*, *Principles of mathematical analysis*. McGraw-hill New York, 1964, vol. 3.
- [48] S. Wold, K. Esbensen, and P. Geladi, “Principal component analysis,” *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.
- [49] X. He, D. Cai, S. Yan, and H.-J. Zhang, “Neighborhood preserving embedding,” in *Tenth IEEE International Conference on Computer Vision (ICCV’05) Volume 1*, vol. 2. IEEE, 2005, pp. 1208–1213.
- [50] X. He and P. Niyogi, “Locality preserving projections,” in *Advances in Neural Information Processing Systems*, 2004, pp. 153–160.
- [51] D. Xu, S. Yan, D. Tao, S. Lin, and H.-J. Zhang, “Marginal fisher analysis and its variants for human gait recognition and content-based image retrieval,” *IEEE Transactions on Image Processing*, vol. 16, no. 11, pp. 2811–2821, 2007.
- [52] F. Nie, G. Cai, J. Li, and X. Li, “Auto-weighted multi-view learning for image clustering and semi-supervised classification,” *IEEE Transactions on Image Processing*, vol. 27, no. 3, pp. 1501–1511, 2017.
- [53] A. Kumar and H. Daumé, “A co-training approach for multi-view spectral clustering,” in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, 2011, pp. 393–400.



**Xiangzhu Meng** received his BS degree from Anhui University, in 2015. Now he is working towards the PHD degree in School of Computer Science and Technology, Dalian University of Technology, China. His research interests include multi-view learning, deep learning and computing vision.



**Lin Feng** received the BS degree in electronic technology from Dalian University of Technology, China, in 1992, the MS degree in power engineering from Dalian University of Technology, China, in 1995, and the PhD degree in mechanical design and theory from Dalian University of Technology, China, in 2004. He is currently a professor and doctoral supervisor in the School of Innovation Experiment, Dalian University of Technology, China. His research interests include intelligent image processing, robotics, data mining, and embedded systems.



**Huibing Wang** Huibing Wang received the Ph.D. degree in the School of Computer Science and Technology, Dalian University of Technology, Dalian, in 2018. During 2016 and 2017, he is a visiting scholar at the University of Adelaide, Adelaide, Australia. Now, he is a postdoctor in Dalian Maritime University, Dalian, Liaoning, China. He has authored and co-authored more than 20 papers in some famous journals or conferences, including TMM, TITS, TSMCS, ECCV, etc. Furthermore, he serves as reviewers for TNNLS, Neurocomputing, PR Letters and MTAP, etc. His research interests include computing vision and machine learning