

# Projective Unsupervised Flexible Embedding with Optimal Graph

Wei Wang<sup>1</sup>

wei.wang@unitn.it

Yan Yan<sup>1</sup>

yan.yan@unitn.it

Feiping Nie<sup>2</sup>

feipingnie@gmail.com

Xavier Alameda Pineda<sup>1</sup>

xavier.alamedapineda@unitn.it

Shuicheng Yan<sup>3</sup>

eleyans@nus.edu.sg

Nicu Sebe<sup>1</sup>

niculae.sebe@unitn.it

<sup>1</sup> DISI,

University of Trento

Trento, Italy

<sup>2</sup> OPTIMAL,

Northwestern Polytechnical University,

Xi'an, China

<sup>3</sup> ECE,

National University of Singapore,

Singapore

## Abstract

Graph based dimensionality reduction techniques have been successfully applied to clustering and classification tasks. The fundamental basis of these algorithms is the constructed graph which dominates their performance. Usually, the graph is defined by the input affinity matrix. However, the affinity matrix is sub-optimal for dimension reduction as there is much noise in the data. To address this issue, we propose the projective unsupervised flexible embedding with optimal graph (PUFE-OG) model. We build an optimal graph by adjusting the affinity matrix. To tackle the out-of-sample problem, we employ a linear regression term to learn a projection matrix. The optimal graph and projection matrix are jointly learned by integrating the manifold regularizer and regression residual into a unified model. An efficient algorithm is derived to solve the challenging model. The experimental results on several public benchmark datasets demonstrate that the presented PUFE-OG outperforms other state-of-the-art methods.

## 1 Introduction

With the popularity of social networks, numerous data are generated everyday, such as images, videos, texts, *etc.* But usually these rich resources are not well organized and bring us big challenges to retrieve them. To address this issue, the most straightforward solution is to assign semantic labels to these data. However, the data we deal with (*e.g.*, images) are always represented by high-dimensional vectors. To better organize and represent the data, many feature selection [1, 2] and probabilistic generative models with feature selection/missing data [3, 4], dimensionality reduction techniques, and clustering algorithms were proposed and benefited many related applications, such as image annotation [5], multi-view

clustering [29], human action recognition [32], human face aging [25], human headpose estimation [33], attribute detection [24, 26] and video annotation [14].

Many dimensionality reduction algorithms seek for an optimal projection matrix to map the data. These algorithms are always linked to clustering and classification algorithms and they work as an intermediate step. For instance, many works perform K-Means clustering after doing Principle Component Analysis (PCA), and spectral clustering (SC) uses Laplacian eigenmap to project the data which is followed by a K-Means clustering. Recently, the denoising auto-encoder [23] has been employed for dimensionality reduction before learning the classifiers. Auto-encoder is very similar to PCA when it only has one linear hidden layer. One main reason to do dimensionality reduction is that the low-dimensional embeddings can better represent the data and accelerate the clustering or classification process.

When the labels are available, the most popular dimensionality reduction algorithm is linear discriminant analysis (LDA). It has excellent performance as LDA utilizes discriminant information to learn the subspace. In addition, simultaneously performing clustering and subspace learning can yield even better clustering result. *Ye et al.* [35] proposed an effective discriminating K-Means (DisKmeans) algorithm by integrating LDA and K-Means. *Yang et al.* [34] proposed a local discriminant model and global integration (LDMGI) model by integrating manifold subspace learning and clustering into a unified framework. However, labeled data are often very costly to obtain.

When the labels are unavailable, unsupervised dimensionality reduction methods become the only choice. For example, PCA is widely used because of its simplicity and efficiency. The unsupervised graph based dimensionality reduction methods (e.g., Neighbor Preserving Embedding [5], Locality Preserving Projections [17]) usually outperform PCA. This is because these methods take advantage of manifold information. Many graph based dimensionality reduction methods have been explored, such as locally linear embedding (LLE) [19], Laplacian eigenmap (LE) [11], and ISOMAP [22]. However, these methods suffer from out-of-sample problem. They can not map the new data points that are not included in the training set. To tackle this problem, many works [15] integrated the manifold regularizer with the ridge regression loss into the subspace learning framework. Similar to other manifold learning algorithms, their performance is also controlled by the graph constructed by the fixed affinity matrix, which might lead to a sub-optimal result [16]. To address this issue, we propose a projective unsupervised flexible embedding with optimal graph (PUFE-OG) framework. Instead of utilizing the fixed affinity matrix to preserve the manifold structure, we construct an optimal graph by adapting the affinity matrix for subspace learning. It is worthwhile to highlight the following contributions of our work: (1) The proposed framework exploits an unsupervised flexible embedding for high-dimensional data. The framework automatically adjusts the graph to an optimal graph based on the local manifold structure; (2) The optimal graph seamlessly integrates manifold learning with subspace learning into a unified framework. Besides, our framework has closed form solution when update its nested variables iteratively.

## 2 Related Work

**Dimensionality Reduction Techniques:** To better organize and represent data, various supervised and unsupervised dimensionality reduction methods are proposed. As a popular unsupervised dimensionality reduction method, PCA considers variance as the most important metric and aims to minimize the data reconstruction error by minimizing the variance.

However, PCA can not guarantee that the learned subspace has discriminant power and it may not improve the clustering result. Different from PCA, spectral clustering [16] is a graph based clustering method. SC utilizes LE for dimensionality reduction based on the affinity matrix. SC is more robust compared with PCA as LE preserves the local manifold structure of the data points. Thus, the data are tailored to low-dimensional space and it accelerates K-Means clustering. To further improve the performance of SC, many works focused on learning a better graph. *Liu et al.* [17] proposed to compress the original large scale graph into a sparse bipartite graph. *Shao et al.* [18] adopted deep networks in SC. But SC can not cluster the data which are not included in the training set as LE can not map new data to the low-dimensional space. Thus, SC suffers from the out-of-sample problem.

In the supervised settings, Fisher LDA is widely used for subspace learning. LDA [19] minimizes the within-cluster scatter matrix while maximizes the between-cluster matrix. Then the cluster centers will be as far as possible while the within-cluster data will be as close as possible. Thus, the learned subspace has greater discriminative power compared with PCA. Dictionary learning can also be employed as a clustering method. *Ramirez et al.* [20] proposed training a dictionary for each cluster. A cross-incoherence promoting term was integrated into the dictionary learning framework. This term encourages the dictionaries to be as independent as possible, which can lead to greater discriminative power. But the performance of these methods highly relies on the sufficiency of the labeled training data which is very time-consuming to obtain.

**Manifold Learning:** When labels are unavailable, graph based dimensionality reduction methods can usually provide a better low-dimensional representation and outperform PCA by utilizing the local manifold structure. Given a dataset  $\mathbf{X}=[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ , where  $n$  is the number of samples and  $\mathbf{x}_i \in \mathbb{R}^d, \forall i$ , manifold information is preserved in the undirected graph  $\mathbf{G}=\{\mathbf{X}, \mathbf{A}\}$  with data matrix  $\mathbf{X}$  and affinity matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . The  $(i, j)^{th}$  element  $A_{i,j}$  in  $\mathbf{A}$  denotes the connectivity between sample  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .  $A_{i,j}$  can be viewed as the weight of the edge in the graph. Let us denote the low-dimensional representations of  $\mathbf{x}_i$  as  $\mathbf{f}_i \in \mathbb{R}^{1 \times c}$ , where  $c$  is the lower dimension and  $\mathbf{f}_i$  is a row vector. The goal of manifold learning is to construct another graph  $\hat{\mathbf{G}}=\{\mathbf{F}, \hat{\mathbf{A}}\}$ , where  $\mathbf{F} \in \mathbb{R}^{n \times c}$  is the low-dimensional representation of  $\mathbf{X}$ . Ideally, we should have  $\hat{\mathbf{A}}=\mathbf{A}$ , which means the connectivity between sample  $\mathbf{f}_i$  and  $\mathbf{f}_j$  should be the same as the connectivity between sample  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Gaussian kernel is widely used to calculate the affinity matrix.  $\mathbf{A}$  is defined as follows:

$$A_{i,j} = \begin{cases} \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}), & \text{if } \mathbf{x}_i, \mathbf{x}_j \text{ are } k \text{ nearest neighbors.} \\ 0, & \text{otherwise.} \end{cases}$$

$\sigma$  controls the weight of the edge. Most manifold learning methods (e.g., LE [21]) represent the manifold structure with this graph, and require the low-dimensional representations to preserve the same manifold structure by employing a manifold regularizer as follows:

$$\min_{\mathbf{F}, \mathbf{F}^T \mathbf{D} \mathbf{F} = \mathbf{I}} \mathbf{F}^T \mathbf{L}_A \mathbf{F} \quad (1)$$

Eq. (1) is the objective function of LE,  $\mathbf{D}$  is a diagonal matrix with the element  $D_{i,i} = \sum_j A_{i,j}$ . The graph Laplacian matrix  $\mathbf{L}_A = \mathbf{D} - \mathbf{A}$ .

To better understand the intuition why this term can preserve the manifold structure, we reformulate Eq. (1) as Eq. (2) ( $\mathbf{A}$  is symmetric).

$$\text{Tr}(\mathbf{F}^T \mathbf{L}_A \mathbf{F}) = \frac{1}{2} \sum_{i,j=1}^n \|\mathbf{f}_i - \mathbf{f}_j\|_2^2 A_{i,j} \quad (2)$$

where  $\mathbf{A}$  can be viewed as a similarity matrix which measures the similarity between every two data points. A larger  $A_{i,j}$  implies a more similar pair of samples. The closer the two samples are, the larger  $A_{i,j}$  will become. Correspondingly,  $\|\mathbf{f}_i - \mathbf{f}_j\|_2^2$  will become smaller to keep the balance.

To deal with the out-of-sample problem and preserve the manifold smoothness, [21] proposed a LapRLS/L method, which exploited ridge regression to learn the subspace and a manifold regularizer term to make use of the local geometry information. Let  $\mathbf{W} \in \mathbb{R}^{d \times c}$  be the projection matrix and  $\mathbf{b} \in \mathbb{R}^c$  be the bias term. The ridge regression function is defined as

$$\frac{1}{n} \|\mathbf{X}^T \mathbf{W} + \mathbf{1} \mathbf{b}^T - \mathbf{Y}\|_F^2 + \gamma \|\mathbf{W}\|_F^2 \quad (3)$$

where  $\mathbf{1} \in \mathbb{R}^n$ , and with all its entries being 1;  $\mathbf{Y}$  is a binary class assignment matrix.  $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n] \in \mathbb{R}^{n \times c}$ , where  $\mathbf{y}_i \in \{0, 1\}^{c \times 1}$  is the class indicator vector for  $\mathbf{x}_i$ . The  $j$ -th element of  $\mathbf{y}_i$  is 1 if  $\mathbf{x}_i$  belongs to the  $j$ -th class, and 0 otherwise. The second term  $\gamma \|\mathbf{W}\|_F^2$  is a regularizer to prevent the entries in the projection matrix  $\mathbf{W}$  becoming too large. The manifold regularizer is

$$\lambda_l \text{Tr}(\mathbf{W}^T \mathbf{X} \mathbf{L}_A \mathbf{X}^T \mathbf{W}) \quad (4)$$

where  $\mathbf{X}^T \mathbf{W}$  is similar to  $\mathbf{F}$  in Eq. (1). After combining Eq. (3) and Eq. (4), the projection matrix and manifold structure can be jointly learned. Then the low-dimensional representation of new data can be obtained using the linear projection function  $h(\mathbf{x}) = \mathbf{x}^T \mathbf{W} + \mathbf{b}^T$ . However, this model still requires labeled training data. Nie *et al.* [15] further extended the LapRLS/L to deal with the unlabeled data. The performance of these graph based methods which integrate manifold regularizer is dominated by the graph.

Thus, learning an effective graph is very essential. However, the graph is always *fixed* by the affinity matrix  $\mathbf{A}$ . The affinity matrix is a hard constraint, and the learned low-dimensional representations are restricted to have the same neighbourhood relationship ( $\hat{\mathbf{A}} = \mathbf{A}$ ) without any violation. But the neighbourhood relationship may be inaccurate as there is much noise in the data. To address this issue, we propose constructing an optimal graph by adjusting the neighbourhood relationships for projective unsupervised *flexible* manifold embedding which can learn the projection matrix and the optimal graph simultaneously.

### 3 Projective Unsupervised Flexible Embedding Framework (PUFE)

To take advantage of manifold information and solve the out-of-sample problem, a unified PUFE framework is proposed by integrating the manifold regularizer with the ridge regression terms. The PUFE is as follows:

$$\min_{\substack{\mathbf{F} \in \mathbb{R}^{n \times c}, \mathbf{F}^T \mathbf{F} = \mathbf{I} \\ \mathbf{W} \in \mathbb{R}^{d \times c}, \mathbf{b} \in \mathbb{R}^c}} \text{Tr}(\mathbf{F}^T \mathbf{L}_A \mathbf{F}) + \mu \|\mathbf{X}^T \mathbf{W} + \mathbf{1} \mathbf{b}^T - \mathbf{F}\|_F^2 + \gamma \|\mathbf{W}\|_F^2 \quad (5)$$

The functionality of each term is described as follows:

- The first term,  $\text{Tr}(\mathbf{F}^T \mathbf{L}_A \mathbf{F})$ , is a manifold regularizer as in Eq. (1), which is in charge of maintaining the manifold smoothness of the data.
- The second term denotes the ridge regression error. This linear regression model learns

a projection matrix  $\mathbf{W}$  to project samples to a low-dimensional space. Usually, the dimensionality reduction methods, such as PCA, LDA, use a linear function to project data matrix  $\mathbf{X}$ :  $\mathbf{F}=\mathbf{X}^T\mathbf{W}$ , and  $\mathbf{F}$  lies in the space spanned by  $\mathbf{X}$ . But the second term relaxes the rigid form  $\mathbf{F}=\mathbf{X}^T\mathbf{W}$  to  $\mathbf{F}\approx\mathbf{X}^T\mathbf{W}+\mathbf{1b}^T$  by minimizing  $\|\mathbf{X}^T\mathbf{W}+\mathbf{1b}^T-\mathbf{F}\|_F^2$  which is more flexible. Moreover, this term allows its interaction with manifold regularizer through the variable  $\mathbf{F}$ .

• The third term is a regularizer which is always combined with the ridge regression term. Its function is to prevent overfitting. The weights,  $\mu$ ,  $\gamma$  balance the importance of each term. **Optimal Graph Construction:** As a graph based method, PUFEE tackles the out-of-sample problem, and the affinity matrix  $\mathbf{A}$  has a significant influence on its performance. The affinity matrix  $\mathbf{A}$  can be viewed as a probability matrix. If  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are similar to each other, a large value will be assigned to  $A_{i,j}$  and it represents that  $\mathbf{x}_i$  has a higher possibility to be the neighbor of  $\mathbf{x}_j$ .

In Eq. (5), graph  $\mathbf{A}$  is a hard constraint, and  $\mathbf{F}$  is forced to embed into the fixed graph  $\mathbf{A}$  by the penalty term  $Tr(\mathbf{F}^T\mathbf{L}_\mathbf{A}\mathbf{F})$  in the loss function. However,  $\mathbf{A}$  is not an optimal graph for dimension reduction, and it may result in a sub-optimal low-dimensional representation. In this paper, we aim to construct an optimal graph  $\mathbf{S}$  ( $S_{i,j} \geq 0$ ) based on the affinity matrix  $\mathbf{A}$ . Similar to  $\mathbf{A}$ , the entry  $S_{i,j}$  in  $\mathbf{S}$  represents the probability that  $\mathbf{x}_j$  being the neighbor of  $\mathbf{x}_i$ . Thus, we have  $\sum_j S_{i,j}=1$ , which can be formulated as  $\mathbf{S}\mathbf{1}=\mathbf{1}$ .  $\mathbf{S}$  works as a soft constraint for  $\mathbf{F}$ . The loss  $\|\mathbf{S}-\mathbf{A}\|_F^2$  can prevent  $\mathbf{S}$  from deviating too far from  $\mathbf{A}$ . At the same time,  $\mathbf{S}$  allows the neighborhood probability to be adjusted for dimension reduction. Then  $\mathbf{F}$  is embedded into the optimal graph  $\mathbf{S}$ .

Let  $\mathbf{L}_\mathbf{S}$  denote the graph Laplacian matrix, we propose the following objective function and the optimal graph for dimension reduction can be derived by solving the following function:

$$\min_{\substack{\mathbf{F} \in \mathbb{R}^{n \times c}, \mathbf{F}^T \mathbf{F} = \mathbf{I} \\ \mathbf{S} \in \mathbb{R}^{n \times n}, \mathbf{S}\mathbf{1} = \mathbf{1}, \mathbf{S} \geq \mathbf{0}}} \|\mathbf{S}-\mathbf{A}\|_F^2 + \lambda Tr(\mathbf{F}^T \mathbf{L}_\mathbf{S} \mathbf{F}) \quad (6)$$

The graph Laplacian matrix  $\mathbf{L}_\mathbf{S} \in \mathbb{R}^{n \times n}$  is denoted as  $\mathbf{L}_\mathbf{S} = \mathbf{D} - \frac{\mathbf{S} + \mathbf{S}^T}{2}$ , where  $\mathbf{D}$  is a diagonal matrix with diagonal elements  $D_{i,i} = \sum_j (S_{i,j} + S_{j,i})/2, \forall i$ . Till now, we have constructed the optimal graph based on the similarity matrix. One very important advantage of this model is that the neighborhood probability can be tuned.

**Projective Unsupervised Flexible Embedding with Optimal Graph (PUFE-OG):** By replacing the affinity matrix  $\mathbf{A}$  with an optimal neighborhood probability matrix  $\mathbf{S}$ . We can learn the optimal graph and projection matrix jointly. The following model is obtained after combining the optimal graph and unsupervised flexing embedding:

$$\begin{aligned} \min_{\mathbf{F}, \mathbf{W}, \mathbf{b}, \mathbf{S}} \quad & \left( \|\mathbf{S}-\mathbf{A}\|_F^2 + \lambda Tr(\mathbf{F}^T \mathbf{L}_\mathbf{S} \mathbf{F}) \right. \\ & \left. + \mu \|\mathbf{X}^T \mathbf{W} + \mathbf{1b}^T - \mathbf{F}\|_F^2 + \gamma \|\mathbf{W}\|_F^2 \right) \\ \text{subject to} \quad & \mathbf{F}^T \mathbf{F} = \mathbf{I}, \mathbf{F} \in \mathbb{R}^{n \times c} \\ & \mathbf{W} \in \mathbb{R}^{d \times c}, \mathbf{b} \in \mathbb{R}^c \\ & \mathbf{S}\mathbf{1} = \mathbf{1}, \mathbf{S} \geq \mathbf{0}, \mathbf{S} \in \mathbb{R}^{n \times n} \end{aligned} \quad (7)$$

The functionality of each term is described as follows:

• The first term,  $\|\mathbf{S}-\mathbf{A}\|_F^2$ , is aimed at constructing an optimal graph  $\mathbf{S} \in \mathbb{R}^{n \times n}$ . The Frobenius norm measures the Euclidean distance between  $\mathbf{S}$  and  $\mathbf{A}$ . This term allows the graph  $\mathbf{S}$  to

adjust itself for the optimal graph construction while prevents it deviating too far away from the similarity matrix  $\mathbf{A}$ .

- The second term,  $\lambda \text{Tr}(\mathbf{F}^T \mathbf{L}_S \mathbf{F})$ , is in charge of maintaining the manifold smoothness of the data. Different from Eq. (5) whose Laplacian matrix  $\mathbf{L}_A$  is derived from the fixed affinity matrix  $\mathbf{A}$ , the Laplacian matrix  $\mathbf{L}_S$  in this term is based on a soft graph  $\mathbf{S}$ .
- The third term denotes the ridge regression residual. The forth term is a regularizer and its function is to prevent overfitting. The weights,  $\lambda$ ,  $\mu$ ,  $\gamma$  represent the importance of each term.

**Optimization of PUFEE-OG:** We solve the model in Eq. (7) by optimizing the variables alternatively which mainly consists of two parts.

**Initialization:** After setting  $\mathbf{X}=[\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n]$  and other parameters,  $\mathbf{A}$  is set to affinity matrix  $\mathbf{A}_0$ .  $\mathbf{S}_0$  is initialized according to Eq. (14) with  $\mathbf{d}_i$  set to  $\mathbf{0}$ .

**(1) Fix  $\mathbf{S}$ , update  $\mathbf{W}, \mathbf{b}, \mathbf{F}$ :** When  $\mathbf{S}$  is fixed, the original optimization model in Eq. (7) can be put in the following form,

$$\min_{\substack{\mathbf{F}^T \mathbf{F} = \mathbf{I} \\ \mathbf{W}, \mathbf{b}, \mathbf{F} \in \mathbb{R}^{n \times c}}} \left( \lambda \text{Tr}(\mathbf{F}^T \mathbf{L}_S \mathbf{F}) + \mu \|\mathbf{X}^T \mathbf{W} + \mathbf{1} \mathbf{b}^T - \mathbf{F}\|_F^2 + \gamma \|\mathbf{W}\|_F^2 \right) \quad (8)$$

Eq. (8) is convex with respect to  $\mathbf{W}, \mathbf{b}, \mathbf{F}$ . The proof is available in [15]. By setting the derivative with respect to  $\mathbf{W}$  and  $\mathbf{b}$  to  $\mathbf{0}$ , we have

$$\begin{aligned} \mathbf{b} &= \frac{1}{n} (\mathbf{F}^T \mathbf{1} - \mathbf{W}^T \mathbf{X} \mathbf{1}) \\ \mathbf{W} &= (\mathbf{X} \mathbf{H}_n \mathbf{X}^T + \frac{\gamma}{\mu} \mathbf{I})^{-1} \mathbf{X} \mathbf{H}_n \mathbf{F} = \mathbf{Q} \mathbf{F} \end{aligned} \quad (9)$$

where  $\mathbf{Q} = (\mathbf{X} \mathbf{H}_n \mathbf{X}^T + \frac{\gamma}{\mu} \mathbf{I})^{-1} \mathbf{X} \mathbf{H}_n$ .

Let  $\mathbf{X}_n = \mathbf{X} \mathbf{H}_n$ . The solution for  $\mathbf{F}$  is

$$\mathbf{F}^* = \arg \min_{\mathbf{F}, \mathbf{F}^T \mathbf{F} = \mathbf{I}} \text{Tr} \mathbf{F}^T (\lambda \mathbf{L}_S + \mu \mathbf{H}_n - \mu \mathbf{N}) \mathbf{F} \quad (10)$$

where

$$\mathbf{N} = \mathbf{X}_n^T (\mathbf{X}_n \mathbf{X}_n^T + \frac{\gamma}{\mu} \mathbf{I})^{-1} \mathbf{X}_n = \mathbf{X}_n^T \mathbf{X}_n (\mathbf{X}_n^T \mathbf{X}_n + \frac{\gamma}{\mu} \mathbf{I})^{-1}$$

A generalized eigenvalue decomposition [16] can be utilized to solve this objective function.

**(2) Fix  $\mathbf{W}, \mathbf{b}, \mathbf{F}$ , update  $\mathbf{S}$ :** When  $\mathbf{W}, \mathbf{b}, \mathbf{F}$  are fixed, the original optimization model Eq. (7) can be put in the following form,

$$\min_{\substack{\mathbf{S} \in \mathbb{R}^{n \times n}, \mathbf{S}^* \mathbf{1} = \mathbf{1}, \mathbf{S} \geq \mathbf{0}}} \left( \|\mathbf{S} - \mathbf{A}\|_F^2 + \lambda \text{Tr}(\mathbf{F}^T \mathbf{L}_S \mathbf{F}) \right) \quad (11)$$

which is equivalent to the following

$$\min_{s_{i,j} \geq 0, \sum_j s_{i,j} = 1} \sum_{i,j=1}^n (s_{i,j} - A_{i,j})^2 + \frac{\lambda}{2} \sum_{i,j=1}^n \|\mathbf{f}_i - \mathbf{f}_j\|_2^2 s_{i,j} \quad (12)$$

The rows of  $\mathbf{S}$  are independent from each other. Thus, we can solve  $\mathbf{S}$  row by row. For each row  $\mathbf{s}_i$ , we have

$$\min_{\mathbf{s}_i \geq 0, \mathbf{s}_i^* \mathbf{1} = 1} \|\mathbf{s}_i - \mathbf{a}_i\|_2^2 + \frac{\lambda}{2} \mathbf{s}_i \mathbf{d}_i^T \quad (13)$$

where  $\mathbf{s}_i \in \mathbb{R}^{1 \times n}$  whose  $j_{th}$  element is  $S_{i,j}$ ,  $\mathbf{a}_i \in \mathbb{R}^{1 \times n}$  whose  $j_{th}$  element is  $A_{i,j}$ , and  $j_{th}$  element of  $\mathbf{d}_i$  is  $\|\mathbf{f}_i - \mathbf{f}_j\|_2^2$ . Eq. (13) can be further formulated as

$$\begin{aligned} & \min_{\mathbf{s}_i \geq 0, \sum_{j=1}^n \mathbf{s}_i = 1} \left\| \mathbf{s}_i - \left( \mathbf{a}_i - \frac{\lambda}{4} \mathbf{d}_i \right) \right\|_2^2 \\ & = \min_{x_j \geq 0, \sum_j x_j = 1} \left\| \mathbf{x} - \mathbf{v} \right\|_2^2 \end{aligned} \quad (14)$$

where  $\mathbf{x} = \mathbf{s}_i$ ,  $\mathbf{v} = \mathbf{a}_i - \frac{\lambda}{4} \mathbf{d}_i$ .  $\|\mathbf{x} - \mathbf{v}\|_2^2 = R^2$  is a Euclidean ball, and  $x_j \geq 0, \sum_j x_j = 1$  defines a hyper-plane. The optimal solution is the tangent point between the Euclidean ball and the hyper-plane or the angle point of the hyper-plane. We summarize the structure of the solution in Algorithm 1.

---

**Algorithm 1** Optimization Algorithm of PUF-OG

---

- 1: *Input: cluster number  $c$ , parameter  $\lambda, \mu, \sigma, \mathbf{X}$ .*
  - 2: *Output: learned projection parameters  $\mathbf{W}, \mathbf{b}$ .*
  - 3: *Initialization:  $\mathbf{A} \leftarrow \mathbf{A}_0, \mathbf{S} \leftarrow \mathbf{S}_0, \mathbf{L}_S = \text{Laplace}(\mathbf{S})$ .*
  - 4: **repeat**
  - 5:   *fix  $\mathbf{S}$ , update  $\mathbf{W}, \mathbf{b}, \mathbf{F}$ :*
  - 6:    *Update  $\mathbf{F}$  according to Eq. (10);*
  - 7:    *Update  $\mathbf{W}, \mathbf{b}$  according to Eq. (9).*
  - 8:   *fix  $\mathbf{W}, \mathbf{b}, \mathbf{F}$ , update  $\mathbf{S}$ :*
  - 9:    *For each row  $\mathbf{s}_i$ , update it according to Eq. (14).*
  - 10:     $\mathbf{L}_S = \mathbf{D}_S - \frac{\mathbf{S} + \mathbf{S}^T}{2}$ ;
  - 11: **until** converges
- 

## 4 Experiments

**Datasets & Setup:** To evaluate our model, we choose three multiview action recognition datasets, three face datasets and one handwritten digit recognition dataset. The action datasets are the IXMAS dataset [47], the newer version of IXMAS dataset referred to as NIXMAS, and the partially occluded dataset OIXMAS [28] dataset. The three face datasets are the JAFFE dataset [43], the UMIST face dataset [9], and the YaleB dataset [9]. We also use the USPS dataset [9] to validate the performance on handwritten digit recognition. Table 1 shows the details of the 7 benchmark datasets, such as the number of instances, the dimension of each instance, as well as the number of classes.

Datasets	Instance No.	Dimension	Class No.
IXMAS [47]	1789	500	12
NIXMAS [47]	1148	500	11
OIXMAS [28]	1800	500	12
JAFFE [43]	213	676	10
UMIST [9]	575	644	20
USPS [9]	9298	256	10
YaleB [9]	575	644	20

Table 1: Description of 7 Benchmark Datasets

In the experiment, we found that 15 is a good neighborhood number to build the graph. For parameters  $\lambda, \mu, \sigma$ , we tune them in the range of  $[10^{-3}; 10^{-2}; \dots 10^3]$ . The dimension of the low-dimensional space is set to  $[10; 20; \dots; 100]$ .

Table 2: Performance comparison (Mean Accuracy%  $\pm$  Standard Deviation)

	IXMAS	NIXMAS	OIXMAS	JAFFE	UMIST	USPS	Yaleb
Raw Feature	20.99 $\pm$ 2.40	20.64 $\pm$ 0.86	18.47 $\pm$ 0.95	73.29 $\pm$ 7.27	41.56 $\pm$ 2.47	65.08 $\pm$ 2.34	11.73 $\pm$ 0.80
PCA	21.11 $\pm$ 2.42	20.50 $\pm$ 0.80	18.63 $\pm$ 0.96	73.18 $\pm$ 7.20	43.29 $\pm$ 2.52	66.31 $\pm$ 2.36	11.71 $\pm$ 0.71
Atuo Enc. [24]	25.59 $\pm$ 2.33	22.30 $\pm$ 0.14	16.39 $\pm$ 1.09	67.81 $\pm$ 5.17	37.58 $\pm$ 2.42	46.90 $\pm$ 0.79	7.33 $\pm$ 0.68
NPE [9]	19.44 $\pm$ 2.32	19.10 $\pm$ 0.68	18.78 $\pm$ 1.05	75.25 $\pm$ 7.92	40.93 $\pm$ 2.44	62.74 $\pm$ 2.71	20.70 $\pm$ 1.18
LPP [17]	22.57 $\pm$ 2.03	19.84 $\pm$ 0.80	18.53 $\pm$ 0.73	79.28 $\pm$ 7.81	44.50 $\pm$ 2.74	68.62 $\pm$ 3.48	28.89 $\pm$ 1.51
PUFE [16]	28.28 $\pm$ 1.71	21.46 $\pm$ 0.54	20.94 $\pm$ 0.58	80.98 $\pm$ 8.39	55.01 $\pm$ 2.65	71.17 $\pm$ 3.64	45.07 $\pm$ 2.47
PUFE-OG	<b>28.59<math>\pm</math>1.72</b>	<b>24.02<math>\pm</math>0.52</b>	<b>21.50<math>\pm</math>0.77</b>	<b>83.30<math>\pm</math>7.24</b>	<b>71.13<math>\pm</math>1.18</b>	<b>75.22<math>\pm</math>2.65</b>	<b>60.05<math>\pm</math>2.43</b>

Table 3: Performance comparison (Mean NMI%  $\pm$  Standard Deviation)

	IXMAS	NIXMAS	OIXMAS	JAFFE	UMIST	USPS	Yaleb
Raw Feature	27.14 $\pm$ 2.73	15.91 $\pm$ 0.97	16.78 $\pm$ 2.11	80.08 $\pm$ 5.06	64.10 $\pm$ 1.76	66.97 $\pm$ 0.86	16.48 $\pm$ 1.09
PCA	27.11 $\pm$ 2.81	15.79 $\pm$ 0.77	16.95 $\pm$ 2.04	80.18 $\pm$ 4.92	63.86 $\pm$ 1.71	61.03 $\pm$ 0.81	16.54 $\pm$ 0.98
Atuo Enc. [24]	25.59 $\pm$ 1.99	20.34 $\pm$ 0.04	13.89 $\pm$ 1.03	75.26 $\pm$ 2.57	54.04 $\pm$ 1.96	42.49 $\pm$ 0.14	9.64 $\pm$ 0.67
NPE [9]	25.09 $\pm$ 2.62	12.60 $\pm$ 0.76	19.03 $\pm$ 2.09	82.53 $\pm$ 4.43	61.01 $\pm$ 2.01	59.42 $\pm$ 1.36	27.06 $\pm$ 0.97
LPP [17]	28.19 $\pm$ 2.62	12.97 $\pm$ 0.96	14.38 $\pm$ 0.11	86.47 $\pm$ 4.33	62.79 $\pm$ 2.25	65.37 $\pm$ 0.99	41.10 $\pm$ 1.00
PUFE [16]	31.75 $\pm$ 2.64	16.90 $\pm$ 0.86	23.53 $\pm$ 0.58	86.44 $\pm$ 5.05	74.51 $\pm$ 1.46	70.82 $\pm$ 1.64	56.85 $\pm$ 1.59
PUFE-OG	<b>33.68<math>\pm</math>2.15</b>	<b>28.83<math>\pm</math>0.65</b>	<b>24.95<math>\pm</math>1.23</b>	<b>88.47<math>\pm</math>4.54</b>	<b>86.71<math>\pm</math>2.26</b>	<b>73.71<math>\pm</math>0.85</b>	<b>73.14<math>\pm</math>1.61</b>

**Experimental Results:** To evaluate the performance with its projection ability, we compare it with raw feature, PCA, denoising auto-encoder (pretraining stage) [24], Neighbor Preserving Embedding (NPE) [9], Locality Preserving Projections (LPP) [17], as well as the PUFE [16] without optimal graph. To evaluate the projection ability of each method, we first learn the projection matrix and obtain the corresponding low-dimensional representations. Then we rely on K-Means to do the clustering. The result is shown in Table 2 and Table 3. We run the K-Means 100 times each time, and report the average value and standard deviation in the table.

The experimental results are listed in Table 2 and Table 3. We can observe that the performance of our method is consistently better than the others under different evaluation metrics. We also observe that:

- The graph based dimensionality reduction methods (NPE, LPP, PUFE, PUFE-OG) usually achieve better performance compared with PCA and raw feature. This observation suggests that low-dimensional representations can have better performance if it is based on a weighted graph which contains the structure information of the local cliques.
- The performance of PCA is always similar to the performance of raw feature. Sometimes it has little improvement or inferior performance compared with the raw feature. The underlying reason of the observation is that PCA does not preserve the manifold smoothness.
- Denoising auto-encoder have relatively better performance for the first two datasets than the raw feature. This observations suggests that denoising auto-encoder is robust to view changes compared with other methods. However, it fails capturing the inter-class difference of other datasets and has the worst performance among all the methods.
- PUFE always outperforms LPP and NPE and it is very stable across all the datasets. This is because PUFE relaxes the rigid projection form, namely,  $\mathbf{F}=\mathbf{X}^T\mathbf{W}$ , to a more flexible version, which is  $\mathbf{F}\approx\mathbf{X}^T\mathbf{W}+\mathbf{1b}^T$ .
- Our PUFE-OG consistently outperforms other graph based algorithms and it is very stable. This advantage is attributed to the optimal neighborhood probability graph construction.
- The data in UMIST is more regular compared with other datasets. Thus, its manifold structure is more label consistent (the samples from the same cluster are lying close by). Then the improvement in accuracy is higher compared with other datasets.

**Parameter Sensitivity Analysis:** We choose UMIST dataset to study the sensitivity of the



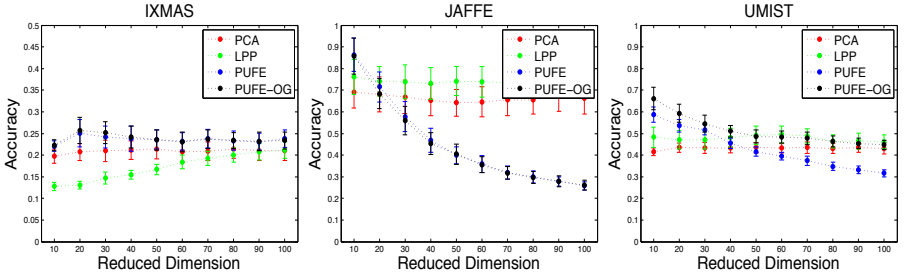


Figure 1: Results of Projection with Different Dimensions on 3 Benchmark Datasets

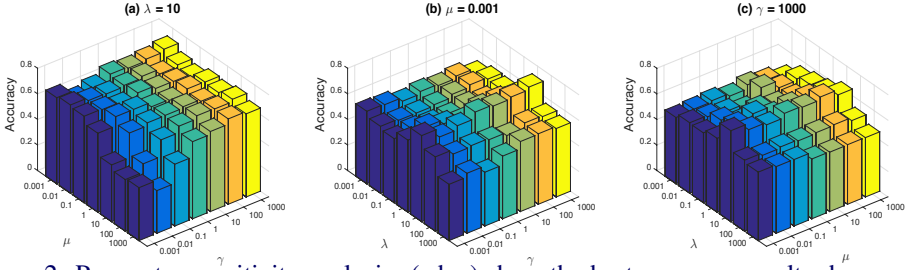


Figure 2: Parameter sensitivity analysis: (a,b,c) show the best accuracy result when one of the three parameters  $\lambda, \mu, \gamma$  is fixed

parameters. UMIST has 20 classes, and it consists of 575 images with the dimension of 644. The best clustering accuracy is obtained when  $\lambda = 10, \mu = 0.001, \gamma = 1000$ .

We study the sensitivity of the three parameters by fixing one of them to the optimal settings. Figure 2 shows the clustering accuracy. From Figure 2(a), we can observe that when  $\lambda$  is fixed, a smaller  $\mu$  and a larger  $\gamma$  can lead to a better performance. Figure 2(b) shows that  $\lambda$  is more sensitive than  $\gamma$  and its optimal setting is  $\lambda = 10$ . A similar result can be observed in Figure 2(c). To sum up,  $\lambda$  is the most sensitive parameter, and a smaller  $\mu$  and larger  $\gamma$  can lead to a better performance.

We also test the sensitivity of the reduced dimensions of the low-dimensional space. As shown in Table 2, the most stable methods across different benchmark datasets are PCA, LPP, PUFE, and PUFE-OG. Thus, we mainly focus on these 4 methods. Figure 1 shows the results on 3 benchmark datasets. We can observe that PUFE and PUFE-OG have better performance when the reduced dimension is low. The performance of PCA and LPP become better as the dimension increases. But different from PCA and LPP, a high dimension degrades the performance of PUFE and PUFE-OG. This is because the quality of the manifold structure has a great influence to the performance. The local manifold structure is often more reliable (the  $k$  nearest neighbors are more likely from the same cluster when  $k$  is small) and the distance can be approximated to be linear for the local structure. When the dimension is low, PUFE-OG will focus on learning the local manifold structure and yield a significant performance gain. However, when dimension is too high, the linear measurement will become inaccurate. Then the performance will drop dramatically. Thus, PUFE-OG is more sensitive to dimension change compared with other methods, and a low dimension is preferable for PUFE-OG.

## 5 Conclusion

To cope with the out-of-sample problem, most unsupervised graph based dimension reduction techniques integrate the ridge regression term with the manifold regularizer. The performance of these graph based methods highly depends on the input affinity matrix. However, the graph might be sub-optimal for dimension reduction. To tackle this problem, we propose a novel projective unsupervised flexible embedding framework by constructing an optimal graph.

Instead of using the affinity matrix to build the graph, we construct a probability graph based on the Gaussian similarity graph. The probability graph is more smooth. Thus, an optimal graph can be learned and the experimental results prove that our PUF-OG model always outperforms the other projective dimension reduction techniques. We derive an efficient algorithm to solve the problem. The extensive experiments demonstrate the effectiveness and the stability of our model. This optimal graph construction method can also be applied to other graphs, such as  $l^1$  graph [80] and LLE graph. In the future, we will explore its application in other graph based methods, as well as its application in the semi-supervised setting.

## 6 Acknowledgement

This work has been partially supported by the EC project QoSTREAM.

## References

- [1] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*, 2001.
- [2] Israel-Dejene Gebru, Xavier Alameda-Pineda, Florence Forbes, and Radu Horaud. EM algorithms for weighted-data clustering with application to audio-visual scene analysis. *TPAMI*, 2016.
- [3] Athinodoros S Georgiades, Peter N Belhumeur, and David J Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *TPAMI*, 2001.
- [4] Daniel B Graham, Nigel Allinson, and M. Characterising virtual eigensignatures for general purpose face recognition. In *Face Recognition*. 1998.
- [5] Xiaofei He, Deng Cai, Shuicheng Yan, and Hong-Jiang Zhang. Neighborhood preserving embedding. In *ICCV*, 2005.
- [6] Jin Huang, Feiping Nie, and Heng Huang. Spectral rotation versus k-means in spectral clustering. In *AAAI*, 2013.
- [7] Jin Huang, Feiping Nie, Heng Huang, and Chris Ding. Supervised and projected sparse coding for image classification. In *AAAI*, 2013.
- [8] Thomas Hueber, Laurent Girin, Xavier Alameda-Pineda, and Gerard Bailly. Speaker-adaptive acoustic-articulatory inversion using cascaded gaussian mixture regression. *TASLP*, 2015.

- [9] B Boser Le Cun, John S Denker, D Henderson, Richard E Howard, W Hubbard, and Lawrence D Jackel. Handwritten digit recognition with a back-propagation network. In *NIPS*, 1990.
- [10] Zechao Li, Yi Yang, Jing Liu, Xiaofang Zhou, and Hanqing Lu. Unsupervised feature selection using nonnegative spectral analysis. In *AAAI*, 2012.
- [11] Jialu Liu, Chi Wang, Marina Danilevsky, and Jiawei Han. Large-scale spectral clustering on graphs. In *IJCAI*, 2013.
- [12] Dijun Luo, Chris HQ Ding, and Heng Huang. Linear discriminant analysis: New formulations and overfit analysis. In *AAAI*, 2011.
- [13] Michael J Lyons, Shigeru Akamatsu, Miyuki Kamachi, Jiro Gyoba, and Julien Budynek. The japanese female facial expression (jaffe) database. 1998.
- [14] Adeel Mumtaz, Emanuele Coviello, Gert RG Lanckriet, and Antoni B Chan. Clustering dynamic textures with the hierarchical em algorithm for modeling video. *TPAMI*, 2013.
- [15] Feiping Nie, Dong Xu, Ivor Wai-Hung Tsang, and Changshui Zhang. Flexible manifold embedding: A framework for semi-supervised and unsupervised dimension reduction. *TIP*, 2010.
- [16] Feiping Nie, Xiaoqian Wang, and Heng Huang. Clustering and projected clustering with adaptive neighbors. In *SIGKDD*. ACM, 2014.
- [17] X Niyogi. Locality preserving projections. In *NIPS*, 2004.
- [18] Ignacio Ramirez, Pablo Sprechmann, and Guillermo Sapiro. Classification and clustering via dictionary learning with structured incoherence and shared features. In *CVPR*, 2010.
- [19] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 2000.
- [20] Ling Shao, Di Wu, and Xuelong Li. Learning deep and wide: A spectral method for learning deep networks. *TNNLS*, 2014.
- [21] Vikas Sindhwani, Partha Niyogi, Mikhail Belkin, and Sathiya Keerthi. Linear manifold regularization for large scale semi-supervised learning. In *ICML Workshop on Learning with Partially Classified Training Data*, 2005.
- [22] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 2000.
- [23] Pascal Vin., Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008.
- [24] Wei Wang, Yan Yan, and Sebe Nicu. Attribute guided dictionary learning. In *ICMR*, 2015.
- [25] Wei Wang, Zhen Cui, Yan Yan, Jiashi Feng, Shuicheng Yan, Xiangbo Shu, and Sebe Nicu. Recurrent face aging. In *CVPR*, 2016.

- [26] Wei Wang, Yan Yan, Stefan Winkler, and Nicu Sebe. Category specific dictionary learning for attribute specific feature selection. *TIP*, 2016.
- [27] Daniel Wein., Remi Ronfard, and Edmond Boyer. Free viewpoint action recognition using motion history volumes. *CVIU*, 2006.
- [28] Daniel Wein., Mustafa Özuysal, and Pascal Fua. Making action recognition robust to occlusions and viewpoint changes. In *ECCV*. 2010.
- [29] Rongkai Xia, Yan Pan, Lei Du, and Jian Yin. Robust multi-view spectral clustering via low-rank and sparse decomposition. In *AAAI*, 2014.
- [30] Shuicheng Yan and Huan Wang. Semi-supervised learning by sparse representation. In *SDM*, 2009.
- [31] Shuicheng Yan, Dong Xu, Benyu Zhang, Hong-Jiang Zhang, Qiang Yang, and Stephen Lin. Graph embedding and extensions: a general framework for dimensionality reduction. *TPAMI*, 2007.
- [32] Yan Yan, Elisa Ricci, Ramanathan Subramanian, Gaowen Liu, and Nicu Sebe. Multi-task linear discriminant analysis for view invariant action recognition. *TIP*, 2014.
- [33] Yan Yan, Elisa Ricci, Ramanathan Subramanian, Gaowen Liu, Oswald Lanz, and Nicu Sebe. A multi-task learning framework for head pose estimation under target motion. *TPAMI*, 2016.
- [34] Yi Yang, Dong Xu, Feiping Nie, Shuicheng Yan, and Yueting Zhuang. Image clustering using local discriminant models and global integration. *TIP*, 2010.
- [35] Jieping Ye, Zheng Zhao, and Mingrui Wu. Discriminative k-means for clustering. In *NIPS*, 2008.