# Graph-adaptive Nonlinear Dimensionality Reduction

Yanning Shen, *Student Member, IEEE,* Panagiotis A. Traganitis, *Student Member, IEEE,*
and Georgios B. Giannakis, *Fellow, IEEE*

*Abstract*—In this era of data deluge, many signal processing and machine learning tasks are faced with high-dimensional datasets, including images, videos, as well as time series generated from social, commercial and brain network interactions. Their efficient processing calls for dimensionality reduction techniques capable of properly compressing the data while preserving task-related characteristics, going beyond pairwise data correlations. The present paper puts forth a nonlinear dimensionality reduction framework that accounts for data lying on known graphs. The novel framework encompasses most of the existing dimensionality reduction methods, but it is also capable of capturing and preserving possibly nonlinear correlations that are ignored by linear methods. Furthermore, it can take into account information from multiple graphs. The proposed algorithms were tested on synthetic as well as real datasets to corroborate their effectiveness.

*Index Terms*—Dimensionality reduction, nonlinear modeling, signal processing over graphs

## I. INTRODUCTION

The massive development of connected devices and highly precise instruments has introduced the world to vast volumes of high-dimensional data. Traditional data analytics cannot cope with these massive amounts, which motivates well investigating dimensionality reduction schemes capable of gleaning out efficiently low-dimensional information from large-scale datasets. Dimensionality reduction is a vital first step to render tractable critical learning tasks, such as large-scale regression, classification, and clustering of high-dimensional datasets. In addition, dimensionality reduction can allow for accurate visualization of high-dimensional datasets.

Dimensionality reduction methods have been extensively studied by the signal processing and machine learning communities [2], [12], [21], [23]. Principal component analysis (PCA) [12] is the 'workhorse' method yielding low-dimensional representations that preserve most of the high-dimensional data variance. Multi-dimensional scaling (MDS) [15] on the other hand, maintains pairwise distances between data when going from high- to low-dimensional spaces, while local linear embedding (LLE) [21] only preserves relationships between neighboring data. Information from non-neighboring data is lost in LLE's low-dimensional representation, which may in turn influence the performance of ensuing tasks such as classification or clustering [7], [30]. It

is also worth stressing that all aforementioned approaches capture and preserve linear dependencies among data. However, for data residing on highly nonlinear manifolds using only linear relations might produce low-dimensional representations that are not accurate. Generalizing PCA, kernel PCA (KPCA) can capture nonlinear relationships between data, for a preselected kernel function. In addition, Laplacian eigenmaps [2] preserve nonlinear similarities between neighboring data.

While all the aforementioned approaches have been successful in reducing the dimensionality of various types of data, they do not consider additional information during the dimensionality reduction process. This prior information may be task specific, e.g. provided by some "expert" or by the physics of the problem, or it could be inferred from alternative views of the data, and can provide additional insights for the desired properties of the low-dimensional representations. In fMRI signals for instance, in addition to time series collected at different brain regions, one may also have access to the structural connectivity patterns among these regions.

At the same time, data may arrive from multiple heterogeneous sources, e.g. in addition to fMRI time courses, electroencephalography time series might be available. While it is desirable to draw inferences from all these multimodal data, their heterogeneous nature inhibits the use of traditional statistical learning tools. Thus, schemes that can generate useful data representations by fusing judiciously the information contained in different data modes are required.

As shown in [10], [11], [24], [25] for PCA, useful additional information can be encoded in a graph, and incorporated into the dimensionality reduction process through *graph-adaptive* regularization. PCA accounting for the graph Laplacian has been advocated in [10], to improve performance by exploiting the underlying graph structure. A low rank matrix factorization method incorporating multiple graph regularizers for linear PCA can be found in [11]. However, a quadratic program must be solved per iteration to optimally combine the adopted graph regularizers. Robust versions of linear graph PCA have also been reported [24]. Multiple graph regularizers were also studied in [33] relying on low-rank matrix matrix factorization.

**Our contributions.** The present manuscript presents a novel *graph-adaptive* (GRAD) *nonlinear* dimensionality reduction approach, to account for prior information on one or multiple graphs. By extending the concept of kernel PCA to graphs, our approach encompasses all aforementioned approaches, while markedly broadening their scope. Compared to our conference precursor in [27], the present manuscript includes GRAD nonlinear dimensionality reduction when domain knowledge is unknown. To this end, a multi-kernel based approach is

developed that uses the data to select the appropriate kernel for the dimensionality reduction task. In addition, we show how our approach can reduce dimensionality of multi-modal datasets, by considering separate graphs induced by different modes. Further, we generalize our approach to semi-supervised scenarios, where labels for a few data are available.

The rest of the paper is organized as follows. Section II provides preliminaries along with notation and background works. Section III introduces the proposed GRAD nonlinear dimensionality reduction scheme, while Section IV provides pertinent generalizations and applications. Section V presents numerical tests conducted to evaluate the performance of the novel dimensionality reduction scheme. Finally, concluding remarks and future research directions are given in Section VI. **Notation:** Unless otherwise noted, lowercase bold letters $x$ denote vectors, uppercase bold letters $\mathbf{X}$ represent matrices, and calligraphic uppercase letters $\mathcal{X}$ stand for sets. The $(i, j)$th entry of $\mathbf{X}$ is denoted by $[\mathbf{X}]_{ij}$; $\mathbf{X}^\top$ denotes the transpose of $\mathbf{X}$, while $\mathbf{X}^\dagger$ denotes the Moore-Penrose pseudo-inverse of matrix $\mathbf{X}$. The $D$-dimensional real Euclidean space is denoted by $\mathbb{R}^D$, the set of positive real numbers by $\mathbb{R}_+$, the positive integers by $\mathbb{Z}_+$, and the $\ell_2$-norm by $\| \cdot \|$.

## II. PRELIMINARIES AND PROBLEM STATEMENT

Consider a dataset with $N$ vectors of dimension $D$ collected as columns of the matrix $\mathbf{Y} := [\mathbf{y}_1, \ldots, \mathbf{y}_N]$. Without loss of generality, it will be assumed that the data are centered, that is the sample mean $N^{-1}\sum_{i=1}^N \mathbf{y}_i$ has been removed from each $\mathbf{y}_i$. For future use, the singular value decomposition (SVD) of the data matrix $\mathbf{Y}$ is $\mathbf{Y} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$. Dimensionality reduction seeks a set of $d < D$-dimensional vectors $\{\boldsymbol{\psi}_i\}_{i=1}^N$, that preserve certain properties of the original data $\{\mathbf{y}_i\}_{i=1}^N$.

The following subsections review popular dimensionality reduction schemes, that can be viewed as special cases of kernel PCA.

### A. Principal component analysis

Given data $\mathbf{Y}$, PCA finds a linear subspace of dimension $d$ such that all the data lie on or close to it, in the Euclidean distance sense. Specifically, PCA solves

$$\min_{\mathbf{U}_d, \{\boldsymbol{\psi}_i\}} \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{U}_d \boldsymbol{\psi}_i\|_2^2 \quad \text{s. to} \quad \mathbf{U}_d^\top \mathbf{U}_d = \mathbf{I} \quad (1)$$

where $\mathbf{U}_d \in \mathbb{R}^{D \times d}$ is an orthonormal matrix whose columns span the sought subspace. The optimal solution of (1) is $\boldsymbol{\psi}_i = \mathbf{U}_d^\top \mathbf{y}_i$, where $\mathbf{U}_d$ is formed by the eigenvectors of $\mathbf{Y}\mathbf{Y}^\top = \mathbf{U}\boldsymbol{\Sigma}^2\mathbf{U}^\top$ corresponding to the $d$ eigenvalues with the largest magnitude, or equivalently to the $d$ leading left singular vectors of $\mathbf{Y}$ [8]. Given $\{\boldsymbol{\psi}_i\}$, the original vectors can be recovered as $\hat{\mathbf{y}}_i = \mathbf{U}_d \boldsymbol{\psi}_i$. PCA has well-documented merits when data lie close to a $d$-dimensional hyperplane. Its complexity is that of eigendecomposing $\mathbf{Y}\mathbf{Y}^\top$, i.e., $\mathcal{O}(ND^2)$, which means PCA is more affordable when $D \ll N$. In contrast, dimensionality reduction of small sets of high-dimensional vectors $(D \gg N)$ becomes more tractable with the dual PCA that we outline next.

### B. Dual PCA and Kernel PCA

Collect all the lower dimensional data representations as columns of the $d \times N$ matrix $\boldsymbol{\Psi} := [\boldsymbol{\psi}_1, \ldots, \boldsymbol{\psi}_N]$. Then using the SVD of $\mathbf{Y}$, we find

$$\boldsymbol{\Psi} = \mathbf{U}_d^\top \mathbf{Y} = \boldsymbol{\Sigma}_d \mathbf{V}_d^\top \quad (2)$$

where $\boldsymbol{\Sigma}_d \in \mathbb{R}^{d \times d}$ is a diagonal matrix containing the $d$ leading singular values of $\mathbf{Y}$, and $\mathbf{V}_d \in \mathbb{R}^{N \times d}$ is the submatrix of $\mathbf{V}$ collecting the corresponding right singular vectors of $\mathbf{Y}$. Since $\mathbf{Y}^\top \mathbf{Y} = \mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^\top$, the low-dimensional representations of data can be obtained through the eigendecomposition of $\mathbf{Y}^\top \mathbf{Y}$. Using this method to find the low-dimensional representations of the data is known as *Dual PCA*. As only eigendecomposition of $\mathbf{Y}^\top \mathbf{Y}$ is required, the complexity of dual PCA is $\mathcal{O}(DN^2)$; therefore, it is preferable when $D \gg N$. Moreover, it can be readily verified that besides (1), $\boldsymbol{\Psi}$ is the optimal solution to the following optimization problem (see Appendix A)

$$\min_{\boldsymbol{\Psi}} \|\mathbf{K}_y - \boldsymbol{\Psi}^\top \boldsymbol{\Psi}\|_F^2 \quad \text{s. to} \quad \boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \boldsymbol{\Lambda}_d \quad (3)$$

where $\mathbf{K}_y := \mathbf{Y}^\top \mathbf{Y}$ is known as the Gram or kernel matrix, and $\boldsymbol{\Lambda}_d$ denotes a $d \times d$ diagonal matrix containing the $d$ largest eigenvalues of $\mathbf{K}_y$. Compared to PCA, dual PCA requires only the inner products $\{\mathbf{y}_i^\top \mathbf{y}_j\}$ in order to obtain the low-dimensional representations. Hence, dual PCA can yield low-dimensional vectors $\{\boldsymbol{\psi}_i\}$ of general (non-metric) objects that are not necessarily expressed using vectors $\{\mathbf{y}_i\}$, so long as inner products (meaning correlations) of the latter are known. On the other hand, the original data $\{\mathbf{y}_i\}$ cannot be recovered from $\{\boldsymbol{\psi}_i\}$ found by the solution of (3).

Consider now expanding the cost in (3), to equivalently express it as

$$\min_{\boldsymbol{\Psi}: \boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \boldsymbol{\Lambda}_d} \text{tr}(\boldsymbol{\Psi}\mathbf{K}_y^{-1}\boldsymbol{\Psi}^\top) \quad (4)$$

Recalling that $\mathbf{K}_y^{-1}$ is symmetric and nonnegative definite with eigenvalues equal to the inverses of the eigenvalues of $\mathbf{K}_y$, we can re-write (4) as

$$\min_{\boldsymbol{\Psi}: \boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \boldsymbol{\Lambda}_d} -\text{tr}(\boldsymbol{\Psi}\mathbf{K}_y\boldsymbol{\Psi}^\top) \quad (5)$$

where now $\boldsymbol{\Lambda}_d$ contains the $d$ largest eigenvalues of $\mathbf{K}_y$.

| Kernel type | $\kappa(\mathbf{y}_i, \mathbf{y}_j)$ | Parameters |
|---|---|---|
| Linear | $\mathbf{y}_i^\top \mathbf{y}_j$ | - |
| Gaussian | $\exp\{\frac{-\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{2\sigma^2}\}$ | $\sigma > 0$ |
| Polynomial kernel | $(\mathbf{y}_i^\top \mathbf{y}_j + c)^p$ | $p > 0, c$ |

TABLE I: Examples of kernels.

While PCA performs well for data that lie close to a hyperplane, this property might not hold for the available data $\mathbf{Y}$ [11]. In such cases one may resort to kernel PCA. Kernel PCA "lifts" $\{\mathbf{y}_i\}$ using a nonlinear function $\phi$, onto a higher (possibly infinite) dimensional space, where the data may lie on or near a linear hyperplane, and then finds low-dimensional

representations $\{\psi_i\}$. Kernel PCA is obtained by solving (3) or (4) with $[\mathbf{K}_y]_{i,j} = \kappa(\mathbf{y}_i, \mathbf{y}_j) = \phi^\top(\mathbf{y}_i)\phi(\mathbf{y}_j)$, where $\kappa(\mathbf{y}_i, \mathbf{y}_j)$ denotes a prescribed kernel function [6]. Table I lists a few popular kernels used in the literature, including the linear kernel which links linear dual PCA with kernel PCA.

## C. Local linear embedding

Another popular method that deals with data that cannot be presumed close to a hyperplane is local linear embedding (LLE) [21]. LLE postulates that $\{\mathbf{y}_i\}$ lie on a smooth manifold, which can be locally approximated by tangential hyperplanes. Specifically, LLE assumes that each datum can be expressed as a linear combination of its neighbors; that is, $\mathbf{y}_i = \sum_{j \in \mathcal{N}_i} w_{ij}\mathbf{y}_j + \mathbf{e}_i$, where $\mathcal{N}_i$ is a set containing the indices of the nearest neighbors of $\mathbf{y}_i$, in the Euclidean distance sense, and $\mathbf{e}_i$ captures unmodeled dynamics.

In order to solve for $\{w_{ij}\}$, the following optimization problem is considered

$$\mathbf{W} = \arg\min_{\check{\mathbf{W}}} \|\mathbf{Y} - \mathbf{Y}\check{\mathbf{W}}\|_F^2$$
$$\text{s. to } \check{w}_{ij} = 0, \quad \forall i \notin \mathcal{N}_j, \quad \sum_i \check{w}_{ij} = 1 \quad (6)$$

where $\check{w}_{ij}$ denotes the $(i,j)$-th entry of $\check{\mathbf{W}}$. Upon obtaining $\mathbf{W}$ as the constrained least-squares solution of (6), LLE finds $\{\psi_i\}$ that best preserve the neighborhood relationships encoded in $\mathbf{W}$ also in the lower dimensional space, by solving

$$\min_{\mathbf{\Psi}} \|\mathbf{\Psi} - \mathbf{\Psi}\mathbf{W}\|_F^2$$
$$\text{s.to } \mathbf{\Psi}\mathbf{\Psi}^\top = \mathbf{\Lambda}_d \quad (7)$$

which is equivalent to

$$\min_{\mathbf{\Psi}} \text{tr}[\mathbf{\Psi}(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^\top \mathbf{\Psi}^\top]$$
$$\text{s. to } \mathbf{\Psi}\mathbf{\Psi}^\top = \mathbf{\Lambda}_d. \quad (8)$$

Conventional LLE adopts $\mathbf{\Lambda}_d = \mathbf{I}$, which is subsumed by the constraint in (7). Nonetheless, the difference is just a scaling of $\{\psi_i\}$ when $\mathbf{\Lambda}_d \neq \mathbf{I}$. If the diagonal of $\mathbf{\Lambda}_d$ collects the $d$ smallest eigenvalues of matrix $(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^\top$, then (7) is a special case of kernel PCA with [cf. (4)]

$$\mathbf{K}_y = [(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^\top]^\dagger. \quad (9)$$

Similarly, other popular dimensionality reduction methods such as multidimensional scaling (MDS) [15], Laplacian eigenmaps [2], and isometric feature mapping (ISOMAP) [31] can also be viewed as special cases of kernel PCA, by appropriately selecting $\mathbf{K}_y$ [5]. Thus, (4) can be viewed as an encompassing framework for nonlinear dimensionality reduction. This will be the foundation of the general GRAD methods we develop in Sec. III.

## D. PCA on graphs

In several application settings, structural information implying or being implied by dependencies is available, and can benefit the dimensionality reduction task. This knowledge can be encoded in a graph and embodied in $\mathbf{\Psi}$ via graph regularization. Specifically, suppose there exists a graph $\mathcal{G}$ over which the data is smooth; that is, vectors $\{\psi_i\}$ that correspond to connected nodes of $\mathcal{G}$ are close to each other in Euclidean distance. With $\mathbf{A}$ denoting the adjacency matrix of $\mathcal{G}$, we have $[\mathbf{A}]_{ij} = a_{ij} \neq 0$ if node $i$ is connected with node $j$. The Laplacian of $\mathcal{G}$ is $\mathbf{L}_{\mathcal{G}} := \mathbf{D} - \mathbf{A}$, where $\mathbf{D}$ is a diagonal matrix with entries $[\mathbf{D}]_{ii} = d_{ii} = \sum_j a_{ij}$. Now consider

$$\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}\mathbf{\Psi}^\top) = \sum_{i=1}^N \sum_{j \neq i}^N a_{ij}\|\psi_i - \psi_j\|_2^2 \quad (10)$$

which is a weighted sum of the distances of adjacent $\psi_i$'s on the graph. By minimizing (10) over $\mathbf{\Psi}$, the low-dimensional representations corresponding to adjacent nodes with large edge weights $a_{ij} > 0$ will be close to each other. Therefore, minimizing (10) promotes the smoothness of $\mathbf{\Psi}$ over the graph.

Augmenting the PCA cost function with the regularizer in (10), yields the graph-regularized PCA [11]

$$\min_{\mathbf{U}_d, \mathbf{\Psi}} \|\mathbf{Y} - \mathbf{U}_d\mathbf{\Psi}\|_F^2 + \lambda\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}\mathbf{\Psi}^\top) \quad (11)$$

where $\lambda > 0$ is the regularization parameter. Building upon (11), robust versions of graph-regularized PCA have also been developed in e.g. [24], [25]. Clearly, (11) accounts only for linear dependencies in the data.

## III. GRAD NONLINEAR DIMENSIONALITY REDUCTION

Other than data correlations, nonlinear dimensionality reduction schemes are not designed to take into account additional prior information. At the same time, PCA on graphs, while able to incorporate prior information in the form of a graph, assumes that data lie near a linear subspace. This section will present a novel approach to graph-adaptive nonlinear dimensionality reduction, that encompasses the aforementioned nonlinear dimensionality reduction schemes, as well as linear PCA on graphs.

## A. Kernel PCA on graphs

Consider the kernel PCA formulation of (3). As with regular PCA, this formulation can be readily augmented with a graph regularizer, to arrive at

$$\min_{\mathbf{\Psi}} \|\mathbf{K}_y - \mathbf{\Psi}^\top\mathbf{\Psi}\|_F^2 + \gamma\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}\mathbf{\Psi}^\top)$$
$$\text{s. to } \mathbf{\Psi}\mathbf{\Psi}^\top = \mathbf{\Lambda}_d \quad (12)$$

where $\gamma$ is a positive scalar, and $\mathbf{\Lambda}_d$ a diagonal matrix. Since the latter only influences the scaling of $\mathbf{\Psi}$, for brevity we will henceforth set $\mathbf{\Lambda}_d = \mathbf{I}_d$. As kernel PCA can be written as a trace minimization problem [cf. (5)], (12) reduces to

$$\min_{\mathbf{\Psi}} -\text{tr}(\mathbf{\Psi}\mathbf{K}_y\mathbf{\Psi}^\top) + \gamma\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}\mathbf{\Psi}^\top)$$
$$\text{s. to } \mathbf{\Psi}\mathbf{\Psi}^\top = \mathbf{I}. \quad (13)$$

Combining the Laplacian regularization with the kernel PCA formulation, (13) is capable of finding $\{\psi_i\}$ that preserve the

---

**Algorithm 1** Kernel PCA on graphs

---
**Input:** $\mathbf{K}_y$, $\mathbf{L}_{\mathcal{G}}$, $\gamma$, $d$
**S1.** Find $r(\mathbf{L}_{\mathcal{G}}) = \mathbf{b}$
**S2.** Find the $d$ largest eigenvalues and corresponding eigenvectors of $\mathbf{K}_y - \gamma r(\mathbf{L}_{\mathcal{G}})$ and collect them in $\mathbf{V}_d$.
**S2.** Find low-dimensional representations $\boldsymbol{\Psi} = \mathbf{V}_d^{\top}$.

---

---

**Algorithm 2** Multi-Kernel PCA on graphs

---
**Input:** $\{\mathbf{K}_y^q\}_{q=1}^Q$, $\mathbf{L}_{\mathcal{G}}$, $\gamma, d$
**while** not converged **do**
    **S1.** Let $\mathbf{K}_y = \sum_{q=1}^Q \theta_q \mathbf{K}_y^{(q)}$
    **S2.** Find $\boldsymbol{\Psi}$ via Algorithm 1
    **S3.** Update $\boldsymbol{\theta}$ using (17).
**end while**

---

"lifted" covariance captured by $\mathbf{K}_y$, while at the same time, promoting the smoothness of the low-dimensional representations over the graph $\mathcal{G}$. As a result of the Courant-Fisher characterization [22], and with $\bar{\mathbf{K}} = \mathbf{K}_y - \gamma \mathbf{L}_{\mathcal{G}} = \bar{\mathbf{V}} \boldsymbol{\Lambda} \bar{\mathbf{V}}^{\top}$, (13) admits a closed-form solution as $\boldsymbol{\Psi} = \bar{\mathbf{V}}_d^{\top}$, which denotes the sub-matrix of $\bar{\mathbf{V}}$ formed by columns corresponding to the $d$ largest eigenvalues. When $\gamma$ is set to 0, one readily obtains the solution of kernel PCA [cf. (2)].

In addition, instead of directly using $\mathbf{L}_{\mathcal{G}}$, a family of graph kernels $r^{\dagger}(\mathbf{L}_{\mathcal{G}}) := \mathbf{U}_{\mathcal{G}} r^{\dagger}(\boldsymbol{\Lambda}) \mathbf{U}_{\mathcal{G}}^{\top}$ can be employed. Here $r(.)$ is a non-decreasing scalar function of the eigenvalues of $\mathbf{L}_{\mathcal{G}}$, while $\mathbf{U}_{\mathcal{G}}$ contains the eigenvectors of $\mathbf{L}_{\mathcal{G}}$. Introducing $r^{\dagger}(\mathbf{L})$ as a kernel matrix, we have

$$\min_{\boldsymbol{\Psi}} -\text{tr}(\boldsymbol{\Psi} \mathbf{K}_y \boldsymbol{\Psi}^{\top}) - \gamma \text{tr}(\boldsymbol{\Psi} r^{\dagger}(\mathbf{L}_{\mathcal{G}}) \boldsymbol{\Psi}^{\top})$$
$$\text{s. to} \quad \boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} = \mathbf{I}. \tag{14}$$

By appropriately selecting $r(.)$, different graph properties can be accounted for. As an example, when $r$ sets eigenvalues above a certain threshold to 0, it acts as a sort of "low pass" filter over the graph. Examples of graph kernels are provided in Table II; see also [19], [20] on graph kernel options, and the graph properties they capture.

| Kernel type | Function | Parameters |
|---|---|---|
| Diffusion [14] | $r(\lambda) = \exp\{\sigma^2 \lambda / 2\}$ | $\sigma^2 \geq 0$ |
| $p$-step random walk [29] | $r(\lambda) = (a - \lambda)^{-p}$ | $a \geq 2, p$ |
| Regularized Laplacian [28], [29] | $r(\lambda) = 1 + \sigma^2 \lambda$ | $\sigma^2 \geq 0$ |
| Bandlimited [20] | $r(\lambda_n) = \begin{cases} 1/\beta & n \leq B \\ \beta & \text{o.w.} \end{cases}$ | $\beta, B > 0$ |

TABLE II: Examples of graph Laplacian kernels.

As with kernel PCA [cf. (3)] the performance of this approach relies critically on the choice of $\mathbf{K}_y$. To circumvent this limitation, the following subsection introduces a multi-kernel based approach to GRAD dimensionality reduction.

*B. Multi-kernel learning based approach*

In several application domains, the appropriate kernel for the dimensionality reduction task might be not known a priori. In such cases, one can resort to multi-kernel approaches. Multi-kernel methods select the appropriate kernel function as a linear combination of a number of preselected kernels [1]. Specifically, $\mathbf{K}_y$ can be formed as a linear combination of $Q$

kernel matrices as

$$\mathbf{K}_y = \sum_{q=1}^Q \theta_q \mathbf{K}_y^{(q)} \tag{15}$$

where $\{\mathbf{K}_y^{(q)}\}_{q=1}^Q$ are predetermined kernel matrices, and $\{\theta_q\}_{q=1}^Q$ are unknown non-negative combination weights. Since $\theta_q$'s are non-negative, the resulting $\mathbf{K}_y$ is also a valid kernel matrix. Multi-kernel methods "learn" the best kernel from the data, by optimizing over the combination weights $\{\theta_q\}_{q=1}^Q$. Incorporating (15) into (13), the pertinent optimization problem becomes

$$\min_{\boldsymbol{\theta}, \boldsymbol{\Psi}} \quad -\text{tr}(\boldsymbol{\Psi}(\sum_{q=1}^Q \theta_q \mathbf{K}_y^{(q)}) \boldsymbol{\Psi}^{\top}) - \gamma \text{tr}(\boldsymbol{\Psi} r^{\dagger}(\mathbf{L}_{\mathcal{G}}) \boldsymbol{\Psi}^{\top})$$
$$\text{s.t.} \quad \boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} = \mathbf{I}_d$$
$$\|\boldsymbol{\theta}\|_2^2 \leq 1, \quad \boldsymbol{\theta} \geq \mathbf{0} \tag{16}$$

where $\boldsymbol{\theta} := [\theta_1, \ldots, \theta_Q]^{\top}$, and the $\ell_2$-norm regularization is introduced to control the model complexity. As (16) is non-convex, it will be solved using alternating optimization. When $\{\theta_q\}$ are fixed, (16) can solved in closed form by eigenvalue decomposition of matrix $\sum_{q=1}^Q \theta_q \mathbf{K}_y^{(q)} + \gamma r^{\dagger}(\mathbf{L}_{\mathcal{G}})$, as in (14). With $\boldsymbol{\Psi}$ fixed, $\boldsymbol{\theta}$ is found as (see Appendix B for the proof)

$$\theta_q = \frac{\text{tr}(\boldsymbol{\Psi} \mathbf{K}_y^{(q)} \boldsymbol{\Psi}^{\top})}{\sqrt{\sum_{q=1}^Q (\text{tr}(\boldsymbol{\Psi} \mathbf{K}_y^{(q)} \boldsymbol{\Psi}^{\top}))^2}}, \quad q = 1 \ldots, Q. \tag{17}$$

The overall GRAD multi-kernel(MK)-PCA scheme is tabulated in Algorithm 2.

Even though only a single graph regularizer is introduced in (13), our method is flexible to include *multiple* graph regularizers based on different graphs. Therefore, the proposed method offers a powerful tool for dimensionality reduction with prior information encoded in the so-called *multi-layer* graphs [13], [32]. Suppose that $L$ such $N$-node graphs $\{\mathcal{G}_{\ell}\}_{\ell=1}^L$ are available, each with corresponding Laplacian matrices $\{\mathbf{L}_{\mathcal{G}}^{\ell}\}_{\ell=1}^L$. If all $L$ graphs are expected to have the same contribution then multiple Laplacian regularizers, say one per graph, can be introduced in the objective function of (16) as

$$\min_{\boldsymbol{\theta}, \boldsymbol{\Psi}} \quad -\text{tr}(\boldsymbol{\Psi}(\sum_{q=1}^Q \theta_q \mathbf{K}_y^{(q)}) \boldsymbol{\Psi}^{\top}) - \gamma \sum_{\ell=1}^L \text{tr}(\boldsymbol{\Psi} r^{\dagger}(\mathbf{L}_{\mathcal{G}}^{\ell}) \boldsymbol{\Psi}^{\top})$$
$$\text{s.t.} \quad \boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} = \mathbf{I}_d$$
$$\|\boldsymbol{\theta}\|_2^2 \leq 1, \quad \boldsymbol{\theta} \geq \mathbf{0}. \tag{18}$$

When the appropriate graph regularizer is unknown, a scheme

| Method | Formulation | Graphs | Factors | Kernels |
|--------|-------------|--------|---------|---------|
| PCA [12] | $\min_{\mathbf{U},\boldsymbol{\Psi}} \|\mathbf{Y}-\mathbf{U}\boldsymbol{\Psi}\|_F^2$ | No | Yes | No |
| GLPCA [10] | $\min_{\mathbf{U},\boldsymbol{\Psi}} \|\mathbf{Y}-\mathbf{U}\boldsymbol{\Psi}\|_F^2 + \lambda\mathrm{tr}(\boldsymbol{\Psi}\mathbf{L}\boldsymbol{\Psi}^\top)$ | Single | Yes | No |
| LapEmb [2] | $\min_{\boldsymbol{\Psi}} \mathrm{tr}(\boldsymbol{\Psi}\mathbf{L}\boldsymbol{\Psi}^\top)$ | Single | Yes | No |
| RPCAG [25] | $\min_{\mathbf{Z},\mathbf{S}} \|\mathbf{Z}\|_* + \gamma\|\mathbf{S}\|_1 + \lambda\mathrm{tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^\top)$ | Single | No | No |
| FRPCAG [24] | $\min_{\mathbf{Z},\mathbf{S}} \|\mathbf{Y}-\mathbf{Z}\|_1 + \gamma\|\mathbf{S}\|_1 + \lambda_1\mathrm{tr}(\mathbf{Z}\mathbf{L}_1\mathbf{Z}^\top) + \lambda_2\mathrm{tr}(\mathbf{Z}^\top\mathbf{L}_1\mathbf{Z})$ | Two | No | No |
| GRAD KPCA | $\min_{\boldsymbol{\theta},\boldsymbol{\Psi},\boldsymbol{\beta}} -\mathrm{tr}(\boldsymbol{\Psi}(\sum_{q=1}^Q \theta_q\mathbf{K}_y^{(q)})\boldsymbol{\Psi}^\top) - \gamma\mathrm{tr}(\boldsymbol{\Psi}(\sum_{m=1}^M \beta_m r^{-1}(\mathbf{L}_\mathcal{G}^m))\boldsymbol{\Psi}^\top)$ | Multiple | Yes | Multiple |

TABLE III: Comparison of graph-regularized PCA methods.

similar to the multi-kernel approach of (16) can be employed to choose the appropriate graph kernel. In this case, the graph kernel can be expressed as a linear combination of the $M$ available graph regularizers, that is $r^\dagger(\mathbf{L}_\mathcal{G}) = \sum_{\ell=1}^L \beta_\ell r^\dagger(\mathbf{L}_\mathcal{G}^\ell)$, where $\beta_\ell$ are unknown non-negative combination weights. Introducing this multi-graph kernel term into (16) yields

$$\min_{\boldsymbol{\theta},\boldsymbol{\Psi},\boldsymbol{\beta}} \quad -\mathrm{tr}(\boldsymbol{\Psi}(\sum_{q=1}^Q \theta_q\mathbf{K}_y^{(q)})\boldsymbol{\Psi}^\top) - \gamma\mathrm{tr}(\boldsymbol{\Psi}(\sum_{\ell=1}^L \beta_\ell r^\dagger(\mathbf{L}_\mathcal{G}^\ell))\boldsymbol{\Psi}^\top)$$
$$\text{s.t.} \quad \boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \mathbf{I}_d$$
$$\|\boldsymbol{\theta}\|_2^2 \leq 1, \quad \boldsymbol{\theta} \geq \mathbf{0}$$
$$\|\boldsymbol{\beta}\|_2^2 \leq 1, \quad \boldsymbol{\beta} \geq \mathbf{0} \tag{19}$$

where $\boldsymbol{\beta} := [\beta_1,\ldots,\beta_L]^\top$. Similar to (16), the non-convex problem in (19) can be solved in an alternating fashion. With $\boldsymbol{\beta}$ fixed, $\boldsymbol{\Psi}$ can be found in closed form by eigenvalue decomposition of matrix $\sum_{q=1}^Q \theta_q\mathbf{K}_y^{(q)} + \gamma\sum_{\ell=1}^L \beta_\ell r^\dagger(\mathbf{L}_\mathcal{G}^\ell)$, while $\boldsymbol{\theta}$ can be obtained using (17). When $\boldsymbol{\Psi}$ and $\boldsymbol{\theta}$ are fixed, $\boldsymbol{\beta}$ can be found in closed form as

$$\beta_\ell = \frac{\mathrm{tr}(\boldsymbol{\Psi} r^\dagger(\mathbf{L}_\mathcal{G}^\ell))\boldsymbol{\Psi}^\top)}{\sqrt{\sum_{\ell=1}^L (\mathrm{tr}(\boldsymbol{\Psi} r^\dagger(\mathbf{L}_\mathcal{G}^\ell))\boldsymbol{\Psi}^\top))^2}}, \quad \ell = 1,\ldots,L. \tag{20}$$

In this section, we put forth a novel scheme for dimensionality reduction over graphs that can also capture nonlinear data dependencies; see also Table III that summarizes how the novel approach fits within the context of prior relevant works. This table showcases the optimization problem solved by each algorithm, as well as whether prior information in the form of a graph can be incorporated. In addition, Table III indicates if low-dimensional representations are directly provided by an algorithm, and whether kernels have been employed to capture data correlations.

## IV. GENERALIZATIONS

The present section showcases three generalizations and applications of our novel GRAD dimensionality reduction scheme. Specifically, the following subsections extend the methods of Sec. III to multi-modal datasets, and semi-supervised settings, as well as generalize the LLE.

### A. Dimensionality reduction for multi-modal datasets

As discussed in Sec. I, many datasets comprise multi-modal data, that is data with features belonging to different types, such as binary, categorical or real-valued features. In this subsection, we demonstrate how our proposed GRAD nonlinear dimensionality reduction approach can readily handle such cases.

Suppose that the collected $N$ data contain $M$ different modes. Vectors of mode $m$ have dimension $D_m$, and are collected in a $D_m \times N$ submatrix $\mathbf{Y}_m$. With these $M$ sets of vectors at hand, $M$ different graphs $\{\mathcal{G}_m\}_{m=1}^M$, each with $N$ nodes can be inferred, based on possibly diverse similarity metrics. These metrics can be different for each mode, e.g. graphs for binary data can be constructed based on the Hamming distance, while graphs for real-valued data can be based on linear or nonlinear correlations. These $M$ graphs can be considered as an $M$-layer multiplex graph [13], on which our proposed scheme can be readily applied. Specifically, given the Laplacian matrices for each of these $M$ graphs $\{\mathbf{L}_\mathcal{G}^1,\cdots,\mathbf{L}_\mathcal{G}^M\}$, lower dimensional representations can be obtained as

$$\min_{\boldsymbol{\Psi}} -\sum_{m=1}^M \mathrm{tr}(\boldsymbol{\Psi} r^\dagger(\mathbf{L}_\mathcal{G}^m)\boldsymbol{\Psi}^\top)$$
$$\text{s. to} \quad \boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \mathbf{I}_d. \tag{21}$$

Therefore, the complexity of GRAD dimensionality reduction is in the order of $\mathcal{O}(DN^2)$, which is the same as the dual PCA. However, the graph-based PCA now can handle data that consist of heterogeneous features, e.g. binary, categorical or real-valued.

This scheme can also be used for dimensionality reduction of very high-dimensional data ($D \gg$). The $D \times N$ data matrix $\mathbf{Y}$ can be split, into $M$ submatrices $\{\mathbf{Y}_m\}_{m=1}^M$ each of dimension $D_m \times N$. These submatrices may contain non-overlapping or overlapping subsets of each data vector $\{\boldsymbol{y}_i\}_{i=1}^N$. Creating a graph for each $\mathbf{Y}_m$, (21) can be used to find lower dimensional representations $\boldsymbol{\Psi}$.

### B. Semi-supervised dimensionality reduction over graphs

In this subsection, we develop our proposed scheme for semi-supervised dimensionality reduction. In addition to data samples $\{\mathbf{y}_i\}_{i=1}^N$, domain knowledge here becomes available in the form of a few pairwise constraints. These constraints specify whether a pair of data samples belong to the same class (must-link constraints), or to different classes (cannot-link constraints). Specifically, let $\mathcal{S}$ be the set containing the tuples

**Algorithm 3** Local nonlinear embedding over graphs
***
**Input:** $\mathbf{Y}$, $\mathbf{L}_{\mathcal{G}}$ $\gamma, d$
**S1.** Estimate $\mathbf{W}$ from $\mathbf{Y}$.
**S2.** Obtain kernel matrix $\mathbf{K}_y$ via (9).
**S3.** Find $\mathbf{\Psi}$ as the leading eigenvectors of $\mathbf{K}_y$.
***

$(i, j)$ for some data belonging to the same class (must-link constraints), and $\mathcal{D}$ the set containing the tuples corresponding to data from different classes (cannot-link constraints). Given these two sets, two graphs can be constructed, one for each constraint set. The graph $\mathcal{G}^{\mathcal{S}}$ is constructed based on the must-link constraints with adjacency matrix $\mathbf{A}^{\mathcal{S}}$ having entries

$$[\mathbf{A}^{\mathcal{S}}]_{ij} = \begin{cases} 1, & \text{if } (i,j) \in \mathcal{S} \\ 0, & \text{otherwise.} \end{cases} \tag{22}$$

Similarly, the graph $\mathcal{G}^{\mathcal{D}}$ is constructed based on the cannot-link constraints and its adjacency matrix $\mathbf{A}^{\mathcal{D}}$ has entries

$$[\mathbf{A}^{\mathcal{D}}]_{ij} = \begin{cases} 1, & \text{if } (i,j) \in \mathcal{D} \\ 0, & \text{otherwise.} \end{cases} \tag{23}$$

Letting $\mathbf{L}_{\mathcal{G}}^{\mathcal{S}}$, $\mathbf{L}_{\mathcal{G}}^{\mathcal{D}}$ denote the graph Laplacians of $\mathcal{G}^{\mathcal{S}}$ and $\mathcal{G}^{\mathcal{D}}$ respectively, the low-dimensional representations of $\mathbf{Y}$ can be obtained as follows

$$\min_{\mathbf{\Psi}} -\text{tr}(\mathbf{\Psi}\mathbf{K}_y\mathbf{\Psi}^{\top}) + \gamma_1\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}^{\mathcal{S}}\mathbf{\Psi}^{\top}) - \gamma_2\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}^{\mathcal{D}}\mathbf{\Psi}^{\top})$$
$$\text{s.t.} \mathbf{\Psi}\mathbf{\Psi}^{\top} = \mathbf{I} \tag{24}$$

where $\gamma_1, \gamma_2 > 0$ are regularization constants.

Clearly, the term $\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}^{\mathcal{S}}\mathbf{\Psi}^{\top})$ forces the low-dimensional representations corresponding to the must-link constraints to be close, while the term $-\text{tr}(\mathbf{\Psi}\mathbf{L}_{\mathcal{G}}^{\mathcal{D}}\mathbf{\Psi}^{\top})$ "pushes" data corresponding to the cannot-link constraints away from each other. The GRAD regularizers effecting these two constraints are well motivated when one is interested in classifying high-dimensional vectors. If only a few of these vectors are labeled, such a semi-supervised setting should be accounted for in obtaining the low-dimensional representations based on which classification is to be performed subsequently.

### C. Local nonlinear embedding on graphs

In this subsection, we develop a major GRAD enhancement of the well appreciated nonlinear dimensionality reduction effected by LLE. We refer to our novel scheme as local nonlinear embedding on graphs (LNEG), because it can capture both linear and nonlinear dependencies among neighboring data, in addition to the structure induced by the graph $\mathcal{G}$. To this end, suppose that each data vector can be represented by its neighbors entry-wise as

$$[\mathbf{y}_i]_m = \sum_{j \in \mathcal{N}_i} h_{ij}([\mathbf{y}_j]_m) + [\mathbf{e}_i]_m, \quad m = 1, \ldots, D \tag{25}$$

where $\{h_{ij}(\cdot)\}_{i,j=1}^{N}$ are prescribed scalar nonlinear functions admitting a $P$th-order expansion

$$h_{ij}(z) = \sum_{p=1}^{P} w_{ij}[p]z^p \tag{26}$$

and coefficients $\{w_{ij}[p]\}$ are to be determined. Taylor's expansion asserts that for $P$ sufficiently large, (26) offers an accurate approximation for all memoryless differentiable nonlinear functions. Such a nonlinear model has been used for graph topology identification [26], but we here employ it as a first step of our LNEG scheme implementing the local nonlinear embedding. In vector form, (25) becomes

$$\bar{\mathbf{y}}_m^{\top} = \tilde{\mathbf{y}}_m^{\top}\tilde{\mathbf{W}} + \mathbf{e}_m \tag{27}$$

where $\bar{\mathbf{y}}_m^{\top} := [y_{1m} \ldots y_{Nm}]$ denotes the $m$-th row of $\mathbf{Y}$; the extended vector on the right hand side of (27) is $\tilde{\mathbf{y}}_m^{\top} := [\tilde{\mathbf{y}}_{1m}^{\top} \ \tilde{\mathbf{y}}_{2m}^{\top} \ \cdots \ \tilde{\mathbf{y}}_{Nm}^{\top}]$ formed with sub-vectors $\tilde{\mathbf{y}}_{im} := [y_{im}, y_{im}^2, \cdots, y_{im}^P]^{\top}$; and, the $N \times NP$ matrix $\tilde{\mathbf{W}}$ is defined as

$$\tilde{\mathbf{W}} := \begin{bmatrix} \mathbf{w}_{11} & \cdots & \mathbf{w}_{1N} \\ \vdots & & \vdots \\ \mathbf{w}_{N1} & \cdots & \mathbf{w}_{NN} \end{bmatrix} \tag{28}$$

where the entries of $\mathbf{w}_{ij} := [w_{ij}[1], \ldots, w_{ij}[P]]^{\top}$ are the coefficients in (26), specifying the nonlinear correlations between data. Upon defining $\hat{\mathbf{Y}} := [\tilde{\mathbf{y}}_1 \ \cdots \ \tilde{\mathbf{y}}_D]^{\top}$, one obtains the following nonlinear matrix model

$$\mathbf{Y} = \tilde{\mathbf{Y}}\tilde{\mathbf{W}} + \mathbf{E}. \tag{29}$$

Matrix $\tilde{\mathbf{W}}$ can be estimated using the least-squares (LS) or sparse regularized LS criteria, e.g.,

$$\mathbf{W}^* = \arg\min \|\mathbf{Y} - \tilde{\mathbf{Y}}\tilde{\mathbf{W}}\|_F^2 + \|\tilde{\mathbf{W}}\|_1 \tag{30}$$

which is convex but non-smooth, and can be solved iteratively to attain the global optimum using proximal splitting methods, see e.g. [26] for details. Using $\mathbf{W}^*$, an $N \times N$ matrix $\mathbf{W}$, similar to the one derived for LLE, can be obtained. Different from LLE, where $\mathbf{W}$ specifies tangential hyperplanes, our generalization here allows the local geometry to be captured by tangential nonlinear manifolds. Since $h$ can also be linear, LNEG is expected to perform at least as well as the LLE. With the estimated $\mathbf{W}$ at hand, the low-dimensional representations can be obtained via (8); see also Algorithm 3.

## V. NUMERICAL TESTS

The performance of our proposed algorithms, as well as their generalizations are tested in the present section. Numerical tests are carried out on both synthetic and real datasets. The performance of the dimensionality reduction task is evaluated through classification and clustering experiments. Specifically, the clustering and classification algorithms used are $K$-means and support vector machines (SVMs), respectively. The software used to conduct all experiments is MATLAB [17]. Reported results represent the averages of 50 independent
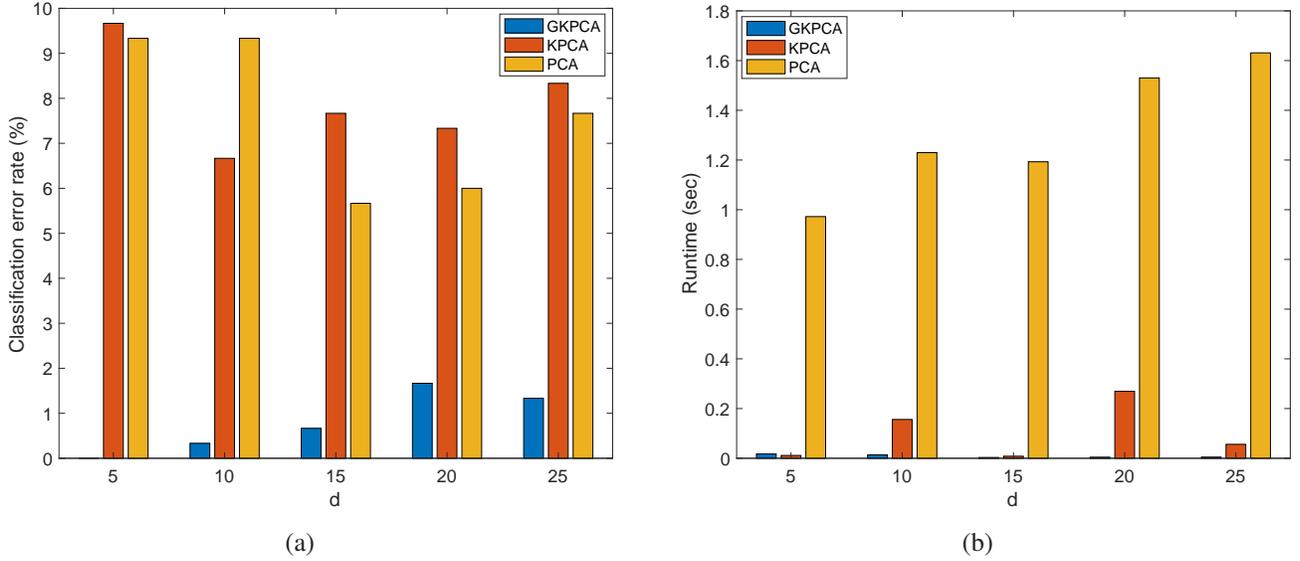
Fig. 1: Classification based on $\{\psi_i\}_{i=1}^N$ assessed by: (a) Classification error rate; and, (b) Running time;
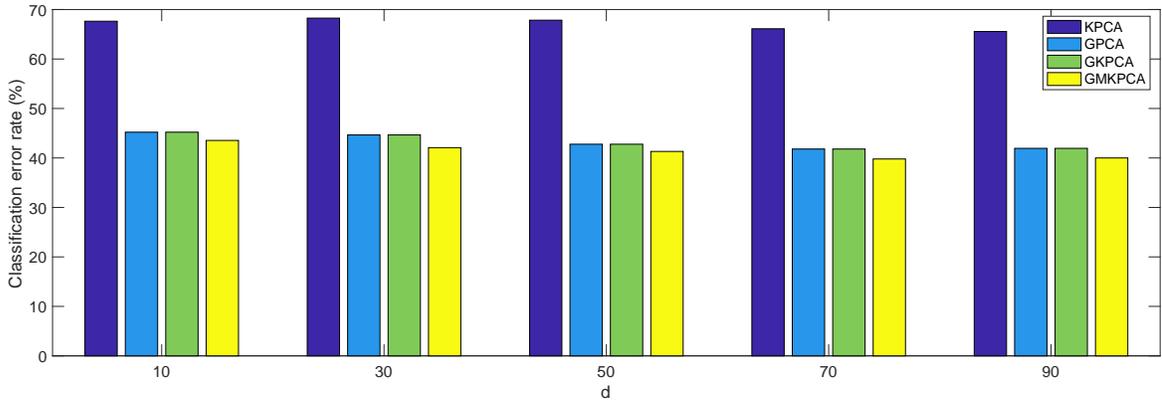


Fig. 2: Clustering results using $\{\psi_i\}_{i=1}^N$ on COIL20 dataset.

Monte Carlo runs. For both clustering and classification tests, performance is measured using the *error rate*, which is defined as the percentage of mis-clustered/ misclassified samples:

$$\text{Error Rate} := 1 - \frac{\text{\# data correctly clustered/classified}}{N} \times 100\%.$$

The datasets used are the following:

- USPS image dataset [9]: This consists of $N = 9,298$ images of size $16 \times 16$. Each image contains a digit scanned from U.S. Postal Service envelopes, and the dataset consists of $K = 10$ classes, one per digit.
- Coil20 image dataset [18]: This contains $N = 1,440$ $32 \times 32$ images of $K = 20$ objects. For each object 72 images are available, each taken under a different angle.
- Drivface image dataset [3]: This consists of $N = 66$ images of size $80 \times 80$, depicting images of drivers from two different angles, front images or side images.

Properties of these datasets are summarized in Table IV.

| Dataset | $N$ samples | $D$ features | $K$ classes |
|---------|-------------|--------------|-------------|
| Drivface | 66 | $6,400$ | 2 |
| USPS | $9,298$ | 256 | 10 |
| COIL20 | $1,440$ | $1,020$ | 20 |

TABLE IV: Datasets description.

### A. Graph Kernel PCA and Graph Multi-kernel PCA

In this section, the performance of the GRAD MK-PCA and GRAD KPCA algorithms is evaluated using both classification and clustering tests.

**Classification experiment.** In this experiment, Algorithm 1 (abbreviated henceforth as *GKPCA*) is tested on the Drivface dataset. The vectorized images $\mathbf{y}_i \in \mathbb{R}^{6,400}$ are used as columns of $\mathbf{Y}$. Here GKPCA is compared to PCA and KPCA. For the novel GKPCA and the KPCA algorithm, a Gaussian
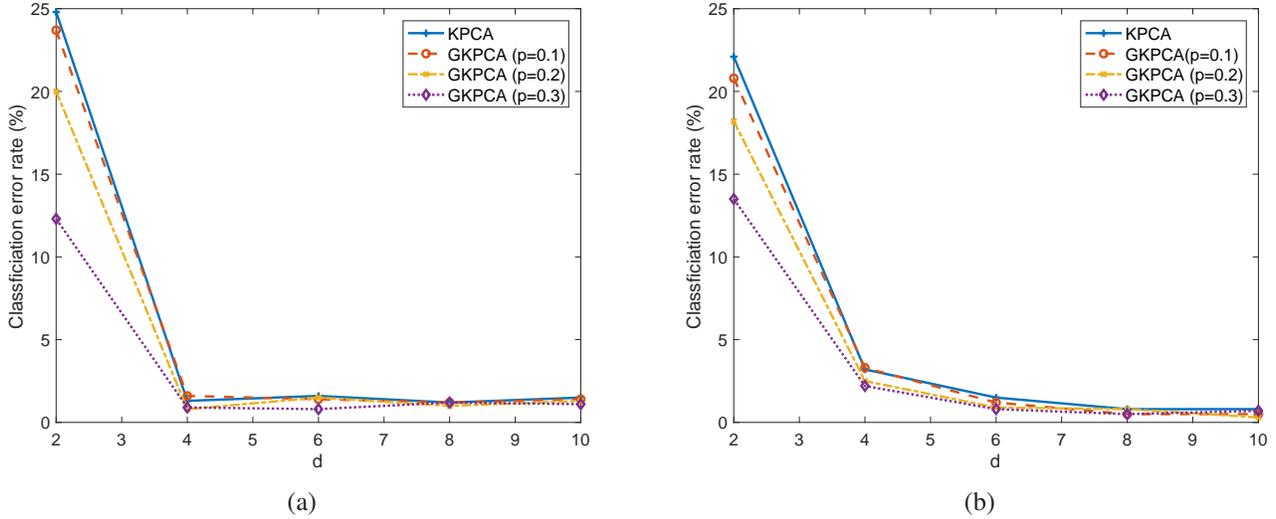
Fig. 3: Classification for USPS dataset using $\{\psi_i\}_{i=1}^N$ for (a) Digits 5 and 6 (b) Digits 7 and 8.
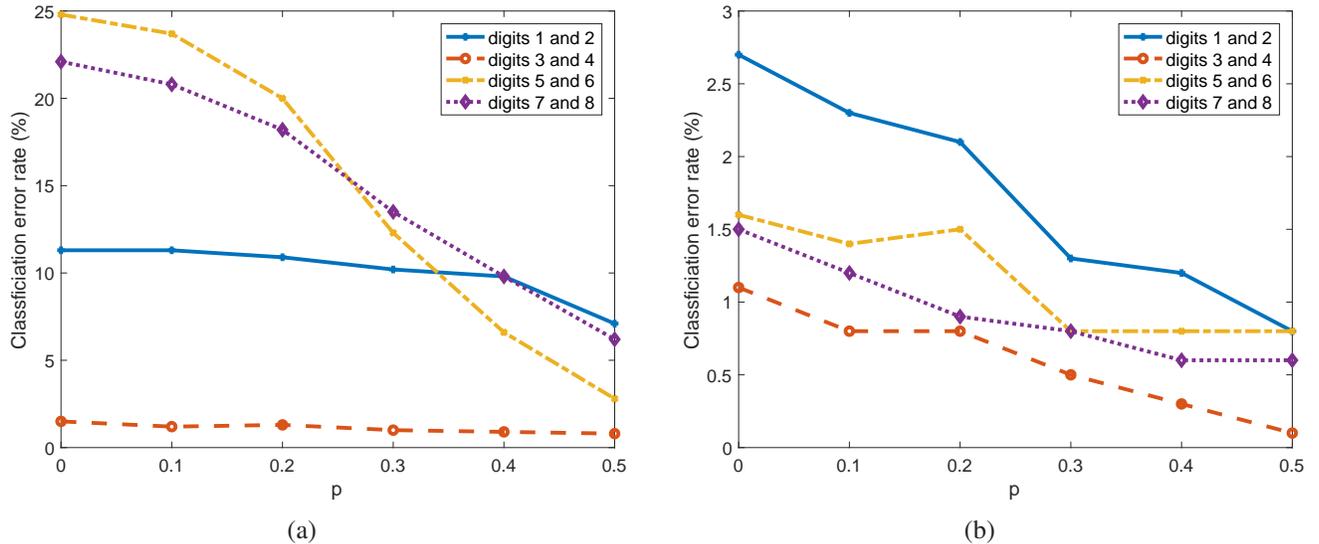


Fig. 4: Classification for USPS dataset using $\{\psi_i\}_{i=1}^N$ with: (a) $d = 2$; and (b) $d = 6$;

kernel with bandwidth $\sigma^2 = 1$ is employed. For GKPCA, the graph is constructed by pairwise linear correlation coefficients of feature vectors $\{\mathbf{f}_i\}$ extracted from the facial images, each $\mathbf{f}_i$ collecting the coordinates of nose, eyes and ears in the picture. Note that, this feature information is provided in the dataset.

Each dimensionality reduction algorithm is applied on $\{\mathbf{y}_i\}$ and low dimensional representations $\{\psi_i\}$ are obtained. Upon obtaining $\{\psi_i\}$, classification is performed using a linear SVM with 5 fold cross validation, with $80\%$ of the data used for training and the remaining $20\%$ for testing.

Figure 1 (a) depicts the testing classification error rate for different values of $d$. Clearly, the novel GKPCA approach outperforms both KPCA and PCA. In addition, Figure 1 (b) shows the runtime of different algorithms and corroborates that the kernel based approaches perform much faster that PCA,

because $D \gg N$. It can also be seen that GKPCA is more computationally efficient than Kernel PCA for most values of $d$. This is due to the graph regularization, which makes the $\bar{\mathbf{K}}$ [cf. Sec. III-A] matrix well-conditioned, and thus speeds up the eigenvalue decomposition.

**Clustering experiment.** In this experiment, the clustering performance was tested using $\{\psi_i\}$ obtained from different algorithms. GKPCA and Alg. 2 (termed henceforth as *GMKPCA*) are compared to KPCA and GPCA. For the GKPCA and KPCA algorithms, a Gaussian kernel with bandwidth $\sigma^2 = 1$ is employed. For GPCA, GKPCA and GMKPCA, the graph used is constructed by finding the pairwise correlation coefficients, $\bar{a}_{ij} = \frac{\mathbf{y}_i^\top \mathbf{y}_j}{\|\mathbf{y}_i\|_2 \|\mathbf{y}_j\|_2}$, and connecting each data sample with its 100 neighbors having the largest $\bar{a}_{ij}$; that is, $a_{ij} = \bar{a}_{ij}$ if $j \in \mathcal{N}_i$, otherwise $a_{ij} = 0$. The GMKPCA uses a dictionary

of Gaussian kernels with bandwidths $\sigma^2$ taking 10 equispaced values from $0.01$ to $1$. The $K$-means algorithm was repeated 50 times and the best result was reported. Figure 2 shows the clustering error rates for different algorithms, and after varying $d$ for the Drivface dataset. Clearly, GMKPCA outperforms the alternatives for clustering tasks.

## B. Semi-supervised graph-based dimensionality reduction

In this subsection, the semi-supervised dimensionality reduction scheme of Sec. IV-B is evaluated on the USPS dataset. For each experiment, a set of $1,000$ images of two different digits is used, with $500$ images for each digit. After obtaining low dimensional representations, a linear SVM classifier with 5-fold cross validation was used to distinguish the two digits.

Let $\Omega$ denote the set containing the indices of data for which labels are available. Two graphs were generated using (22) and (23) based on the known labels. Figure 3 showcases the performance of Algorithm 1 as a function of $d$ and for different numbers of available labels $|\Omega| = p \times N$, when classifying the digits $5$ and $6$ or $7$ and $8$. The available labeled data are chosen uniformly at random for each experiment. As the number of known labels increases, the classification error rate also decreases. Figure 4 depicts the classification error rate for a variable number of available labels used to classify different digits with $d = 2$ and $d = 6$, respectively. During all experiments, $\gamma_1$ and $\gamma_2$ were both set at $0.5$, and a Gaussian kernel with $\sigma^2 = 1$ was used. Clearly, this semi-supervised scheme can successfully incorporate label information into the graph-adaptive nonlinear dimensionality reduction task, such that the ensuing classification performance is improved.

## C. Local Nonlinear Embedding

Algorithm 3 is tested using $\mathbf{K}_y$ as in (9) for the locally nonlinear embedding (LNE) without and with graph regularization (the latter abbreviated as LNEG), both also compared with LLE and PCA. For all experiments, the graph $\mathcal{G}$ is constructed with adjacency matrix $\mathbf{A}$ with $(i, j)$th entry $a_{ij} = \mathbf{y}_i^\top \mathbf{y}_j / \|\mathbf{y}_i\| \|\mathbf{y}_j\|$. Two types of tests are carried out in order to: a) evaluate embedding performance for a single manifold; and b) assess how informative the low-dimensional embeddings are for distinguishing different manifolds.
**Embedding experiment.** In this experiment, we test the embedding performance of the proposed method. A 3-dimensional Swiss roll manifold is generated, and $1,000$ data are randomly sampled from the manifold as shown in Figure 5 (a). Figure 5 (b) illustrates the 2-dimensional embeddings obtained from PCA, while Figs. 5 (c) and (d) illustrate the resulting embeddings from LLE and LNEG respectively, where neighborhoods of $k = 20$ data are considered. Figs. 5 (e) and (f) depict embeddings obtained by LLE and LNEG with $k = 40$. The regularization parameter of LNEG is set to $\gamma = 0.1$, and the polynomial order is set to $P = 2$. Clearly, by exploiting the nonlinear relationships between data, the resulting low-dimensional representations are capable of better preserving the structure of the manifold, thus allowing for more accurate visualization.
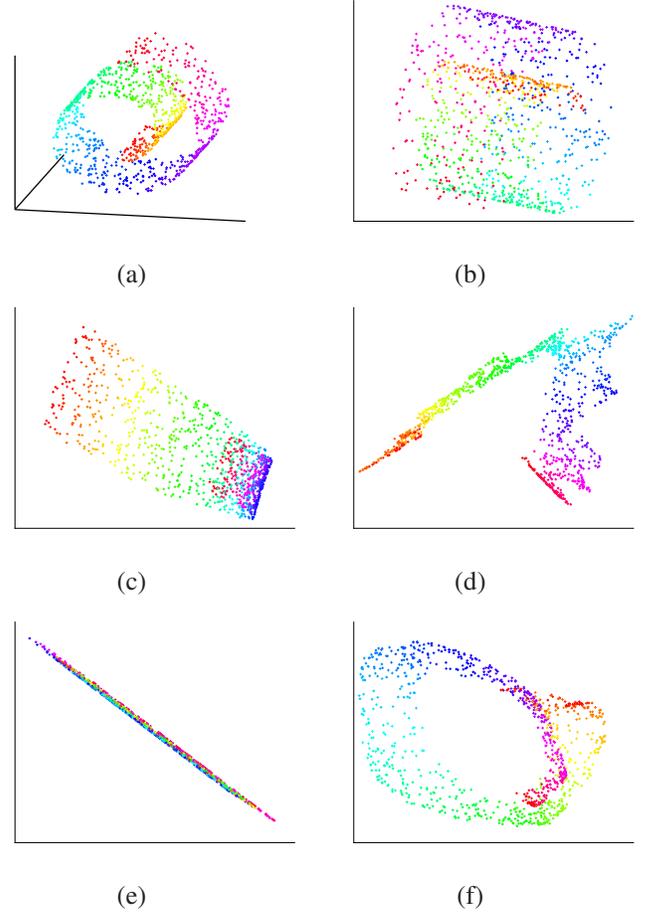


(a)  (b)

(c)  (d)

(e)  (f)

Fig. 5: Embedding results of two manifolds: linear hyperplane and trefoil (a) visualization of $\{\mathbf{y}_i\}_{i=1}^N$; and $\{\boldsymbol{\psi}_i\}_{i=1}^N$ obtained from (b) PCA; (c) LLE with $K = 20$; (d) LNEG with $K = 20$; (e) LLE with $K = 40$; (f) LNEG with $K = 40$.

**Clustering experiment.** In this experiment, the ability of Algorithm 3 to provide meaningful embeddings for clustering of different manifolds is assessed. Two 3-dimensional manifolds, a linear hyperplane with a hole around the origin and a trefoil are generated on the same ambient space as per [4], and 200 and 400 data are sampled from them. Here each manifold corresponds to a different cluster. Figure 6(a) illustrates the sampled points from the generated manifolds. Matrices $\mathbf{Z}_1 \in \mathbb{R}^{3 \times 200}$ and $\mathbf{Z}_2 \in \mathbb{R}^{3 \times 400}$ contain the data generated from the linear hyperplane and the trefoil. Both manifolds are then linearly embedded in $\mathbb{R}^{100}$, that is $\mathbf{Y}_i = \mathbf{P}\mathbf{Z}_i + \mathbf{E}_i$, where $\mathbf{P} \in \mathbb{R}^{100 \times 3}$ is an orthonormal matrix, and $\mathbf{E}$ is a noise matrix with entries sampled from a zero mean Gaussian distribution with variance $0.01$. Afterwards, the 100-dimensional data in $\mathbf{Y} := [\mathbf{Y}_1 \ \mathbf{Y}_2]$ are embedded into 2-dimensional representations $\mathbf{\Psi} \in \mathbb{R}^{2 \times 600}$ using LLE, LNEG and PCA. Figures. 6(b), (c), and (d) depict the 2-dimensional embeddings $\mathbf{\Psi}$ provided by LLE, LNEG, and PCA, respectively. Similarly, Figure 7 illustrates the resulting embeddings when $\mathbf{Z}_2$ is sampled from a nonlinear sphere. In

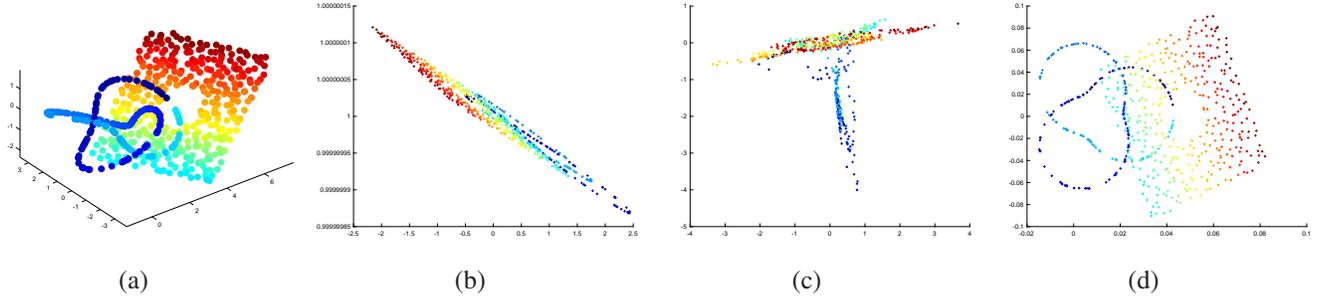(a)       (b)       (c)       (d)

Fig. 6: Embedding results of two manifolds: a linear hyperplane with hole and a trefoil (a) visualization of two manifolds; and $\{\psi_i\}_{i=1}^N$ obtained from (b) LLE with $K = 40$; and, (c) LNEG with $K = 40$ and $P = 2$; (d) PCA.
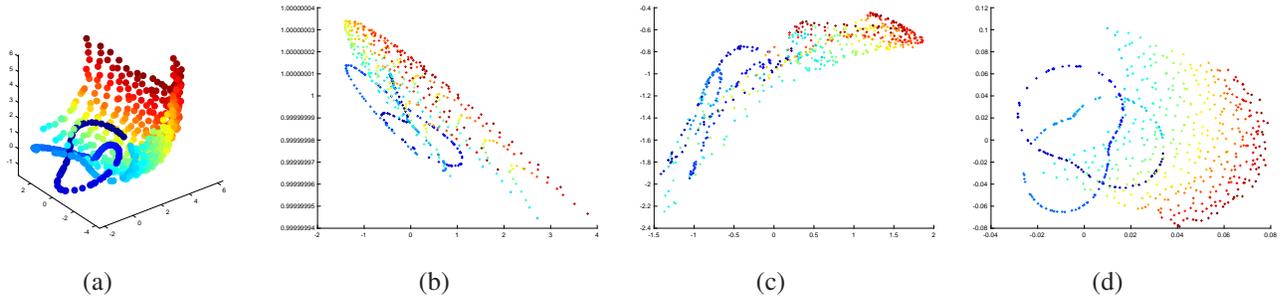


(a)       (b)       (c)       (d)

Fig. 7: Embedding results of two manifolds: a nonlinear sphere and a trefoil (a) visualization of two manifolds; and $\{\mathbf{y}_i\}_{i=1}^N$ obtained from (b) LLE with $K = 40$; (c) LNEG with $K = 40$ and $P = 3$; and (d) PCA.

| $K$ neighbours | Plane-hole-trefoil | | | Sphere-trefoil | | |
|---|---|---|---|---|---|---|
| | **LLE** | **LNE** | **LNEG** | **LLE** | **LNE** | **LNEG** |
| 5 | 0.25 | 0.18 | 0.18 | 0.29 | 0.21 | 0.17 |
| 10 | 0.44 | 0.21 | 0.18 | 0.39 | 0.27 | 0.16 |
| 20 | 0.15 | 0.13 | 0.17 | 0.48 | 0.26 | 0.14 |
| 30 | 0.21 | 0.20 | 0.17 | 0.46 | 0.28 | 0.19 |
| 40 | 0.36 | 0.20 | 0.17 | 0.39 | 0.21 | 0.20 |
| **PCA** | 0.49 | | | 0.43 | | |

TABLE V: Clustering error rate on low-dimensional representations obtained from: LLE, LNE, LNEG and PCA.

both cases, the nonlinear methods result in embeddings that separate the two manifolds. To further assess the performance, $K$-means is carried out on the resulting $\boldsymbol{\Psi}$ [16]. Table V shows the clustering error when running $K$-means on the low-dimensional embeddings given by PCA, LLE, LNE and LNEG, across different values of $k$. The proposed approaches provide embeddings that enhance separability of the two manifolds, resulting in lower clustering error compared to LLE and PCA. In addition, greater performance gain is observed when both manifolds are nonlinear, as in the case of Figure 7. The graph regularized method performs slightly better than that without regularization.

## VI. CONCLUSIONS

This paper introduced a general framework for nonlinear dimensionality reduction over graphs. By leveraging nonlinear relationships between data, low-dimensional representations were obtained to preserve these nonlinear correlations. Graph regularization was employed to account for additional prior knowledge when seeking the low-dimensional representations. An efficient algorithm that admits closed-form solution was developed along with a multi-kernel based algorithm that can handle settings where the nonlinear relationship between data is unknown. Furthermore, pertinent generalizations of the proposed schemes were provided. Several tests were conducted on simulated and real data to demonstrate the effectiveness of the proposed approaches. To broaden the scope of this study, several intriguing directions open up: a) online implementations that can handle streaming data; and b) generalizations to cope with large-scale graphs and high-dimensional datasets.

Consider the objective function of (3), and define $\mathbf{B} := \boldsymbol{\Psi}^\top \boldsymbol{\Psi}$. Then (3) can be rewritten as

$$\min_{\mathbf{B}:\text{rank}(\mathbf{B})=d} \|\mathbf{Y}^\top \mathbf{Y} - \mathbf{B}\|_F^2 \qquad (31)$$

where the $\text{rank}(\mathbf{B}) = d$ constraint comes from the fact that $\mathbf{B} = \boldsymbol{\Psi}^\top \boldsymbol{\Psi}$ and $\boldsymbol{\Psi}$ is a $d \times N$ matrix with $d \leq N$. The optimal solution $\mathbf{B}^*$ of (31) is given by the $d$ leading singular values and corresponding singular vectors of $\mathbf{Y}^\top \mathbf{Y}$ [22]. Since $\mathbf{Y} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$ we have $\mathbf{Y}^\top \mathbf{Y} = \mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^\top$, and consequently

$$\mathbf{B}^* = \mathbf{V}_d \boldsymbol{\Sigma}_d^2 \mathbf{V}_d^\top \qquad (32)$$

where $\mathbf{V}_d$ is a sub-matrix of $\mathbf{V}$ containing the $d$ singular vectors corresponding to the leading $d$ eigenvalues. It follows from (32) and $\mathbf{B} = \boldsymbol{\Psi}^\top \boldsymbol{\Psi}$ that

$$\boldsymbol{\Psi} = \boldsymbol{\Sigma}_d \mathbf{V}_d^\top \qquad (33)$$

which is the low-dimensional representation matrix provided by dual PCA [cf. (2)]. To complete the proof, just recall that

$$\boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \boldsymbol{\Sigma}_d^2 = \boldsymbol{\Lambda}_d \qquad (34)$$

where $\boldsymbol{\Lambda}_d$ contains the leading $d$ eigenvalues of $\mathbf{B}^*$.

Here, we will show that having found $\boldsymbol{\Psi}$, the coefficients $\{\theta_q\}$ in (15) can be obtained as in (17). Specifically, when the $\boldsymbol{\Psi}$ is available, $\boldsymbol{\theta}$ can be obtained by

$$\min_{\boldsymbol{\theta}} \quad -\text{tr}(\boldsymbol{\Psi}(\sum_{q=1}^{Q} \theta_q \mathbf{K}_y^{(q)})\boldsymbol{\Psi}^\top)$$
$$\text{s.t.} \ \ \|\boldsymbol{\theta}\|_2^2 \leq 1, \quad \boldsymbol{\theta} \geq \mathbf{0}. \qquad (35)$$

The Lagrangian of (35) is

$$\mathcal{L}(\boldsymbol{\theta},\lambda) = \text{tr}(\boldsymbol{\Psi}(\sum_{q=1}^{Q} \theta_q \mathbf{K}_y^{(q)})\boldsymbol{\Psi}^\top) + \lambda(\boldsymbol{\theta}^\top \boldsymbol{\theta} - 1). \qquad (36)$$

where $\lambda > 0$ is the Lagrange multiplier. Taking the gradient of $\mathcal{L}(\boldsymbol{\theta},\lambda)$ with respect to $\theta_q$ and equating it to zero we have

$$-\text{tr}(\boldsymbol{\Psi}\mathbf{K}_y^{(q)}\boldsymbol{\Psi}^\top) + \lambda\theta_q = 0, \quad \forall q = 1, \ldots, Q \qquad (37)$$

which yields

$$\theta_q = \frac{1}{\lambda}\text{tr}(\boldsymbol{\Psi}\mathbf{K}_y^{(q)}\boldsymbol{\Psi}^\top). \qquad (38)$$

Taking the gradient of $\mathcal{L}(\boldsymbol{\theta},\lambda)$ with respect to $\lambda$ and setting it to 0 we obtain

$$\sum_{q=1}^{Q} \theta_q^2 = 1. \qquad (39)$$

Substituting (38) into (39), we arrive at

$$\lambda = \sqrt{\sum_{q=1}^{Q}(\text{tr}(\boldsymbol{\Psi}\mathbf{K}_y^{(q)}\boldsymbol{\Psi}^\top))^2}. \qquad (40)$$

Combining (38) with (40) leads to (17).

# REFERENCES

[1] F. R. Bach, G. R. Lanckriet, and M. I. Jordan, "Multiple kernel learning, conic duality, and the SMO algorithm," in *Proc. Intl. Conf. on Machine Learning*, New York, USA, 2004, pp. 6–13.

[2] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural Computation*, vol. 15, no. 6, pp. 1373–1396, 2003.

[3] K. Diaz-Chito, A. Hernández-Sabaté, and A. M. López, "A reduced feature set for driver head pose estimation," *Applied Soft Computing*, vol. 45, pp. 98–107, 2016.

[4] E. Elhamifar and R. Vidal, "Sparse manifold clustering and embedding," in *Advances in Neural Information Processing Systems*, Granada, Spain, 2011, pp. 55–63.

[5] A. Ghodsi, "Dimensionality reduction -A short tutorial," *Department of Statistics and Actuarial Science, Univ. of Waterloo, Ontario, Canada*, vol. 37, p. 38, 2006.

[6] J. Ham, D. D. Lee, S. Mika, and B. Schölkopf, "A kernel view of the dimensionality reduction of manifolds," in *Proc. Intl. Conf. on Machine Learning*. Alberta, Canada: ACM, Jul. 2004, p. 47.

[7] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A K-means clustering algorithm," *Journal of the Royal Statistical Society*, vol. 28, no. 1, pp. 100–108, Jan. 1979.

[8] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2009.

[9] J. J. Hull, "A database for handwritten text recognition research," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 16, no. 5, pp. 550–554, 1994.

[10] B. Jiang, C. Ding, B. Luo, and J. Tang, "Graph-laplacian PCA: Closed-form solution and robustness," in *Proc. Intl. Conf. on Computer Vision and Pattern Recognition*, Columbus, Ohio, USA, June 2013, pp. 3492–3498.

[11] T. Jin, J. Yu, J. You, K. Zeng, C. Li, and Z. Yu, "Low-rank matrix factorization with multiple hypergraph regularizer," *Pattern Recognition*, vol. 48, no. 3, pp. 1011–1022, Mar. 2015.

[12] I. Jolliffe, *Principal Component Analysis*. Wiley Online Library, 2002.

[13] M. Kivelä, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter, "Multilayer networks," *Journal of Complex Networks*, vol. 2, no. 3, pp. 203–271, 2014.

[14] R. I. Kondor and J. Lafferty, "Diffusion kernels on graphs and other discrete structures," Sydney, Australia, Jul. 2002, pp. 315–322.

[15] J. B. Kruskal and M. Wish, *Multidimensional Scaling*. Sage, 1978, vol. 11.

[16] S. Lloyd, "Least-squares quantization in PCM," *IEEE Trans. Info. Theory*, vol. 28, no. 2, pp. 129–137, 1982.

[17] MATLAB, *version 9.1.0 (R2016b)*. Natick, Massachusetts: The Math-Works Inc., 2016.

[18] S. A. Nene, S. K. Nayar, H. Murase *et al.*, "Columbia object image library (coil-20)," 1996.

[19] D. Romero, V. N. Ioannidis, and G. B. Giannakis, "Kernel-based Reconstruction and Kalman Filtering of Space-time Functions on Dynamic Graphs," *IEEE Journal on Special Topics in Signal Processing*, vol. 11, no. 6, pp. 856 – 869, Sep.

[20] D. Romero, M. Ma, and G. B. Giannakis, "Kernel-based reconstruction of graph signals," *IEEE Transactions on Signal Processing*, vol. 65, no. 3, pp. 764–778, Feb. 2017.

[21] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, Dec. 2000.

[22] Y. Saad, *Numerical Methods for Large Eigenvalue Problems*. Manchester University Press, 1992.

[23] B. Schölkopf, A. Smola, and K.-R. Müller, "Kernel principal component analysis," in *Proc. Intl. Conf. on Artificial Neural Networks*, Lausanne, Switzerland, Oct. 1997, pp. 583–588.

[24] N. Shahid, N. Perraudin, V. Kalofolias, G. Puy, and P. Vandergheynst, "Fast robust PCA on graphs," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 740–756, Feb. 2016.

[25] F. Shang, L. Jiao, and F. Wang, "Graph dual regularization non-negative matrix factorization for co-clustering," *Pattern Recognition*, vol. 45, no. 6, pp. 2237–2250, 2012.

[26] Y. Shen, B. Baingana, and G. B. Giannakis, "Kernel-based structural equation models for topology identification of directed networks," *IEEE Trans. Sig. Proc.*, vol. 65, no. 10, pp. 2503–2516, May 2017.

[27] Y. Shen, P. A. Traganitis, and G. B. Giannakis, "Nonlinear dimensionality reduction on graphs," in *Proc. of CAMSAP*, Dutch Antilles, Dec. 2017.

[28] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 83–98, May 2013.

[29] A. J. Smola and R. I. Kondor, "Kernels and regularization on graphs," in *Learning Theory and Kernel Machines*. Springer, 2003, pp. 144–158.

[30] J. A. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural Processing Letters*, vol. 9, no. 3, pp. 293–300, Jun. 1999.

[31] J. B. Tenenbaum, V. d. Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000. [Online]. Available: http://science.sciencemag.org/content/290/5500/2319

[32] P. A. Traganitis, Y. Shen, and G. B. Giannakis, "Topology inference of multilayer networks," in *Intl. Workshop on Network Science for Comms.*, Atlanta, GA, May 2017.

[33] J. J.-Y. Wang, H. Bensmail, and X. Gao, "Multiple graph regularized nonnegative matrix factorization," *Pattern Recognition*, vol. 46, no. 10, pp. 2840–2847, 2013.