

Neural Processing Letters manuscript No.
(will be inserted by the editor)

Local Tangent Space Discriminant Analysis

Yang Zhou · Shiliang Sun

Received: date / Accepted: date

Abstract We propose a novel supervised dimensionality reduction method named *local Tangent Space Discriminant analysis* (TSD) which is capable of utilizing the geometrical information from tangent spaces. The proposed method aims to seek an embedding space where the local manifold structure of the data belonging to the same class is preserved as much as possible, and the marginal data points with different class labels are better separated. Moreover, TSD has an analytic form of the solution and can be naturally extended to non-linear dimensionality reduction through the kernel trick. Experimental results on multiple real-world data sets demonstrate the effectiveness of the proposed method.

Keywords Dimensionality reduction · Supervised learning · Manifold learning · Tangent space

1 Introduction

Dimensionality reduction is a learning task that aims to find a low-dimensional representation of high-dimensional data, while preserving data information as much as possible. Processing data in the low-dimensional space can reduce computational cost and suppress noises. Provided that dimensionality reduction is performed appropriately, the discovered low-dimensional representations of data will benefit subsequent tasks, e.g., classification, clustering, data visualization.

Yang Zhou

Shanghai Key Laboratory of Multidimensional Information Processing, Department of Computer Science and Technology, East China Normal University, 500 Dongchuan Road, Shanghai 200241, P. R. China

Shiliang Sun (Corresponding Author)

Shanghai Key Laboratory of Multidimensional Information Processing, Department of Computer Science and Technology, East China Normal University, 500 Dongchuan Road, Shanghai 200241, P. R. China

Tel.: +86-21-54345183

Fax: +86-21-54345119

E-mail: slsun@cs.ecnu.edu.cn, shiliangsun@gmail.com

PCA, as an unsupervised dimensionality reduction method, seeks a set of orthogonal projection directions along which the sum of variances of data is maximized. Some other popular unsupervised methods are geometrically motivated, which aim to discover the geometrical structure of the underlying manifold, such as Laplacian Eigenmaps [2], Hessian Eigenmaps [5], Locally Linear Embedding [9], Locality Preserving Projections [7], and Local Tangent Space Alignment [17], etc. Although unsupervised approaches can reveal the underlying data manifold, they may not be the best choices for some learning scenarios because they are not able to utilize the discriminative information from data labels.

LDA is a supervised dimensionality reduction method. It finds a subspace in which data points from different classes are projected far away from each other, while those belonging to the same class are projected as close as possible. However, LDA tends to get undesirable results when data are multimodal [6] or are mainly characterized by their variances. The reason why this happens lies in the assumption adopted by LDA that data points belonging to each class are generated from the multivariate Gaussian distributions with the same covariance matrix but different means. This assumption is invalid in dealing with the data formed by several separate clusters or those living on an underlying manifold.

To solve this problem, Subclass Discriminant Analysis [18] approximates the potential distribution of data with a mixture of Gaussian distributions. More specifically, it first divides each class into a set of subclasses through clustering and then performs LDA on the divided data. Another way to overcome the drawback of LDA is preserving the data structure locally. Marginal Fisher Analysis (MFA) [15] aims to gather the nearby examples of the same class, and separate the marginal examples belonging to different classes. Locality Sensitive Discriminant Analysis (LSDA) [3] maps data points into a subspace where the examples with the same label at each local area are close, while the nearby examples from different classes are apart from each other. Local Fisher Discriminant Analysis (LFDA) [12] also focuses on discovering the local data structure. It is equivalent to operate LDA in the local scope around each example. In fact, these local structure oriented methods actually fall into the same graph Laplacian based framework. All of them employ the Laplacian matrix to preserve the local geometry of the data manifold. However, this framework fails to discover the local manifold information from tangent spaces which could be very useful and can enhance the performance of dimensionality reduction in some situations [11, 13].

In this paper, we present a novel supervised dimensionality reduction method named *local Tangent Space Discriminant analysis* (TSD). Unlike previous approaches using the graph Laplacian to discover the data manifold, our method uses the first-order Taylor expansion to represent the geometry of the local area around each data point. This strategy provides us with a natural way to utilize the information from tangent spaces. Then we seek a linear transformation to preserve the local manifold structure of the data belonging to the same class as much as possible, while maximizing the marginal data points with different class labels. As a result, the geometrical information from tangent spaces can be readily incorporated into the proposed method to improve the performance of dimensionality reduction. Moreover, the objective function of our method can be optimized analytically by solving a generalized eigenvalue problem. This also leads to a natural extension for non-linear dimensionality reduction through the kernel trick.

The rest of this paper is organized as follows. We briefly review some related work including MFA, LSDA and LFDA, and show how these methods can be considered in the same framework in Section 2. Then the *local Tangent Space Discriminant analysis* (TSD) along with its kernelization are introduced in Section 3. Section 4 discusses the connection and difference between the proposed method and related work. In Section 5, experimental results are presented. Finally, we give concluding remarks and discuss some future work in Section 6.

2 Related Work

Many dimensionality reduction methods have been proposed in recent years. Although they have different names and are derived from various motivations, a large portion of them fall into an unified graph Laplacian based framework. For example, Yan et al. [15] proposed a dimensionality reduction framework called Graph Embedding which can group together many popular dimensionality reduction approaches into a general formulation. In this section, we first introduce the Graph Embedding framework, and then briefly review some related work and show how they fall into the same framework. Finally, we provide a brief summary to discuss the strength and weakness of this framework.

2.1 Graph Embedding

Given a data set X consisting of n examples and labels, $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ denotes a d -dimensional example, $y_i \in \{1, 2, \dots, C\}$ denotes the class label corresponding to \mathbf{x}_i , and C is the total number of classes. The relationship between each example can be easily characterized by a undirected weighted graph $G = \{X, W\}$, where each example is served as the vertex of G and a symmetric weight matrix $W \in \mathbb{R}^{n \times n}$ records the weight on the edge of each pair of vertices. W measures the similarity between each example, and its characteristic varies as the criterion of similarity changes. Generally, if two examples \mathbf{x}_i and \mathbf{x}_j are “close”, the corresponding weight W_{ij} is large, whereas if they are “far away”, then the W_{ij} is small. Provided a certain W , the intrinsic geometry of graph G can be represented by the Laplacian matrix [4], which is defined as

$$L = D - W, \quad (1)$$

where D is a diagonal matrix with the i -th diagonal element being $D_{ii} = \sum_{j \neq i} W_{ij}$. The Laplacian matrix is capable of representing certain geometry of data according to a specific weight matrix, and thus can be used for dimensionality reduction.

To find a good low-dimensional embedding $\mathbf{b} = (b_{x_1}, b_{x_2}, \dots, b_{x_n})^\top$ from high-dimensional data, we have to preserve the intrinsic geometry of the original data as much as possible. Therefore, it is natural to seek the embedding preserving the most information from G in each dimension. This graph-preserving criterion can be formulated as follows [15]:

$$\begin{aligned} \mathbf{b}^* &= \arg \min_{\mathbf{b}^\top L \mathbf{b} = a} \sum_{i \neq j} W_{ij} \|\mathbf{b}_i - \mathbf{b}_j\|^2 \\ &= \arg \min_{\mathbf{b}^\top L \mathbf{b} = a} \mathbf{b}^\top L \mathbf{b}, \end{aligned} \quad (2)$$

where L is the Laplacian matrix of G defined in (1), L_p is the penalty constraint matrix, and a is a constant defined to avoid a trivial solution of the objective function. Note that L_p can have multiple forms, which is usually a diagonal matrix or the Laplacian matrix of a penalty graph $G_p = \{X, W^p\}$ constructed by the same vertices X yet a different weight matrix W^p .

In this paper, we mainly focus on the linear dimensionality reduction, and the embedding of each example \mathbf{x}_i is computed as $b_{\mathbf{x}_i} = \mathbf{t}^\top \mathbf{x}_i$ where \mathbf{t} is a projection vector. In this case, (2) becomes:

$$\mathbf{t}^* = \arg \min_{\mathbf{t}^\top X L^p X^\top \mathbf{t} = a} = \mathbf{t}^\top X L X^\top \mathbf{t}. \quad (3)$$

The objective function (3) can be converted to a generalized eigenvalue problem:

$$X L X^\top \mathbf{t} = \lambda X L^p X^\top \mathbf{t}. \quad (4)$$

whose solution can be easily given by the eigenvector with respect to the smallest eigenvalue.

Given the above results, many local structure oriented dimensionality reduction approaches, which closely related to our proposed method, can be grouped into the unified framework. Next, we mainly discuss three of them including MFA, LSDA and LFDA.

2.2 Marginal Fisher Analysis

Marginal Fisher Analysis (MFA) is a dimensionality reduction approach that is directly derived from the Graph Embedding framework [15]. The main idea of MFA is to preserve the intraclass compactness represented by an intrinsic graph under the constraint that the interclass separability characterized by a penalty graph should be kept. For this purpose, each example is connected to its k_1 -nearest neighbors belonging to the same class in the intrinsic graph, and the penalty graph is built by connecting k_2 -nearest pairs of the marginal point in different classes.

For each example, let $N_{k_1}(i)$ be the set of the k_1 -nearest neighbors of \mathbf{x}_i in the same class, and $P_{k_2}(i)$ indicates the set of the k_2 -nearest pairs among the set $\{(i, j), y_i \neq y_j\}$. The intraclass compactness is formulated by a local within-class scatter matrix \bar{S}_w :

$$\begin{aligned} \bar{S}_w &= \sum_i \sum_{i \in N_{k_1}(j) \text{ or } j \in N_{k_1}(i)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top \\ &= 2X(\bar{D} - \bar{W})X^\top \\ &= 2X\bar{L}X^\top, \end{aligned}$$

where $\bar{L} = \bar{D} - \bar{W}$ is the Laplacian matrix, \bar{D} is a diagonal matrix with the i -th diagonal element being $\bar{D}_{ii} = \sum_{j \neq i} \bar{W}_{ij}$, and the weight matrix \bar{W} is defined as follows:

$$\bar{W}_{ij} = \begin{cases} 1 & \text{if } i \in N_{k_1}(j) \text{ or } j \in N_{k_1}(i) \\ 0 & \text{else.} \end{cases}$$

Similarly, the interclass separability can be characterized by a local between-class scatter matrix \bar{S}_b :

$$\begin{aligned}\bar{S}_b &= \sum_i \sum_{(i,j) \in P_{k_2}(i) \text{ or } (i,j) \in P_{k_2}(j)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top \\ &= 2X(\bar{D}^p - \bar{W}^p)X^\top \\ &= 2X\bar{L}^pX^\top,\end{aligned}$$

where $\bar{L}^p = \bar{D}^p - \bar{W}^p$ is the Laplacian matrix, \bar{D}^p is a diagonal matrix with the i -th diagonal element being $\bar{D}_{ii}^p = \sum_{j \neq i} \bar{W}_{ij}^p$ and the weight matrix \bar{W}^p is defined as follows:

$$\bar{W}_{ij}^p = \begin{cases} 1 & \text{if } (i,j) \in P_{k_2}(i) \text{ or } (i,j) \in P_{k_2}(j) \\ 0 & \text{else.} \end{cases}$$

Following the graph-preserving criterion presented in (2), the objective function of MFA can be written as follows:

$$\begin{aligned}\mathbf{t}^{MFA} &= \arg \min_{\mathbf{t}^\top \bar{S}_b \mathbf{t} = a} \mathbf{t}^\top \bar{S}_w \mathbf{t} \\ &= \arg \min_{\mathbf{t}^\top X\bar{L}^pX^\top \mathbf{t} = a} \mathbf{t}^\top X\bar{L}X^\top \mathbf{t}.\end{aligned}\quad (5)$$

2.3 Locality Sensitive Discriminant Analysis

Another related method is Locality Sensitive Discriminant Analysis (LSDA) [3] which assumes that data live on or close to a manifold. It aims to preserve the local geometrical structure of the manifold while maximizing the local margin between different classes.

LSDA seeks a linear projection \mathbf{t} optimizing

$$\begin{aligned}\min & \sum_{ij} \bar{W}_{ij}^w (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j)^2, \\ \max & \sum_{ij} \bar{W}_{ij}^b (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j)^2\end{aligned}$$

under the constraint that $\mathbf{t}^\top X\bar{D}^wX^\top \mathbf{t} = 1$. Let $N^w(i)$ be the set of the k -nearest neighbors of \mathbf{x}_i sharing the same label y_i , and $N^b(i)$ be the set of the k -nearest neighbors of \mathbf{x}_i having the labels different from y_i . Then the weight matrices \bar{W}^w and \bar{W}^b are defined as:

$$\begin{aligned}\bar{W}_{ij}^w &= \begin{cases} 1 & \text{if } i \in N^w(j) \text{ or } j \in N^w(i) \\ 0 & \text{else,} \end{cases} \\ \bar{W}_{ij}^b &= \begin{cases} 1 & \text{if } i \in N^b(j) \text{ or } j \in N^b(i) \\ 0 & \text{else.} \end{cases}\end{aligned}$$

The objective function of LSDA described above can be formulated as follows:

$$\begin{aligned}\mathbf{t}^{LSDA} &= \arg \max_{\mathbf{t}^\top X\bar{D}^wX^\top \mathbf{t} = 1} \mathbf{t}^\top X(\alpha\bar{L}^b + (1-\alpha)\bar{W}^w)X^\top \mathbf{t} \\ &= \arg \min_{\mathbf{t}^\top X\bar{D}^wX^\top \mathbf{t} = 1} \mathbf{t}^\top X((\alpha-1)\bar{W}^w - \alpha\bar{L}^b)X^\top \mathbf{t},\end{aligned}\quad (6)$$

where α is a trade-off parameter, \bar{L}^b is the Laplacian matrix constructed by \bar{W}^b , and \bar{D}^w is a diagonal matrix with the i -th diagonal element being $\bar{D}_{ii}^w = \sum_{j \neq i} \bar{W}_{ij}^w$. LSDA follows the framework defined in (3) and the solution \mathbf{t}^{LSDA} is given by (4) with $L = (\alpha - 1)\bar{W}^w - \alpha\bar{L}^b$ and $L^p = \bar{D}^w$.

2.4 Local Fisher Discriminant Analysis

Local Fisher Discriminant Analysis (LFDA) [12] combines the ideas of LDA and LPP [7] to overcome the problem that LDA [6] can not appropriately handle the data with multimodality. More specifically, it evaluates the levels of the between-class scatter and the within-class scatter in a local manner, and tries to attain the local between-class separation and the local within-class structure preservation at the same time [12].

Let \tilde{S}_w and \tilde{S}_b be the local within-class scatter matrix and the local between-class scatter matrix defined by

$$\begin{aligned}\tilde{S}_w &= \frac{1}{2} \sum_{ij} \tilde{W}_{ij}^w (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top = X \tilde{L}^w X^\top, \\ \tilde{S}_b &= \frac{1}{2} \sum_{ij} \tilde{W}_{ij}^b (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top = X \tilde{L}^b X^\top,\end{aligned}$$

where \tilde{L}^w and \tilde{L}^b are the Laplacian matrices constructed by the weight matrices \tilde{W}_{ij}^w and \tilde{W}_{ij}^b with

$$\begin{aligned}\tilde{W}_{ij}^w &= \begin{cases} A_{ij}/n_c & \text{if } y_i = y_j \\ 0 & \text{if } y_i \neq y_j, \end{cases} \\ \tilde{W}_{ij}^b &= \begin{cases} A_{ij}(1/n - 1/n_c) & \text{if } y_i = y_j \\ 1/n & \text{if } y_i \neq y_j. \end{cases}\end{aligned}$$

n_c denotes the number of examples from the c -th class, and A_{ij} is a weight that indicates the similarity between \mathbf{x}_i and \mathbf{x}_j , whose definition is given as follows:

$$A_{ij} = \begin{cases} \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma_i \sigma_j}) & \text{if } i \in N_k(j) \text{ or } j \in N_k(i) \\ 0 & \text{else,} \end{cases}$$

where σ_i is set to be the distance between \mathbf{x}_i and its k -th nearest neighbor.

LFDA seeks a linear projection so that \tilde{S}_w is minimized and \tilde{S}_b is maximized. Essentially, this strategy is equivalent to find a projection which fits the Fisher criterion in the local area around each example. The optimization problem of LFDA is given as follows:

$$\begin{aligned}\mathbf{t}^{LFDA} &= \arg \max_{\mathbf{t}} \frac{\mathbf{t}^\top \tilde{S}_b \mathbf{t}}{\mathbf{t}^\top \tilde{S}_w \mathbf{t}} = \arg \max_{\mathbf{t}} \frac{\mathbf{t}^\top X \tilde{L}^w X^\top \mathbf{t}}{\mathbf{t}^\top X \tilde{L}^b X^\top \mathbf{t}} \\ &= \arg \min_{\mathbf{t}^\top X \tilde{L}^b X^\top \mathbf{t} = 1} = \mathbf{t}^\top X \tilde{L} X^\top \mathbf{t}.\end{aligned}\quad (7)$$

According to (7), it is easy to find that LFDA also falls into the Graph Embedding framework defined in (3) with $L = \tilde{L}^w$, $L^p = \tilde{L}^b$ and $a = 1$.

2.5 A Brief Summary

With the above results, it is easy to find that all methods mentioned above can be considered in the same graph Laplacian based framework and the main difference among them only lies in the different graphs adopted by each method. For different methods, such graphs are constructed by certain weight matrices to incorporate specific neighborhood information of the data set. One important merit of this framework is that it not only takes advantage of the facility of the Laplacian matrix to preserve the local geometry of the data manifold, but also benefits from the elegant formulation which can be easily optimized through the generalized eigenvalue decomposition.

Although the graph Laplacian provides us with a powerful and flexible tool to discover the underlying data manifold, it fails to discover the local geometrical structure from tangent spaces, and thus may lose much useful information whose effectiveness has been shown in many applications especially in the handwritten digit recognition [11]. Moreover, the graph Laplacian may fail to capture meaningful manifold structures, when data are sparsely distributed in the original space. In this case, the graph Laplacian constructed by sparsely distributed data in the high-dimensional space may not be able to discover the correct underlying manifold, since it can hardly connect sparse data points into a smooth manifold. On the other hand, the tangent spaces of the underlying manifold, which are low-dimensional in nature, can reflect the manifold structure in each local area. This implies that tangent spaces are very useful for learning the data manifold. Then how to develop a dimensionality reduction algorithm which is capable of combining the flexibility of the graph Laplacian with the utility of tangent spaces? To solve this problem, we present our algorithm which can readily use the structural information from tangent spaces for supervised dimensionality reduction.

3 Local Tangent Space Discriminant Analysis

In this section, we present the *local Tangent Space Discriminant analysis* (TSD) algorithm and its non-linear extension. As a supervised dimensionality reduction method, TSD aims to seek an embedding space where the local manifold structure of the data belonging to the same class is preserved as much as possible, and the marginal data points with different class labels are better separated. Compared with the methods discussed in Section 2, the key advantage of our algorithm is that it is capable of capturing the local manifold structure from tangent spaces without losing the analytic form of the solution.

3.1 Preliminaries

To begin with, we briefly introduce the concepts of the tangent space and tangent vector. In differential geometry, one can attach to every point \mathbf{x} of a differentiable manifold \mathcal{M} a tangent space $T_{\mathbf{x}}\mathcal{M}$ in which every vector tangentially passes through \mathbf{x} . The elements of the tangent space are called tangent vectors at \mathbf{x} , which is a vector that is tangent to a curve or surface at \mathbf{x} (see Fig. 1 for the

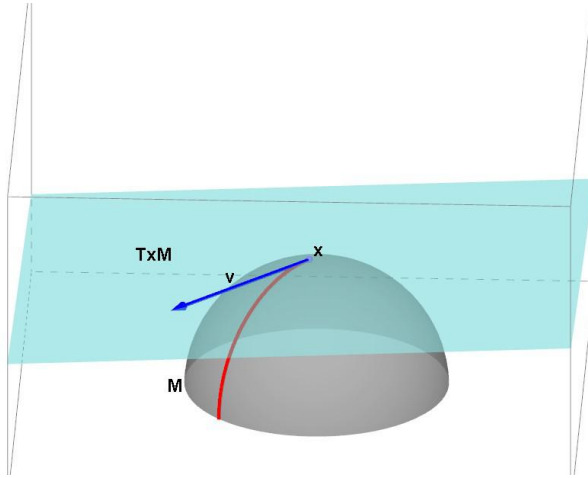


Fig. 1 The tangent space $T_x\mathcal{M}$ and a tangent vector $v \in T_x\mathcal{M}$, along a curve travelling through $x \in \mathcal{M}$.

illustration). All the tangent spaces of a connected manifold have the same dimension, equal to the dimension of the manifold. In practice, if the manifold is smooth enough, the subspace constructed by performing PCA on the neighborhood of x can be a good approximation of the tangent space at x [14], since the nearby data points of x can be viewed as approximately lying in a subspace which is tangent to the data manifold. Once tangent spaces and tangent vectors have been introduced, they can serve to characterize a differentiable curve on the manifold whose derivative at any point is equal to the tangent vector attached to that point. This is a crucial property that plays a key role in deriving our TSD algorithm.

In recent years, some tangent space based dimensionality reduction methods have been proposed by using the above property. They learn the data manifold by estimating a function whose value can serve as a low-dimensional representation of the manifold. Tangent Space Intrinsic Manifold Regularization (TSIMR) [13] estimates a local linear function on the manifold which has constant manifold derivatives. Parallel Vector Field Embedding (PFE) [8] represents a function along the manifold from the perspective of vector fields and requires the vector field at each data point to be as parallel as possible. Although they are effective to preserve the manifold geometry, these tangent space based methods are unsupervised in nature, which have no ability to utilize the discriminant information of class labels. Therefore, they are not optimal for the supervised case. To solve this problem, we propose the TSD algorithm which partly shares the same spirit with TSIMR and PFE but is optimal for supervised dimensionality reduction.

3.2 The Algorithm

Suppose that data are sampled from an m -dimensional smooth manifold \mathcal{M} in a d -dimensional space. Let $\mathcal{T}_z\mathcal{M}$ denotes the tangent space attached to \mathbf{z} , where $\mathbf{z} \in \mathcal{M}$ is a fixed data point on the \mathcal{M} . Motivated by Tangent Space Intrinsic Manifold Regularization (TSIMR) [13], we represent the local manifold structure of data by means of tangent spaces. According to the first-order Taylor expansion at \mathbf{z} , any function f defined on the manifold \mathcal{M} can be expressed as:

$$f(\mathbf{x}) = f(\mathbf{z}) + \mathbf{w}_z^\top \mathbf{u}_z(\mathbf{x}) + O(\|\mathbf{x} - \mathbf{z}\|^2),$$

where $\mathbf{x} \in \mathbb{R}^d$ is a d -dimensional data point and $\mathbf{u}_z(\mathbf{x}) = T_z^\top(\mathbf{x} - \mathbf{z})$ is an m -dimensional tangent vector which gives the m -dimensional representation of \mathbf{x} in $\mathcal{T}_z\mathcal{M}$. T_z is a $d \times m$ matrix formed by the orthonormal bases of $\mathcal{T}_z\mathcal{M}$, which can be estimated through local PCA, i.e., performing standard PCA on the neighborhood of \mathbf{z} . And \mathbf{w}_z is an m -dimensional vector representing the directional derivative of f at \mathbf{z} with respect to $\mathbf{u}_z(\mathbf{x})$ on the manifold \mathcal{M} .

In the scenario of dimensionality reduction, $f(\mathbf{x})$ denotes a one-dimensional embedding of \mathbf{x} . If there are two data points \mathbf{z} and \mathbf{z}' have a small Euclidean distance, by using the first-order Taylor expansion at \mathbf{z}' , the embedding $f(\mathbf{z})$ can be represented as:

$$f(\mathbf{z}) = f(\mathbf{z}') + \mathbf{w}_{z'}^\top T_{z'}^\top(\mathbf{z} - \mathbf{z}') + O(\|\mathbf{z} - \mathbf{z}'\|^2). \quad (8)$$

Suppose that the data can be well characterized by a linear function on the underlying manifold \mathcal{M} . Then we can omit the remainders in (8) because the second-order derivatives of f vanishes. Therefore, provided \mathbf{z} and \mathbf{z}' are close enough, any embedding $f(\mathbf{z})$ can be well approximated by a linear function as follows:

$$f(\mathbf{z}) \approx f(\mathbf{z}') + \mathbf{w}_{z'}^\top T_{z'}^\top(\mathbf{z} - \mathbf{z}'). \quad (9)$$

Based on the above results, we know that every data point in a local area should satisfies (9), which leads to a natural criterion of preserving the local manifold structure of data. Suppose that the training data include n examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ belonging to C classes where $\mathbf{x}_i \in \mathbb{R}^d$ is a d -dimensional example, and $y_i \in \{1, 2, \dots, C\}$ is the class label associated with the example \mathbf{x}_i . Consider a linear projection $\mathbf{t} \in \mathbb{R}^d$ which maps the data to a one-dimensional embedding. Then the embedding of \mathbf{x} can be expressed as $f(\mathbf{x}) = \mathbf{t}^\top \mathbf{x}$. We aim to find a projection \mathbf{t} to minimize the difference between the l.h.s and the r.h.s of (9) for every example and its neighbors belonging to the same class, and to better separate the marginal data points in different classes.

In order to minimize the difference between the l.h.s and the r.h.s of (9) for nearby intraclass data, we need to construct the within-class graph G^w . For the within-class graph G^w , if \mathbf{x}_i is among the k_1 -nearest neighbors of \mathbf{x}_j with $y_i = y_j$, an edge is added between \mathbf{x}_i and \mathbf{x}_j , and the elements of the weight matrix W^w are set to $W_{ij}^w = W_{ji}^w = 1$. Then we can formulate a within-class objective function as follows:

$$\min \sum_{ij} W_{ij}^w (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j - \mathbf{w}_{x_j}^\top T_{x_j}^\top(\mathbf{x}_i - \mathbf{x}_j))^2, \quad (10)$$

$$W_{ij}^w = \begin{cases} 1 & \text{if } i \in N_{k_1}(j) \text{ or } j \in N_{k_1}(i) \\ 0 & \text{else,} \end{cases} \quad (11)$$

where $N_{k_1}(i)$ denotes the set of the k -nearest neighbors of \mathbf{x}_i sharing the same label y_i , and the orthonormal base matrix $T_{\mathbf{x}_i}$ of the tangent space $\mathcal{T}_{\mathbf{x}_i}\mathcal{M}$ at each data point \mathbf{x}_i are computed by performing PCA on the k_1 -nearest neighborhood of \mathbf{x}_i .

To separate the marginal data points in different classes, we need to maximize and the distances of the embeddings of nearby interclass data points. To this end, we need to construct the between-class graph G^b . For the between-class graph G^b , if the pair (i, j) is among the k_2 -shortest pairs in the set $\{(i, j), y_i \neq y_j\}$, an edge is added between \mathbf{x}_i and \mathbf{x}_j , and the elements of the weight matrix W^b are set to $W_{ij}^b = 1$. Similar to the above objective function, we can write a between-class objective function as follows:

$$\max \sum_{ij} W_{ij}^b (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j)^2, \quad (12)$$

$$W_{ij}^b = \begin{cases} 1 & \text{if } (i, j) \in P_{k_2}(i) \text{ or } (i, j) \in P_{k_2}(j) \\ 0 & \text{else,} \end{cases} \quad (13)$$

where $P_{k_2}(i)$ indicates the set of the k_2 -nearest pairs among the set $\{(i, j), y_i \neq y_j\}$.

Note that the terms $\mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j)$ ($i, j = 1, \dots, n$) in (10) distinguish our strategy of preserving the local data structure from the graph Laplacian based one, where $\mathbf{w}_{\mathbf{x}_j}$ is a coefficient vector and should be optimized with \mathbf{t} simultaneously. These terms characterize how well two different examples \mathbf{x}_i and \mathbf{x}_j fit into the local linear approximation of f , which leads to an appropriate way to preserve the local intraclass geometry along the manifold \mathcal{M} . Therefore, our strategy can extract more geometrical information from the data than the graph Laplacian based one. Moreover, any valid weight matrix W^w , such as the one used in LFDA, can be used to preserve specific geometrical structure of the data manifold. This free-form property of the weight matrix is of great importance when we want to apply dimensionality reduction to various types of data.

The objective function (10) can be reformulated as a canonical matrix quadratic form as follows:

$$\begin{aligned} & \sum_{ij} W_{ij}^w (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j - \mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 \\ &= \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} XS_1X^\top & XS_2 \\ S_2^\top X^\top & S_3 \end{pmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top S \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}, \end{aligned} \quad (14)$$

where we have defined $\mathbf{w} = (\mathbf{w}_{\mathbf{x}_1}^\top, \mathbf{w}_{\mathbf{x}_2}^\top, \dots, \mathbf{w}_{\mathbf{x}_n}^\top)^\top$, $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ is the data matrix, and S is a $(d + mn) \times (d + mn)$ positive semi-definite matrix constructed by four blocks, i.e., XS_1X^\top , XS_2 , $S_2^\top X^\top$ and S_3 . For simplicity, we omit the detailed derivation of S here, which is available in the Appendix A.

Recall that $\mathbf{w}_{\mathbf{x}_i}$ is the directional derivative of f at \mathbf{x}_i . Note that the linear projection vector \mathbf{t} is under the influences of both the direction and the length of each $\mathbf{w}_{\mathbf{x}_i}$. To make within-class examples further compacted, we hope that the projection \mathbf{t} is more effected by $\mathbf{w}_{\mathbf{x}_i}$'s direction than its length. Therefore, it makes sense to regularize the length of $\mathbf{w}_{\mathbf{x}_i}$ ($i = 1, \dots, n$). This can be achieved by adding

a regularizer $\|\mathbf{t}\|^2 + \sum_i \|\mathbf{w}_{x_i}\|^2$ to (14). Define $\mathbf{f} = (\mathbf{t}^\top, \mathbf{w}^\top)^\top$, the optimization function turns out to be:

$$\begin{aligned} & \mathbf{f}^\top S \mathbf{f} + \gamma (\|\mathbf{t}\|^2 + \sum_i \|\mathbf{w}_{x_i}\|^2) \\ &= \mathbf{f}^\top S \mathbf{f} + \gamma \|\mathbf{f}\|^2 = \mathbf{f}^\top (S + \gamma I) \mathbf{f}, \end{aligned} \quad (15)$$

where $\gamma > 0$ is a trade-off parameter. In fact, the extra term $\gamma \|\mathbf{f}\|^2$ often refers to as the Tikhonov regularizer, which is commonly used with a small γ to keep matrices from being singular. However, the value of γ is crucial for our method, because it also controls the influence of \mathbf{w} 's length on the projection \mathbf{t} .

With simple algebraic formulation, the objective function (12) becomes

$$\begin{aligned} & \sum_{ij} W_{ij}^b (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \mathbf{x}_j)^2 \\ &= 2\mathbf{t}^\top X (D^b - W^b) X^\top \mathbf{t} = 2\mathbf{t}^\top X L^b X^\top \mathbf{t} \\ &= \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} 2XL^bX^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top S_b \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}, \end{aligned} \quad (16)$$

where $L^b = D^b - W^b$ is the Laplacian matrix, and D^b is a diagonal matrix with the i -th diagonal element being $D_{ii}^b = \sum_{j \neq i} W_{ij}^b$.

Finally, by integrating (15) and (16), the objective function of TSD can be written as follows:

$$\mathbf{f}^* = \arg \max_{\mathbf{f}} \frac{\mathbf{f}^\top S_b \mathbf{f}}{\mathbf{f}^\top (S + \gamma I) \mathbf{f}}. \quad (17)$$

The optimization of (17) can be achieved by solving a generalized eigenvalue problem:

$$S_b \mathbf{f} = \lambda (S + \gamma I) \mathbf{f} \quad (18)$$

whose solution can be easily given by the eigenvector with respect to the largest eigenvalue. In order to obtain a one-dimensional embedding of an example \mathbf{x} , we just use the first part of $\mathbf{f}^* = (\mathbf{t}^{*\top}, \mathbf{w}^{*\top})^\top$ and compute $b = \mathbf{t}^{*\top} \mathbf{x}$. Suppose that we want to project d -dimensional data into an r -dimensional subspace. Let $\mathbf{f}_1, \dots, \mathbf{f}_r$ be the solutions of (17) corresponding to the r largest eigenvalues $\lambda_1 > \dots > \lambda_r$. Then the r -dimensional embedding \mathbf{b} of \mathbf{x} is computed as follows:

$$\mathbf{b} = T^\top \mathbf{x}, \quad T = (\mathbf{t}_1, \dots, \mathbf{t}_r).$$

Algorithm 1 gives the pseudo-code for TSD. It is worth noting that although \mathbf{w}^* seems not to be used in computing the low-dimensional embeddings, as the parameter which is simultaneously optimized with \mathbf{t}^* , it exerts a crucial influence on the resultant transformation matrix T . This is means that both \mathbf{t}^* and \mathbf{w}^* determine the final results of TSD.

The main computational cost of TSD lies in building tangent spaces for n data points and solving the generalized eigenvalue problem (18). Our algorithm has a time complexity of $O((d^2m + m^2d) \times n)$ for the construction of n tangent spaces and $O(r^2 \times (d + mn))$ for finding r eigenvectors with respect to the r

Algorithm 1 TSD

Input: Labeled examples $\{\{\mathbf{x}_i, \mathbf{y}_i\} | \mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in \{1, 2, \dots, C\}\}_{i=1}^n$;
Dimensionality of embedding space m ($1 \leq m \leq d$);
Trade-off parameters γ ($\gamma > 0$).

Output: $d \times r$ transformation matrix T .

Construct the within-class graph G^w and the between-class graph G^b ;
Calculate the weight matrices W^w and W^b with (11) and (13);
for $i = 1$ **to** n **do**
 Construct $T_{\mathbf{x}_i}$ by performing PCA on the intraclass neighborhood of \mathbf{x}_i ;
end for
Compute the eigenvectors $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_r$ of (18) with respect to the top r eigenvalues;
 $T = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_r)$.

largest eigenvalues. For comparison, we also give the time complexities of some classical and related methods. The time complexity of PCA is $O(n^2d)$ and that of LDA is also $O(n^2d)$ [16]. As we have discussed in Section 2, MFA, LSDA, LFDA fall into the same framework with different graphs, which implies that they have the same computational cost. Since LDA also falls into the graph Laplacian based framework [15], their time complexities turn out to be $O(n^2d)$. The above analysis suggests that TSD is more time consuming compared with other methods. However, since local tangent spaces are estimated by local PCA, we can obtain at most $k_1 + 1$ meaningful orthonormal bases for each tangent space¹, where k_1 is the size of within-class neighborhood. This implies that the dimensionality m of the directional derivative $\mathbf{w}_{\mathbf{x}_i}$ ($i = 1, \dots, n$) is always less than $k_1 + 1$. In practice, k_1 is usually small to ensure the locality. This makes sure that m is actually a small constant. To conclude, the overall time complexity of TSD is $O((d^2m + m^2d) \times n + r^2 \times (d + mn))$. Since m is usually small, TSD has an acceptable computational cost.

3.3 Kernel TSD

TSD is a linear dimensionality reduction method. In this section, we propose Kernel TSD which can be performed in a Reproducing Kernel Hilbert Space (RKHS) for non-linear dimensionality reduction.

Consider a feature space \mathcal{F} induced by a non-linear mapping $\phi : \mathcal{X} \rightarrow \mathcal{F}$, where \mathcal{X} is an input domain. We can construct an RKHS $H_{\mathcal{K}}$ by defining a kernel function $\mathcal{K}(\cdot, \cdot)$ using the inner product operation $\langle \cdot, \cdot \rangle$, such that $\mathcal{K}(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$. Given a data set $\{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^n$, we can define the data matrix in the feature space \mathcal{F} as $\Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))$. Then one can use the orthogonal projection to decompose any projection vector $\mathbf{t} \in H_{\mathcal{K}}$ into a sum of two functions: one lying in the $span\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)\}$, and the other lying in the orthogonal complementary space. Therefore, there exist a set of coefficients α_i ($i = 1, 2, \dots, n$) satisfying

$$\mathbf{t} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) + \mathbf{v} = \Phi \boldsymbol{\alpha} + \mathbf{v}, \quad (19)$$

¹ That's because there are only $k_1 + 1$ examples as the inputs of local PCA.

where $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)^\top$ and $\langle \mathbf{v}, \phi(\mathbf{x}_i) \rangle = 0$ for all i . Note that although we set $\mathbf{f} = (\mathbf{t}^\top, \mathbf{w}^\top)^\top$ and optimize \mathbf{t} and \mathbf{w} together, we should estimate tangent spaces in \mathcal{F} through local Kernel PCA [10] rather than reparametrize \mathbf{w} like \mathbf{t} .

Let $T_{\mathbf{x}_i}^\phi$ be the matrix formed by the orthonormal bases of the tangent space attached to $\phi(\mathbf{x}_i)$. By replacing $T_{\mathbf{x}_i}$ with $T_{\mathbf{x}_i}^\phi$ ($i = 1, 2, \dots, n$) and substituting (19) into (15), the objective function (15) turns out to be:

$$\begin{aligned} & \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} XS_1X^\top & XS_2 \\ S_2^\top X^\top & S_3 \end{pmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix} + \gamma \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} KS_1K & KS_2 \\ S_2^\top K & S_3 \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix} + \gamma \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} K & \mathbf{0} \\ \mathbf{0} & \bar{I} \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} KS_1K + \gamma K & KS_2 \\ S_2^\top K & S_3 + \gamma \bar{I} \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top S^\phi \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}, \end{aligned}$$

where K is a kernel matrix with $K_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ and \bar{I} is an identity matrix sized $mn \times mn$. Similarly, the objective function (16) becomes

$$\begin{aligned} & \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} 2XL^bX^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top \begin{pmatrix} 2KL^bK & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}^\top S_b^\phi \begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{w} \end{pmatrix}. \end{aligned}$$

Note that due to $\langle \mathbf{v}, \phi(\mathbf{x}_i) \rangle = 0$ for all i , every term of \mathbf{v} vanishes from the above formulations. Finally, Kernel TSD can be converted to a generalized eigenvalue problem as follows:

$$S_b^\phi \boldsymbol{\varphi} = \lambda S^\phi \boldsymbol{\varphi}, \quad (20)$$

where we have defined $\boldsymbol{\varphi} = (\boldsymbol{\alpha}^\top, \mathbf{w}^\top)^\top$.

Given the eigenvectors $\boldsymbol{\varphi}_1, \dots, \boldsymbol{\varphi}_r$ with respect to the r largest eigenvalues of (20), the resultant transformation matrix can be written as $\Gamma = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_r)$. Then, the embedding \mathbf{b} of an original example \mathbf{x} is computed as:

$$\mathbf{b} = \Gamma^\top \Phi^\top \phi(\mathbf{x}) = \Gamma^\top (\mathcal{K}(\mathbf{x}_1, \mathbf{x}), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}))^\top.$$

4 Discussion

For developing a good graph-based dimensionality reduction method, one of the most important problems is how to construct a good graph so that the preferred data structure can be preserved. As we have discussed in Section 2, many existing methods such as MFA, LSDA and LFDA aim to design specific graphs to enhance the local compactness of the data in the same class and separate the data points with different class labels. However, none of them breaks the graph Laplacian based

framework. Our method mainly focuses on developing a new strategy to extract more information of the data manifold from a given graph. More specifically, we use the first-order Taylor expansion to incorporate the structural information from tangent spaces, i.e., the terms $\mathbf{w}_{x_j}^\top T_{x_j}^\top (\mathbf{x}_i - \mathbf{x}_j)^2$ ($i, j = 1, \dots, n$) in (10), into the scatter matrix S . Moreover, it is worth noting that although we specify a certain weight matrix (11) to construct the within-class graph G^w , our method is flexible enough to utilize the information from any valid graphs such as the one used in LFDA.

Local Tangent Space Alignment (LTSA) [17] is a popular dimensionality reduction method which also use the information from tangent spaces. The main idea of LTSA is to align every local tangent space to construct global coordinates. Although both TSD and LTSA utilize local tangent spaces, there are mainly two differences between them: 1) they actually solve different problems in essence. TSD is a linear supervised dimensionality reduction method, while LTSA is a non-linear unsupervised one. As a result, our method not only considers the class labels to make use of discriminant information, but can obtain an explicit transformation matrix to compute the mappings for out-of-sample data. 2) TSD is a graph-based method which can adopt any valid graph for training, whereas LTSA is not. This property provides TSD with much more flexibility to handle various types of data for different applications.

Our method shares the same spirit with TSIMR [13], and both of them employ tangent spaces to discover the geometrical structure of the data manifold. However, our approach and TSIMR differ in two key aspects: 1) Like LTSA, TSIMR is a non-linear unsupervised method, and thus has no ability to capture the discriminant information or give an explicit transformation matrix. 2) They have totally different objective functions. It should be noted that TSD employs (10) to construct the scatter matrix S , while the objective function of TSIMR has other terms $\|\mathbf{w}_{x_i} - T_{x_i}^\top T_{x_j} \mathbf{w}_{x_j}\|_2^2$ ($i, j = 1, \dots, n$). And we find these terms are not much beneficial for discriminant analysis.

The effect of the Tikhonov regularizer $\gamma \|\mathbf{f}\|^2 = \gamma(\|\mathbf{t}\|^2 + \sum_i \|\mathbf{w}_{x_i}\|^2)$ in (15) should be highlighted, since it plays a key role in our method. Generally, Tikhonov regularization is a common technique employed by many dimensionality reduction methods to deal with the singularity problem of the matrix, where the parameter γ is always set to a very small value. However, our method needs an appropriate large γ to penalize large $\|\mathbf{w}_{x_i}\|$ ($i = 1, \dots, n$) so that the within-class compactness can be enhanced.

5 Experiments

To evaluate the proposed method, related dimensionality reduction methods including PCA, LDA, MFA, LSDA and LFDA are compared with TSD on multiple real-world data sets from the UCI Machine Learning Repository [1], the Protein Sequence data set² from glycosylation database Uniprot (v8.0), and the USPS data set. Specifically, we first perform dimensionality reduction to map all examples into a subspace, and then carry out classification using the nearest neighbor classifier (1-NN) in the subspace. This experimental setting is the same as the one

² This Protein Sequence data set is available at <http://www.ebi.ac.uk/uniprot>.

Table 1 List of the classification data sets used in our experiments.

Data Set	Dimensionality	# of examples	# of classes	Rates of training
Satellite	36	6435	2	10%
Theorem Prove	51	3059	2	10%
Breast Cancer	9	263	2	50%
Column2C	6	310	2	50%
Image	18	2086	2	10%
Ionosphere	34	351	2	50%
Protein Sequence	420	2000	2	20%
Semeion Handwritten	256	1593	10	20%
USPS	256	2007	10	20%

adopted in [12]. Moreover, we also compare the baseline method that just employs the 1-NN classifier in the original space without performing dimensionality reduction.

Seven UCI data sets (Satellite, Theorem Prove, Breast Cancer, Column2C, Image, Ionosphere, Semeion Handwritten), the Protein Sequence data set, and the USPS data set are used to conduct our experiments. Originally, the Theorem Prove and USPS data sets are divided into a training set and a test set. For simplicity, we just use their test sets to carry on the experiments. For the Protein Sequence data set, we use a subset of the Uniprot database which contains only 99 mammalian protein entries. For each data set, we randomly split certain rates of the data as the training set and the rest as the test set. Furthermore, all the parameters for MFA, LSDA, LFDA and TSD are selected by three-fold cross-validation. The configuration of each data set is shown in Table 1.

The performance of PCA and graph-based dimensionality methods including MFA, LSDA, LFDA and TSD depend on the dimensionality of the discovered embedding subspace. Thus we show the best results obtained by those methods. Every experimental result is obtained from the average over 20 splits. We give the mean values and standard deviations of the error rates (%) on the employed data sets, where the best method is highlighted in bold font and the best and comparable ones based on the t-test with the significance level 5% are marked with ‘ Δ ’.

The experimental results on the Satellite, Theorem Prove, Breast Cancer and Column2C data sets are shown in Table 2. In most cases, classification with dimensionality reduction is statistically better than the baseline. However, LDA perform well on none of the four data sets, probably because the implicit assumption adopted by LDA mismatches the distributions of these data sets. On the other hand, all the graph-based methods get reasonable well results, because they aim to preserve the local structure of data. PCA also works well for the purpose of separating data from different classes. Although it does not attain the best performance, our method achieves comparable good results.

Table 3 describes the classification performance of each method on the Image, Ionosphere, Protein Sequence, Semeion Handwritten and USPS data sets. Again, LDA gets worse results. Surprisingly, the counterparts of TSD including MFA, LSDA, LFDA, fail to perform well for the Protein Sequence and the Semeion Handwritten data sets. In the case of Semeion Handwritten data set, these methods are even worse than the baseline. The characteristics of the feature vectors in the

Table 2 Mean values and standard deviations of the error rates (%) on the Satellite, Theorem Prove, Breast Cancer and Column2C data sets. The best method is highlighted by bold font. The best and comparable ones based on the t-test with the significance level 5% are marked with ‘ Δ ’.

Methods	Satellite	Theorem Prove	Breast Cancer	Column2C
Baseline	13.31 \pm 0.72	30.86 \pm 1.66	24.73 \pm 4.08	18.77 \pm 2.70
PCA	12.97\pm0.62Δ	30.35 \pm 1.51 Δ	30.88\pm3.15Δ	18.10 \pm 2.26 Δ
LDA	17.56 \pm 0.45	37.60 \pm 1.10	35.80 \pm 4.01	24.35 \pm 3.70
MFA	14.34 \pm 0.62	29.90\pm1.22 Δ	31.15 \pm 2.52	23.23 \pm 3.69
LSDA	13.27 \pm 0.48	31.79 \pm 1.42	32.33 \pm 2.70	22.26 \pm 3.74
LFDA	13.13 \pm 0.52 Δ	30.19 \pm 1.52 Δ	31.15 \pm 2.98 Δ	17.84\pm2.52 Δ
TSD	13.29 \pm 0.90 Δ	30.14 \pm 1.55 Δ	31.30 \pm 3.03 Δ	18.55 \pm 2.74 Δ

two data sets probably explain why this happens. Every example in the two data sets has a sparse and binary feature vector with high dimensionality in which only a small number of elements are one, and the rest are zero. For instance, the Semeion Handwritten data set contains 1593 binary images of handwritten digits consisting 16×16 pixels. In this case, the graph Laplacian based methods may not be able to capture the meaningful local geometry of the data manifold any more. On the other hand, TSD achieves the best results probably because it can capture extra geometrical information from tangent spaces. This suggests that our proposed method makes good use of the information from tangent spaces and thus can correctly discover the data structure. In addition, the limitation inherited from the graph Laplacian based framework rather than the choice of graphs in each graph-based method should be responsible for the undesirable results, since the adopted graph in TSD is similar to those in MFA, LSDA and LFDA. Moreover, even when the feature vectors are no longer sparse and binary, TSD can also get the lowest error rates compared with the other methods with high level of statistical significance in the Image and Ionosphere data sets. This demonstrates that due to utilizing the structural information from tangent spaces, TSD can not only improve the performance of dimensionality reduction, but be applied to the data sets on which the graph Laplacian based counterparts fail to perform effectively.

Table 4 gives the time consumptions of different methods. As can be seen, TSD is relatively less efficient than its counterparts, because it has to estimate the tangent space and tangent vector at each data point. In fact, this is also the weakness of other tangent space based methods [8, 13]. Therefore, proposing a strategy to make tangent space based methods more scalable can be an interesting research direction.

6 Conclusion

In this paper, we have proposed a novel supervised dimensionality reduction method named *local Tangent Space Discriminant analysis* (TSD), which differs from the methods based on the graph Laplacian framework. By introducing tangent spaces and using the first-order Taylor expansion, we develop a new strategy to utilize the information from tangent spaces, which leads to a natural way of preserving the geometrical structure of the data manifold. The proposed method aims to seek an embedding space where the local manifold structure of the data belonging to

Table 3 Mean values and standard deviations of the error rates (%) on the Image, Ionosphere, Protein Sequence, Semeion Handwritten and USPS data sets. The best method is highlighted by bold font. The best and comparable ones based on the t-test with the significance level 5% are marked with ‘ Δ ’.

Methods	Image	Ionosphere	Protein Sequence	Semeion	USPS
Baseline	9.25 \pm 1.00	14.00 \pm 1.98	33.17 \pm 1.05	15.98 \pm 1.16	12.07 \pm 0.72
PCA	9.09 \pm 1.00	10.77 \pm 1.98	22.43 \pm 0.83 Δ	13.03 \pm 1.09 Δ	11.62 \pm 0.64
LDA	25.99 \pm 2.31	16.94 \pm 2.86	38.97 \pm 2.99	42.94 \pm 2.15	20.26 \pm 1.47
MFA	7.57 \pm 1.03	10.66 \pm 2.16	30.21 \pm 3.20	44.95 \pm 2.26	19.46 \pm 1.23
LSDA	8.58 \pm 1.44	10.71 \pm 2.86	30.76 \pm 2.24	65.36 \pm 3.82	71.23 \pm 1.81
LFDA	7.77 \pm 0.98	13.94 \pm 1.75	40.51 \pm 6.91	45.27 \pm 2.47	12.07 \pm 1.57
TSD	6.73\pm0.95Δ	9.34\pm1.37 Δ	22.06\pm1.19 Δ	12.78\pm1.04Δ	10.94\pm0.66Δ

Table 4 Computation time (in seconds) of each method for dimensionality reduction.

Methods	PCA	LDA	MFA	LSDA	LFDA	TSD
Satellite	0.0014	0.0021	0.1198	0.0971	0.0191	5.4878
Theorem Prove	0.0015	0.0028	0.0457	0.0265	0.0098	3.8428
Breast Cancer	0.0002	0.0005	0.0039	0.0043	0.0022	0.4427
Column2C	0.0002	0.0004	0.0051	0.0043	0.0024	0.5647
Image	0.0004	0.0006	0.0232	0.0080	0.0037	1.0081
Ionosphere	0.0006	0.0010	0.0080	0.0071	0.0035	1.0006
Protein Sequence	0.1376	0.2093	0.2789	0.3219	0.2559	47.901
Semeion Handwritten	0.0400	0.0470	0.1152	0.0931	0.0914	17.103
USPS	0.0571	0.0555	0.1131	0.0974	0.1123	24.704

the same class is preserved as much as possible, while the marginal data points with different class labels are better separated. Moreover, TSD has the analytic solution by solving a generalized eigenvalue problem and can be easily extended to non-linear dimensionality reduction through the kernel trick.

The effectiveness of the proposed method has been demonstrated by comparing with related work on multiple real-world data sets including the UCI data sets and the Protein Sequence data set. The experimental results show that TSD works well on the data sets which can hardly be well handled by its counterparts, and attains better performance of classification due to utilizing the extra information from tangent spaces. Future work directions include extending our method to different learning scenarios such as semi-supervised learning and developing the sparse algorithm of TSD for large-scale learning tasks.

Acknowledgements This work is supported by the National Natural Science Foundation of China under Projects 61370175 and 61075005, and Shanghai Knowledge Service Platform Project (No.ZF1213).

References

1. K. Bache and M. Lichman. UCI machine learning repository, 2013.
2. M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15(6):1373–1396, 2003.
3. D. Cai, X. He, K. Zhou, J. Han, and H. Bao. Locality sensitive discriminant analysis. In *Proceedings of the 20rd International Joint Conference on Artificial Intelligence (IJCAI)*, pages 708–713, 2007.

4. F. R. K. Chung. *Spectral Graph Theory*. American Mathematical Society, Rhode Island, 1997.
5. D. L. Donoho and C. Grimes. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100(10):5591–5596, 2003.
6. K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, 2nd edition, 1990.
7. X. He and P. Niyogi. Locality preserving projections. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems 16*, pages 1–8. MIT Press, Cambridge, MA, 2004.
8. B. Lin, X. He, C. Zhang, and M. Ji. Parallel vector field embedding. *Journal of Machine Learning Research*, 14(1):2945–2977, 2013.
9. S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.
10. B. Schölkopf, A. Smola, and K. R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
11. P. Simard, Y. LeCun, and J. S. Denker. Efficient pattern recognition using a new transformation distance. In S.J. Hanson, J.D. Cowan, and C.L. Giles, editors, *Advances in Neural Information Processing Systems 5*, pages 50–58. Morgan-Kaufmann, 1993.
12. M. Sugiyama. Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis. *Journal of Machine Learning Research*, 8:1027–1061, 2007.
13. S. Sun. Tangent space intrinsic manifold regularization for data representation. In *Proceedings of the IEEE China Summit and International Conference on Signal and Information Processing (ChinaSIP)*, pages 179–183, 2013.
14. H. Tyagi, E. Vural, and P. Frossard. Tangent space estimation for smooth embeddings of riemannian manifolds. *Information and Inference*, 2(1):69–114, 2013.
15. S. Yan, D. Xu, B. Zhang, H. Zhang, Q. Yang, and S. Lin. Graph embedding and extensions: a general framework for dimensionality reduction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(1):40–51, 2007.
16. J. Ye and Q. Li. A two-stage linear discriminant analysis via QR-decomposition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(6):929–941, 2005.
17. Z. Zhang and H. Zha. Principal manifolds and nonlinear dimension reduction via local tangent space alignment. *SIAM Journal on Scientific Computing*, 26(1):313–338, 2004.
18. M. Zhu and A. M. Martinez. Subclass discriminant analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(8):1274–1286, 2006.

Appendix A. Detailed Derivation of S

In order to fix S , we decompose (10) into three additive terms as follows:

$$\begin{aligned}
 \mathbf{f}^\top S \mathbf{f} &= \underbrace{\sum_{i,j=1}^n W_{ij}^w ((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{t})^2}_{\text{term one}} + \\
 &\quad \underbrace{\sum_{i,j=1}^n W_{ij}^w (\mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j))^2}_{\text{term two}} + \\
 &\quad \underbrace{\sum_{i,j=1}^n W_{ij}^w [-2((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{t}) \mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j)]}_{\text{term three}},
 \end{aligned}$$

and then examine their separate contributions to the whole S .

Term One

$$\begin{aligned} & \sum_{i,j=1}^n W_{ij}^w ((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{t})^2 \\ &= 2\mathbf{t}^\top X(D^w - W^w)X^\top \mathbf{t} = 2\mathbf{t}^\top XL^wX^\top \mathbf{t}, \end{aligned}$$

where D^w is a diagonal weight matrix with $D_{ii}^w = \sum_{j=1}^n W_{ij}^w$, and $L^w = D^w - W^w$ is the Laplacian matrix. Then we have $S_1 = 2(D^w - W^w) = 2L^w$. And term one contributes to XS_1X^\top in (14).

Term Two Define $B_{ji} = T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j)$, then

$$\begin{aligned} & \sum_{i,j=1}^n W_{ij}^w (\mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 \\ &= \sum_{i,j=1}^n W_{ij}^w (\mathbf{w}_{\mathbf{x}_j}^\top B_{ji})^2 \\ &= \sum_{i,j=1}^n W_{ij}^w \mathbf{w}_{\mathbf{x}_j}^\top B_{ji} B_{ji}^\top \mathbf{w}_{\mathbf{x}_j} \\ &= \sum_{j=1}^n \mathbf{w}_{\mathbf{x}_j}^\top \left(\sum_{i=1}^n W_{ij}^w B_{ji} B_{ji}^\top \right) \mathbf{w}_{\mathbf{x}_j} = \sum_{i=1}^n \mathbf{w}_{\mathbf{x}_i}^\top H_i \mathbf{w}_{\mathbf{x}_i}, \end{aligned}$$

where we have defined matrices $\{H_j\}_{j=1}^n$ with $H_j = \sum_{i=1}^n W_{ij}^w B_{ji} B_{ji}^\top$.

Now suppose we define a block diagonal matrix S_3 sized $mn \times mn$ with block size $m \times m$. Set the (i, i) -th block ($i = 1, \dots, n$) of S_3 to be H_i . Then the resultant S_3 is the contribution of term two for S in (14).

Term Three Define vectors $\{F_j\}_{j=1}^n$ with $F_j = \sum_{i=1}^n W_{ij}^w B_{ji}$, then term three can be rewritten as:

$$\begin{aligned} & \sum_{i,j=1}^n W_{ij}^w [-2((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{t}) \mathbf{w}_{\mathbf{x}_j}^\top T_{\mathbf{x}_j}^\top (\mathbf{x}_i - \mathbf{x}_j)] \\ &= \sum_{i,j=1}^n 2W_{ij}^w [((\mathbf{x}_j - \mathbf{x}_i)^\top \mathbf{t}) \mathbf{w}_{\mathbf{x}_j}^\top B_{ji}] \\ &= \sum_{i,j=1}^n W_{ij}^w (-\mathbf{t}^\top \mathbf{x}_i B_{ji}^\top \mathbf{w}_{\mathbf{x}_j}) + \sum_{i=1}^n \mathbf{t}^\top \mathbf{x}_i F_i^\top \mathbf{w}_{\mathbf{x}_i} + \\ & \quad \sum_{i,j=1}^n W_{ij}^w (-\mathbf{w}_{\mathbf{x}_j}^\top B_{ji} \mathbf{x}_i^\top \mathbf{t}) + \sum_{i=1}^n \mathbf{w}_{\mathbf{x}_i}^\top F_i \mathbf{x}_i^\top \mathbf{t}. \end{aligned}$$

From this expression, we can give the formulation of S_2 . Then the S_2^\top in (14), which is its transpose, is ready to get.

Suppose we define two block matrices S_2^1 and S_2^2 sized $n \times mn$ each where the block size is $1 \times m$, and S_2^2 is a block diagonal matrix. Set the (i, j) -th block ($i, j = 1, \dots, n$) of S_2^1 to be $-W_{ij}^w B_{ji}^\top$, and the (i, i) -th block ($i = 1, \dots, n$) of S_2^2 to be F_i^\top . Then, term three can be rewritten as: $\mathbf{t}^\top X(S_2^1 + S_2^2)\mathbf{w} + \mathbf{w}^\top (S_2^1 + S_2^2)^\top X^\top \mathbf{t}$. It is clear that $S_2 = S_2^1 + S_2^2$.