



Learning dictionary on manifolds for image classification

Bao-Di Liu^{a,*}, Yu-Xiong Wang^a, Yu-Jin Zhang^a, Bin Shen^b

^a Department of Electronic Engineering, Tsinghua University, Beijing 100084, China

^b Department of Computer Science, Purdue University, West Lafayette, IN 47907, USA

ARTICLE INFO

Available online 23 November 2012

Keywords:

Sparse coding
Image classification
Locally linear embedding
Coordinate descent
Manifold

ABSTRACT

At present, dictionary based models have been widely used in image classification. The image features are approximated as a linear combination of bases selected from the dictionary in a sparse space, resulting in compact patterns. The features applied to image classification usually reside on low dimensional manifolds embedded in a high dimensional ambient space; traditional sparse coding algorithm, however, does not consider this topological structure. It can be characterized naturally by linear coefficients that reconstruct each data point from its neighbors. One of the central issues here is how to determine the neighbors and learn the coefficients. In this paper, the geometrical structures are encoded in two situations. In simple cases when data points distribute on a single manifold, it is explicitly modeled by locally linear embedding algorithm combined with k -nearest neighbors. Nevertheless, in real-world scenarios, complex data points often lie on multiple manifolds. Sparse representation algorithm combined with k -nearest neighbors is instead utilized to construct the topological structures, because it is capable of approximating the data point by selecting its homogenous neighbors adaptively to guarantee the smoothness of each manifold. After obtaining the local fitting relationship, these two topological structures are then embedded into sparse coding algorithm as regularization terms to formulate the corresponding objective functions of dictionary learning on single manifold (DLSM) and dictionary learning on multiple manifolds (DLMM), respectively. Upon this, a coordinate descent scheme is proposed to solve the unified optimization problems. Experimental results on several benchmark data sets, such as Caltech-256, Caltech-101, Scene 15, and UIUC-Sports, show that our proposed algorithms equal or outperform other state-of-the-art image classification algorithms.

© 2012 Elsevier Ltd. All rights reserved.

1. Introduction

Recently, image classification, which aims at associating images with semantic labels automatically, has become quite a significant topic in computer vision. The most common framework is the discriminative model [1,2,3,4]. Typically, learning vocabulary by applying k -means clustering on image patches (also called bag-of-words model) [5] and combining with hard-assignment vector quantization (VQ) is the most popular method at the early stage. Hard-assignment VQ method treats an image as a collection of “Visual words” (vocabulary), and each image patch is mapped to one word in the vocabulary. After that, several variants of vocabulary combined with vector quantization have been proposed in improving the classification performance recently. Jurie and Triggs [6] considered the densest point rather than the center of uniform region as the clustering center to learn vocabulary. Lazebnik et al. [1] extended the hard-assignment VQ

method with spatial pyramid matching (SPM) kernel to compensate the loss of spatial information. Wu and Rehg [3] proposed histogram intersection kernel to form the vocabulary and applied one-class SVM to retrain the vocabulary and achieved satisfying results. In the year 2009, Yang et al. [2] used sparse coding algorithm for learning dictionary and coding images, resulting in state-of-the-art performance in image classification. Compared with hard-assignment VQ method, sparse coding algorithm can achieve sparse approximations with lower reconstruction.

On the other hand, some recent research work suggested that image space is actually a smooth low dimensional sub-manifold embedded in a high dimensional ambient space. Many manifold learning algorithms, such as locally linear embedding (LLE) [7], ISOMAP [8], and Laplacian Eigenmaps [9] were proposed to explicitly explore the intrinsic topological structure, which could significantly enhance the dimensionality reduction performance. All these algorithms consider the locality property and preserve it when learning the patterns. As suggested in [10], locality was more essential than sparsity, since locality can lead to sparsity while sparsity cannot cause locality. Therefore, more and more researchers focused on locality-preserving during dictionary learning for image classification. van Gemert et al. [11] proposed kernel codebook and

* Corresponding author. Tel.: +86 10 62798336; fax: +86 10 62770317.

E-mail addresses: lbd08@mails.tsinghua.edu.cn (B.-D. Liu), albertwyx@gmail.com (Y.-X. Wang), zhang-yj@mail.tsinghua.edu.cn (Y.-J. Zhang), stanshenbin@gmail.com (B. Shen).

soft-assignment vector quantization to preserve the local properties. Wang et al. [12] considered that each word in the vocabulary was on a manifold, and utilized locally linear coding [7] for vector quantization to preserve the local information on vocabulary. Yan et al. [13] proposed a general framework for dimensionality reduction called graph embedding and claimed that most of graph embedding methods could be unified within this framework. Gao et al. [14] proposed to incorporate the histogram intersection kernel based Laplacian matrix into the objective function of sparse coding to enforce the consistence in sparse representation of similar local features. Yang et al. [15] proposed intrinsic graph and penalty graph to preserve the intrinsic graph and laid off the penalty graph in a supervised way. Zheng et al. [16] proposed to incorporate the vector quantization based Laplacian matrix into the objective function of sparse coding. Lu and Peng [17] proposed to incorporate the hypergraph (vertex, hyperedge, incidence matrix and hyperedge weights) regularization term into the objective function of sparse coding. Ramamurthy et al. [18] assumed that data points distributed on the same manifold and proposed a manifold projection to improve traditional sparse coding. Shen and Si [19] proposed to construct multiple manifolds structures by sparse representation algorithm, however, locality was not explicitly considered. Liu et al. [20] proposed a discriminant sparse coding scheme to incorporate the label information into sparse coding algorithm.

Inspired by both the superior performance of sparse coding based dictionary learning for image classification and the enhancement of dimension reduction on manifold, dictionary learning on single manifold (DLSM) and multiple manifolds (DLMM) are integrated and proposed in this paper. For DLSM, the intrinsic topological structure of the original data is explicitly modeled under the assumption that the original data could be fit by the linear combination of its all k -nearest neighbors which is then embedded into the objective function of sparse coding algorithm as the regularization term, while for DLMM, the original data could be fit by the linear combination of its neighbors, which only lie on the same manifold with the original data, selected from k -nearest neighbors. After that, a coordinate descent scheme [21] with guaranteed convergence is proposed to solve the unified optimization problems. The proposed dictionary learning methods on manifolds for image classification are evaluated on four benchmark data sets and achieve higher classification accuracy than traditional sparse coding algorithm.

The major contributions of this paper are as follows. First, two types of topological structures are proposed and analyzed on why and how these two graph models construct the manifold structures. Second, a coordinate descent scheme is proposed to solve the DLSM and DLMM as unified optimization problems. Third, the proposed algorithms are evaluated in image classification task; as shown experimentally in Section 6, the performances of our algorithms equal or outperform other state-of-the-art image classification algorithms on several benchmark data sets.

The rest of the paper is organized as follows. Section 2 overviews some related work contributing to image classification. Dictionary learning methods on manifolds are proposed in Section 3. The solution to the minimization of the objective function and guaranteed convergence are elaborated in Section 4. And implementation for image coding and spatial pooling is given in Section 5. Then, experimental results and analysis are shown in Section 6. Finally, discussions and conclusions are drawn in Section 7.

2. Related work

In this section, some abbreviations and notations are given first. Then several popular dictionary learning and image coding methods are listed.

Table 1
Some abbreviations.

Abbreviations	Full name
DLMM	Dictionary learning on multiple manifolds
DLSM	Dictionary learning on single manifold
HIK	Histogram intersection kernel
HIKVQ	VQ method combining with HIK
KC	Kernel codebook
KScSPM	Kernel ScSPM
KSPM	Nonlinear kernel SPM
LLC	Locality-constrained linear coding
LLE	Locally linear embedding
LScSPM	Laplacian ScSPM
OCSVM	One-class SVM for generating codebook
SC	Sparse coding
ScSPM	Sparse coding based linear SPM
SPM	Spatial pyramid matching
VQ	Vector quantization

2.1. Some abbreviations and notations

For convenience, some abbreviations are as shown in Table 1.

Let $\mathbf{X} \in \mathbb{R}^{D \times N}$ represent the local descriptors extracted from training images for learning dictionary, where D is the dimension of \mathbf{X} , and N is the number of samples in \mathbf{X} . Let $\mathbf{B} \in \mathbb{R}^{D \times K}$ be the dictionary, and $\mathbf{S} \in \mathbb{R}^{K \times N}$ be the corresponding codes, where K is the size of the dictionary. Let $\mathbf{Y} \in \mathbb{R}^{D \times M}$ be the features extracted from an image \mathbf{J} , and $\mathbf{V} \in \mathbb{R}^{K \times M}$ represent the corresponding codes. Let $\mathbf{Z} \in \mathbb{R}^{(K \times L) \times 1}$ be a vector after spatial pooling to represent the image, where L is the total number of regions in each layer split by SPM ($L=1+4+16=21$). Let $\|\cdot\|_F$ represent the Frobenius norm. Let $\mathbf{A}_{\bullet n}$ and $\mathbf{A}_{m \bullet}$ denote the n th column and m th row vectors of matrix \mathbf{A} , respectively. Let $\text{tr}\{\mathbf{A}\}$ represent the trace of matrix \mathbf{A} . Let knn represent the number of neighbors for each descriptor.

2.2. Several popular dictionary learning and image coding methods

The aim of dictionary learning [22] is to find the optimum dictionary to make $\mathbf{X} \approx \mathbf{BS}$. After obtaining the dictionary, for each image \mathbf{J} , we assume that f_c and f_p denote coding and pooling operators. The image coding and pooling step can be formulated as

$$\mathbf{V} = f_c(\mathbf{Y}), \quad \mathbf{Z} = f_p(\mathbf{V}) \quad (1)$$

Usually, image coding is related to dictionary learning, and pooling strategy is max [2] or average. LIBSVM [23] is adopted for classifier training. Therein, dictionary learning and image coding are the most focused steps. Several popular dictionary learning and image coding methods are as follows.

Hard-assignment vector quantization for image classification: Hard-assignment vector quantization method is the most popular method, which solves the following optimization problem:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} \quad & f(\mathbf{B}, \mathbf{S}) = \|\mathbf{X} - \mathbf{BS}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{S}_{\bullet i}\|_0 = 1, \|\mathbf{S}_{\bullet i}\|_1 = 1, \mathbf{S}_{\bullet i} \geq 0, \quad \forall i \end{aligned} \quad (2)$$

The above optimization problem is usually solved by k -means algorithm. $\|\mathbf{S}_{\bullet i}\|_0 = 1$ means that there is only one nonzero element in each $\mathbf{S}_{\bullet i}$, i.e. each $\mathbf{X}_{\bullet i}$ is represented by only one basis in dictionary \mathbf{B} . $\|\mathbf{S}_{\bullet i}\|_1 = 1$ and $\mathbf{S}_{\bullet i} \geq 0$ mean that the weight for the corresponding basis is 1 to represent $\mathbf{X}_{\bullet i}$.

The distance between local feature and dictionary (also called codebook) utilizes histogram intersection kernel in [3]. The objective function and kernel function are defined as (3) and (4),

respectively

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} \quad & f(\mathbf{B}, \mathbf{S}) = \|\phi(\mathbf{X}) - \phi(\mathbf{B})\mathbf{S}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{S}_{\bullet i}\|_0 = 1, \|\mathbf{S}_{\bullet i}\|_1 = 1, \mathbf{S}_{\bullet i} \geq 0, \quad \forall i \end{aligned} \quad (3)$$

$$\kappa_{\text{hik}}(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet j}) = \phi(\mathbf{X}_{\bullet i})^T \phi(\mathbf{B}_{\bullet j}) = \sum_{d=1}^D \min(\mathbf{X}_{di}, \mathbf{B}_{dj}) \quad (4)$$

Hard-assignment vector quantization method is simple and intuitive. However, the reconstruction error would be too high due to too much constraints, thus leading to lower performance. Notably, Hard-assignment vector quantization obtains better performance with histogram intersection kernel than that with linear kernel.

Soft-assignment vector quantization for image classification: The soft-assignment vector quantization method also adopts the optimization problem in Eq. (2) to build dictionary. To reduce the reconstruction error, soft-assignment vector quantization method loosens the restrictions of descriptor's coding, which uses a linear combination of multiple bases to approximate $\mathbf{X}_{\bullet i}$

$$\mathbf{S}_{ki} = \frac{\kappa(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet k})}{\sum_{m=1}^K \kappa(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet m})} \quad (5)$$

The kernel function κ here can be any types of kernels, [11] adopted gaussian kernel

$$\kappa_{\text{gauss}}(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet j}) = \exp(-\beta \|\mathbf{X}_{\bullet i} - \mathbf{B}_{\bullet j}\|_2^2) \quad (6)$$

where β is the soft assignments factor. The soft-assignment vector quantization achieves good performance because the Gaussian kernel is performed as a local constraint, i.e. each bin of $\mathbf{S}_{\bullet i}$ is corresponding to the weight of bases, and the smaller distance between descriptor $\mathbf{X}_{\bullet i}$ and dictionary \mathbf{B} will get a higher weight.

Locality-constrained linear coding (LLC) for image classification: Locality-constrained linear coding [12] for image classification uses the following criteria:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} \quad & f(\mathbf{B}, \mathbf{S}) = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + \lambda \|\kappa_{\text{gauss}}(\mathbf{X}, \mathbf{B}) \odot \mathbf{S}\|_F^2 \\ \text{s.t.} \quad & \mathbf{1}^T \mathbf{B}_{\bullet k} = 1, \quad \forall k \end{aligned} \quad (7)$$

where

$$\kappa_{\text{gauss}}(\mathbf{X}, \mathbf{B}) = [\kappa_{\text{gauss}}(\mathbf{X}_{\bullet 1}, \mathbf{B}), \dots, \kappa_{\text{gauss}}(\mathbf{X}_{\bullet N}, \mathbf{B})]$$

$$\kappa_{\text{gauss}}(\mathbf{X}_{\bullet i}, \mathbf{B}) = [\exp(\text{dist}(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet 1})/\sigma), \dots, \exp(\text{dist}(\mathbf{X}_{\bullet i}, \mathbf{B}_{\bullet K})/\sigma)]^T$$

\odot denotes the element-wise multiplication, λ is used for adjusting the locality. The larger the λ is, the more locality is considered. LLC has several advantages. First, it can better reconstruct the raw descriptors than VQ method. Second, the raw descriptor and its neighboring bases construct a local coordinate system, which would lead to local smooth sparsity, i.e. the explicit locality adopted in LLC ensures that similar patches have similar codes. Third, this optimization problem has the analytical solution. The LLC method obtains better performance than VQ-based method.

Sparse coding for image classification: Sparse coding (SC) can be considered as methods of rearranging the structure of the original data in order to make the energy compact under over-complete or non-orthogonal bases. Hence, the data point can be represented as a linear combination of only few active bases that possess overwhelmingly majority energy of the data. Sparse coding method was firstly introduced into image classification in [2]. Unlike VQ method, sparse coding loosens the constraint on the codes, each code can be represented by the linear combination of several bases to minimize the reconstruction error. The sparse coding algorithm can be written as the following formula:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} \quad & f(\mathbf{B}, \mathbf{S}) = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha \|\mathbf{S}\|_1 \\ \text{s.t.} \quad & \|\mathbf{B}_{\bullet i}\|_2 = 1, \quad \forall i = 1, 2, \dots, K \end{aligned} \quad (8)$$

The regularization term is to control sparsity in \mathbf{S} , where α is a regularization parameter to control the tradeoff between fitting goodness and sparseness. Sparse coding method has several advantages. First, it can learn over-complete bases, i.e. $K > D$. That is to say, each feature has opportunity to choose better base. Second, it can achieve less reconstruction error than VQ method and LLC method. Third, it can capture salient patterns of local features. So image classification by sparse coding framework achieves state-of-the-art performance on several benchmarks such as Caltech-101, Scene 15, etc.

Locality-constrained sparse coding for image classification: Sparse coding for image classification shows good performance, however, due to the over-complete dictionary and the independent coding process [14], the locality or the geometrical structure among the instances to be encoded are lost. Gao et al. [14] and Zheng et al. [16] proposed locality-constrained sparse coding to preserve the local manifold structure of the instances by embedding the Laplacian matrix \mathbf{G} into sparse coding algorithm. The objective function for locality-constrained sparse coding is as follows:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} \quad & f(\mathbf{B}, \mathbf{S}) = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha \|\mathbf{S}\|_1 + \eta \text{tr}\{\mathbf{S}\mathbf{G}\mathbf{S}^T\} \\ \text{s.t.} \quad & \|\mathbf{B}_{\bullet i}\|_2 = 1, \quad \forall i = 1, 2, \dots, K \end{aligned} \quad (9)$$

where η is the regularization parameter balancing the weight between the fitting goodness and locality preservation. The Laplacian matrix $\mathbf{G} \in \mathbb{R}^{N \times N}$ can be obtained by $\mathbf{G} = \mathbf{W} - \mathbf{U}$, where [16] constructed the k -nearest neighbors' distance matrix \mathbf{U} as follows:

$$\mathbf{U}_{ij} = \begin{cases} 1 & \text{if } \mathbf{X}_{\bullet j} \in \mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\} \\ 0 & \text{if } \mathbf{X}_{\bullet j} \notin \mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\} \end{cases} \quad (10)$$

Gao et al. [14] constructed the k -nearest neighbors' distance matrix \mathbf{U} as follows,

$$\mathbf{U}_{ij} = \begin{cases} \sum_{d=1}^D \min(\mathbf{X}_{di}, \mathbf{X}_{dj}) & \text{if } \mathbf{X}_{\bullet j} \in \mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\} \\ 0 & \text{if } \mathbf{X}_{\bullet j} \notin \mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\} \end{cases} \quad (11)$$

$\mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\}$ represents k -nearest neighbors of $\mathbf{X}_{\bullet i}$. The degree matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$ is a diagonal matrix and $\mathbf{W}_{ii} = \sum_{j=1}^N \mathbf{U}_{ij}$.

3. Proposed dictionary learning methods on manifolds

In this section, dictionary learning methods on manifolds are proposed for image classification. First, the geometrical structure of the single manifold is constructed by locally linear embedding algorithm combined with k -nearest neighbors. Then, sparse representation combined with k -nearest neighbors is adopted to construct the geometrical structure of the multiple manifolds. After that, we preserve these two structures to form dictionary learning on single manifold (DLSM) method and dictionary learning on multiple manifolds (DLMM) method, respectively.

3.1. Image features lie on a single manifold

In the classification application, the image data are often sampled from a nonlinear low dimensional manifold embedding in a high dimensional space (see in Fig. 1 (a)). In the situation of a single manifold, each data point and its neighboring points are close to or even lie on a locally flat patch of the manifold. A very straightforward idea in modeling this global nonlinear structure is to learn from locally linear fitting.

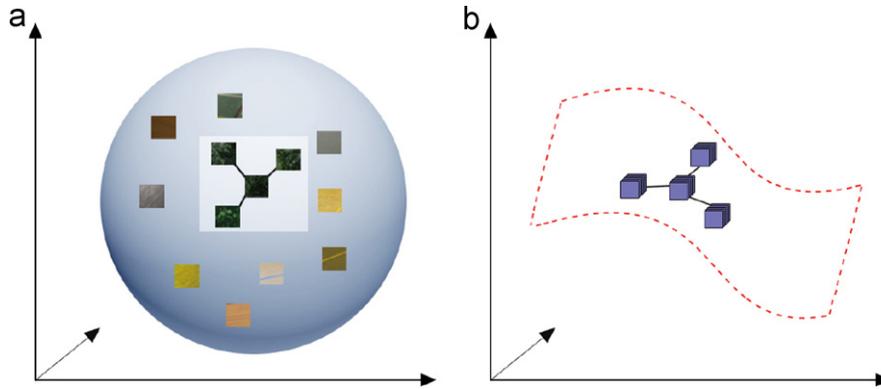


Fig. 1. Data in the features space and codes space. (a) The original image features lying on the manifold and (b) the corresponding feature codes.

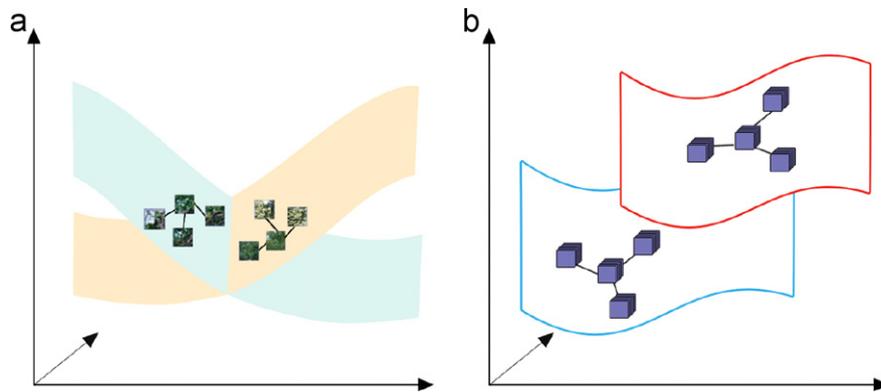


Fig. 2. Data in the features space and codes space. (a) The original image features lying on the two manifolds and (b) the corresponding feature codes.

Let $\mathbf{X}_{\bullet i}$ represent the i th data point, and $\mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\}$ represent its k -nearest neighbors. For most cases, five nearest neighbors are considered. The local topology around $\mathbf{X}_{\bullet i}$ can be characterized by the linear combination of its neighbors as follows:

$$\begin{aligned} \mathbf{X}_{\bullet i} &\approx \mathbf{X}\mathbf{U}_{\bullet i} \\ \text{s.t. } \sum_{j=1}^N \mathbf{U}_{ij} &= 1 \& \mathbf{U}_{ij} = 0 \text{ for } \mathbf{X}_{\bullet j} \notin \mathcal{N}_{knn}\{\mathbf{X}_{\bullet i}\} \end{aligned} \quad (12)$$

where $\mathbf{U} \in \mathbb{R}^{N \times N}$ is the coefficient matrix.

In this paper, k is set to 5 and locally linear embedding (LLE) [7] algorithm is used to obtain the coefficient matrix \mathbf{U} of the single manifold's graph.

3.2. Image features lie on multiple manifolds

The majority of real-world data often reside on multiple manifolds [24], whose manifold structures, such as curvature variation, may differ a lot. They may also overlap or intersect (Fig. 2(a)). This makes the whole scenario much more complicated. By manifold learning, we hope to analyze these manifolds separately. We still model the underlying geometrical structure patchwise. Thus for a certain data point, the true neighborhood on the same manifold instead of the entire space is needed firstly. However, existing manifold learning algorithms, such as locally linear embedding algorithm and Laplacian eigenmaps, define the neighborhood relationship only according to the Euclidean distance, which is insufficient. In other words, data points on a nearby manifold may be involved mistakenly as the neighboring point in the fix-sized neighborhood used in LLE. To preserve the topology of multiple manifolds, we assume that each manifold is smooth [25], and the local data points lying on the same manifold are homogenous in a certain sense such as direction, while

overlaps or intersections under different manifolds are not ruled out. We propose to utilize sparse representation algorithm combined with k -nearest neighbors to fit the topology of these multiple manifolds, with the consideration that sparse representation is capable of approximating the data point by selecting the homogenous neighbors adaptively and thus guarantees the smoothness of the manifold. This is helpful in identifying the neighborhood relationship on the same single manifold. And in fact we want to preserve the structure of the same single manifold in the following processing.

To demonstrate our conclusion, we carry out an experiment on ORL face database. The ORL face database contains 10 images for each of forty human faces with each lying on the same manifold. Each data point and its k -nearest neighbors (5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 100, 150, 200, 250, 300, 350, 399) are executed twice by locally linear embedding algorithm and sparse representation algorithm, respectively. For each data point, given the initial size of neighborhood (characterized by the number of selected neighboring points), the percentages of number and absolute coefficients summation of neighbors on the same manifold (belongs to the same human) are calculated. And the means of the percentages of number and absolute coefficients summation of neighbors lying on the same manifold for all data points are given. For locally linear embedding, the initial nearby points are all used as the neighboring points; whereas, for sparse representation algorithm, the active points with nonzero coefficients after sparse coding, the subset of the initial points, are selected as the neighboring points. Fig. 3 shows the percentage of number of neighbors lying on the same manifold with varying number of neighbors. From Fig. 3, we can see that most data points and their five neighbors lie on the same manifold. The percentage of number of neighbors on the same manifold is 83.45% among k -nearest neighbors, and 83.33% selected after

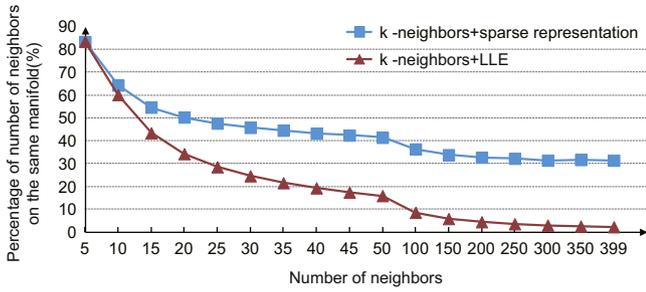


Fig. 3. The percentage of neighbors on the same manifold with varying number of neighbors. The percentage is the ratio of the number of the neighbors lying on the same manifold with nonzero fitting coefficients to the number of the overall neighbors with nonzero fitting coefficients.

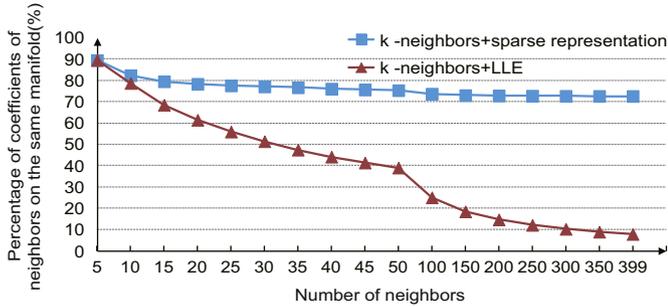


Fig. 4. The percentage of summation of absolute fitting coefficients of neighbors on the same manifold with varying number of neighbors. The percentage is the ratio of the sum of absolute nonzero fitting coefficients of neighbors lying on the same manifold to the sum of absolute nonzero fitting coefficients of neighbors.

sparse representation combined with k -nearest neighbors. With the increasing number of neighbors, this value decreases till 2.26% for 399 among k -nearest neighbors. This indicates that almost all data points' neighbors are not on the same manifold. However, this percentage achieves 31.37% selected after sparse representation. So sparse representation algorithm combined with k -nearest neighbors is capable of selecting effective ones from neighbors and removing the neighbors that do not lie on the same manifold especially when the initial nearby points distribute on the different manifold structures. Another statistical quantity, the percentage of summation of absolute fitting coefficients of neighbors on the same manifold in all the coefficients constructing the graph model, is much more convincing than the percentage of number of neighbors on the same manifold. From Fig. 4, we can see that even for 399 neighbors, (i.e. only nine neighbors are on the same manifold, and other 390 neighbors distribute on the other 39 manifolds), the percentage of summation of absolute fitting coefficients of neighbors on the same manifold by sparse representation combined with k -nearest neighbors is 72.54%, much higher than 7.89% by locally linear embedding combined with k -nearest neighbors. And with the number of neighbors increasing, this percentage by sparse representation combined with k -nearest neighbors maintains at 72%. This result is beneficial for complex scenario with multiple manifolds. From the comparison, sparse representation combined with k -nearest neighbors is capable of capturing the similar response as neighbors on the same manifold and therefore constructing the graph model for multiple manifolds in an unsupervised way. In other words, it is sort of resistant to the initial k -nearest neighboring points. Here, we use l_1 -norm minimization technique:

$$f(\mathbf{U}_i^{knn}) = \|\mathbf{X}_i - \mathcal{N}_{knn}(\mathbf{X}_i)\mathbf{U}_i^{knn}T\|^2 + 2\lambda|\mathbf{U}_i^{knn}| \quad (13)$$

where $\mathbf{X}_i \in \mathbb{R}^{D \times 1}$ represents the i th data point of \mathbf{X} . $\mathcal{N}_{knn}(\mathbf{X}_i) \in \mathbb{R}^{D \times knn}$ represents the data point's k -neighbors.

$\mathbf{U}_i^{knn} \in \mathbb{R}^{1 \times knn}$ is the corresponding sparse coefficient under the bases $\mathcal{N}_k(\mathbf{X}_i)$.

$$\mathbf{U}_{ij} = \begin{cases} \mathbf{U}_{ij}^{knn} & \text{if } \mathbf{X}_j \in \mathcal{N}_{knn}(\mathbf{X}_i) \\ 0 & \text{if } \mathbf{X}_j \notin \mathcal{N}_{knn}(\mathbf{X}_i) \end{cases} \quad (14)$$

In this section and Section 6.6.3, feature-sign search algorithm [22] is used, and the λ in Eq. (13) is set to 80. In the later part of the paper for image classification, knn is set to 10, feature-sign search algorithm [22] is adopted to obtain the coefficient matrix \mathbf{U} of multiple manifolds graph, and λ is set to 0.05.

3.3. The objective function of dictionary learning on manifolds

The geometrical structure represented by $\mathbf{U} \in \mathbb{R}^{N \times N}$ is preserved during the sparse coding procedure. Specifically, when the features \mathbf{X} are transformed to the new features $\mathbf{S} \in \mathbb{R}^{K \times N}$, the local geometrical structure is preserved through the relationship \mathbf{U} for \mathbf{S} , which can be obtained by minimizing the following term (see in Figs. 1 and 2):

$$\sum_i \|\mathbf{S}_i - \mathbf{S}\mathbf{U}_i\|_F^2 = \|\mathbf{S} - \mathbf{S}\mathbf{U}\|_F^2 = \text{tr}(\mathbf{S}(\mathbf{I} - \mathbf{U})(\mathbf{I} - \mathbf{U})^T \mathbf{S}^T) = \text{tr}(\mathbf{G}\mathbf{S}\mathbf{S}^T) \quad (15)$$

where $\mathbf{G} \in \mathbb{R}^{N \times N}$, $\text{tr}(\mathbf{G}\mathbf{S}\mathbf{S}^T)$ is the preservation of local geometry structure in the sparse space.

We incorporate the preservation term of local geometry structure $\text{tr}(\mathbf{G}\mathbf{S}\mathbf{S}^T)$ into sparse coding algorithm as a regularization term. The dictionary learning on single manifold (DLSM) algorithm and dictionary learning on multiple manifolds (DLMM) can be unified into the same framework and gained from minimizing the following objective function:

$$\mathcal{O}_{DL} = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha\|\mathbf{S}\|_1 + \eta \text{tr}(\mathbf{G}\mathbf{S}\mathbf{S}^T) \quad (16)$$

where the first term is the fitting error, the second term encourages the sparseness, and the third term is the locality constraint. η is the regularization parameter balancing the weight between the fitting goodness and geometrical structure. α is the regularization parameter adjusting the sparsity. When applying the \mathbf{U} in Eqs. (12)–(15), dictionary learning on single manifold (DLSM) can be obtained by solving Eq. (16). When applying the \mathbf{U} in Eqs. (13)–(15), dictionary learning on multiple manifolds (DLMM) can be obtained by solving Eq. (16). It is important to note that the form of our objective function is similar to that of [14,16]; however, it is quite different as for representing the topology. For the graph model preservation in [14,16], the author proposed to use Laplacian matrix to preserve the distance between the data point and its neighbors. Our DSLM method utilizes the weight matrix obtained by LLE algorithm combined with k -nearest neighbors to preserve the data point's fitting coefficients. And our DLMM method adopts the weight matrix obtained by sparse representation algorithm combined with k -nearest neighbors to preserve the data point's fitting coefficients. Our DLMM method has the advantage of handling multiple manifolds scenario.

4. Optimization of the objective function

In this section, we focus on solving the minimization of the objective function proposed in (16). This optimization problem is not jointly convex in both \mathbf{B} and \mathbf{S} , while it is separately convex in either \mathbf{B} or \mathbf{S} with \mathbf{S} or \mathbf{B} fixed. So it can be decoupled into the following two optimization subproblems which can be solved by alternating minimizations.

Finding the sparse codes is as follows:

$$\min f(\mathbf{S}) = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha\|\mathbf{S}\|_1 + \eta \text{tr}(\mathbf{G}\mathbf{S}\mathbf{S}^T) \quad (17)$$

Learning bases are as follows:

$$\begin{aligned} \min f(\mathbf{B}) &= \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 \\ \text{s.t. } \|\mathbf{B}_{\bullet i}\|_2 &= 1, \quad \forall i = 1, 2, \dots, K \end{aligned} \quad (18)$$

In the following subsection, a coordinate descent algorithm is introduced to resolve these two optimization problems.

4.1. Finding sparse codes

By fixing \mathbf{B} , we update \mathbf{S} to decrease the value of the objective function

$$\begin{aligned} f(\mathbf{S}) &= \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha\|\mathbf{S}\|_1 + \eta \operatorname{tr}\{\mathbf{S}\mathbf{G}\mathbf{S}^T\} \\ &= \operatorname{tr}\{\mathbf{X}^T\mathbf{X} - 2\mathbf{X}^T\mathbf{B}\mathbf{S} + \mathbf{S}^T\mathbf{B}^T\mathbf{B}\mathbf{S}\} + 2\alpha\|\mathbf{S}\|_1 + \eta \operatorname{tr}\{\mathbf{S}\mathbf{G}\mathbf{S}^T\} \\ &= \operatorname{tr}\{\mathbf{X}^T\mathbf{X}\} - 2 \sum_{n=1}^N [\mathbf{X}^T\mathbf{B}]_{n\bullet} \mathbf{S}_{\bullet n} + \sum_{n=1}^N \mathbf{S}_{\bullet n}^T \mathbf{B}^T \mathbf{B} \mathbf{S}_{\bullet n} \\ &\quad + 2\alpha \sum_{k=1}^K \sum_{n=1}^N |\mathbf{S}_{kn}| + \eta \sum_{k=1}^K \mathbf{S}_{k\bullet} \mathbf{G} \mathbf{S}_{k\bullet}^T \end{aligned} \quad (19)$$

Ignoring the constant term $\operatorname{tr}\{\mathbf{X}^T\mathbf{X}\}$, the objective function of \mathbf{S}_{kn} reduces to (20) with \mathbf{B} and $\{\mathbf{S}_{1n}, \mathbf{S}_{2n}, \dots, \mathbf{S}_{kn}\} / \mathbf{S}_{kn}$ fixed

$$\begin{aligned} f(\mathbf{S}_{kn}) &= \mathbf{S}_{kn}^2 \{[\mathbf{B}^T \mathbf{B}]_{kk} + \eta \mathbf{G}_{nn}\} + 2\alpha |\mathbf{S}_{kn}| \\ &\quad + 2\mathbf{S}_{kn} \left\{ \sum_{l=1, l \neq k}^K [\mathbf{B}^T \mathbf{B}]_{kl} \mathbf{S}_{ln} + \eta \sum_{r=1, r \neq n}^N \mathbf{G}_{nr} \mathbf{S}_{kr} - [\mathbf{B}^T \mathbf{X}]_{kn} \right\} \\ &= \mathbf{S}_{kn}^2 \{[\mathbf{B}^T \mathbf{B}]_{kk} + \eta \mathbf{G}_{nn}\} + 2\alpha |\mathbf{S}_{kn}| - 2\mathbf{S}_{kn} \mathbf{H}_{kn} \end{aligned} \quad (20)$$

where $\mathbf{H}_{kn} = [\mathbf{B}^T \mathbf{X}]_{kn} - \sum_{l=1, l \neq k}^K [\mathbf{B}^T \mathbf{B}]_{kl} \mathbf{S}_{ln} - \eta \sum_{r=1, r \neq n}^N \mathbf{G}_{nr} \mathbf{S}_{kr}$.

Here, $[\mathbf{B}^T \mathbf{B}]_{kk} = 1$, and $\mathbf{G}_{nn} > 0$. So $f(\mathbf{S}_{kn})$ is piece-wise parabolic function that opens up. Based on the convexity and monotonic property of the parabolic function, it is not difficult to know that $f(\mathbf{S}_{kn})$ reaches the minimum at the unique point

$$\mathbf{S}_{kn} = \{\max\{\mathbf{H}_{kn}, \alpha\} + \min\{\mathbf{H}_{kn}, -\alpha\}\} / (1 + \eta \mathbf{G}_{nn}) \quad (21)$$

4.2. Learning bases

Without the sparseness regularization term in (17) and additional constraints in (18), $\mathbf{S}_{k\bullet}$ and $\mathbf{B}_{\bullet k}$ are dual in objective function $\|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2$ for $\forall k \in \{1, 2, \dots, K\}$. Hence, $\forall d \in \{1, 2, \dots, D\}$, $k \in \{1, 2, \dots, K\}$, with $\{\mathbf{B}_{pq}, p = 1, 2, \dots, D, q = 1, 2, \dots, K\} / \mathbf{B}_{dk}$ and \mathbf{S} fixed, the unconstrained single variable minimization problem of (18) has the closed-form solution

$$\mathbf{B}_{dk} = \arg \min_{\mathbf{B}_{dk}} \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 = \frac{[\mathbf{X}\mathbf{S}^T]_{dk} - \sum_{l=1, l \neq k}^K \mathbf{B}_{dl} [\mathbf{S}\mathbf{S}^T]_{lk}}{[\mathbf{S}\mathbf{S}^T]_{kk}} \quad (22)$$

while $\|\mathbf{S}_{k\bullet}\|_1 > 0$.

Since the optimal value for \mathbf{B}_{dk} does not depend on the other entries in the same column, the objective function of $\mathbf{B}_{\bullet k}$ reduces to (23) with \mathbf{S} fixed

$$f(\mathbf{B}_{\bullet k}) = [\mathbf{S}_{k\bullet} [\mathbf{S}_{k\bullet}]^T] [\mathbf{B}_{\bullet k}^T \mathbf{B}_{\bullet k}] + 2[\mathbf{B}_{\bullet k}]^T \{\tilde{\mathbf{B}}^k \mathbf{S} [\mathbf{S}_{k\bullet}]^T - \mathbf{X} [\mathbf{S}_{k\bullet}]^T\} \quad (23)$$

where

$$\tilde{\mathbf{B}}^k = \begin{cases} \mathbf{B}_{\bullet p}, & p \neq k \\ \mathbf{0}, & p = k \end{cases}$$

When imposing the norm constraint, i.e. $\|\mathbf{B}_{\bullet k}\|_2 = [\mathbf{B}_{\bullet k}]^T \mathbf{B}_{\bullet k} = 1$, (23) becomes (24)

$$f(\mathbf{B}_{\bullet k}) = 2[\mathbf{B}_{\bullet k}]^T \{\tilde{\mathbf{B}}^k \mathbf{S} [\mathbf{S}_{k\bullet}]^T - \mathbf{X} [\mathbf{S}_{k\bullet}]^T\} + \mathbf{S}_{k\bullet} [\mathbf{S}_{k\bullet}]^T \quad (24)$$

Hence, the original constrained minimization problem becomes a linear programming under a unit norm constraint,

whose solution is as follows:

$$\mathbf{B}_{\bullet k} = \frac{\mathbf{X} [\mathbf{S}_{k\bullet}]^T - \tilde{\mathbf{B}}^k \mathbf{S} [\mathbf{S}_{k\bullet}]^T}{\|\mathbf{X} [\mathbf{S}_{k\bullet}]^T - \tilde{\mathbf{B}}^k \mathbf{S} [\mathbf{S}_{k\bullet}]^T\|_2} \quad (25)$$

4.3. Convergence analysis

Assuming that $(\mathbf{B}^t, \mathbf{S}^t)$ is the result after the t th iteration, and $(\mathbf{B}^{t+1}, \mathbf{S}^{t+1})$ is the result after the $(t+1)$ th iteration. Since the exact minimum point is obtained by (21) and (25), each update operation will monotonically decrease the value of corresponding objective function. Considering that the objective function is obviously bounded below, and $f(\mathbf{B}^t, \mathbf{S}^t) \geq f(\mathbf{B}^{t+1}, \mathbf{S}^{t+1}) \geq f(\mathbf{B}^{t+1}, \mathbf{S}^{t+1})$, it converges.

4.4. Overall algorithm

Our algorithm for learning dictionary is shown in Algorithm 1.

Algorithm 1. Learning dictionary on manifolds.

Require Data matrix $\mathbf{X} \in \mathbb{R}^{D \times N}$ and K

- 1: $\mathbf{B} \leftarrow \text{rand}(D, K), \mathbf{B}_{\bullet k} = \frac{\mathbf{B}_{\bullet k}}{\|\mathbf{B}_{\bullet k}\|_2} \forall k, \mathbf{S} \leftarrow \text{zeros}(K, N)$
- 2: $\text{iter} = 0$
- 3: **while** $(f(\text{iter}) - f(\text{iter} + 1)) / f(\text{iter}) > 1e-5$ **do**
- 4: $\text{iter} \leftarrow \text{iter} + 1$
- 5: **Update S:**
- 6: Compute $\mathbf{A} = (\mathbf{B}^T \mathbf{B}) \odot (\mathbf{1} - \mathbf{I}), \mathbf{C} = \mathbf{G} \odot (\mathbf{1} - \mathbf{I})$ and $\mathbf{E} = \mathbf{B}^T \mathbf{X}$
- 7: **for** $n = 1; n \leq N; n++$ **do**
- 8: **for** $k = 1; k \leq K; k++$ **do**
- 9: $\mathbf{S}_{kn} = \{\max\{\mathbf{E}_{kn} - \mathbf{A}_{k\bullet} \mathbf{S}_{\bullet n} - \eta \mathbf{S}_{k\bullet} \mathbf{C}_{\bullet n}, \alpha\} + \min\{\mathbf{E}_{kn} - \mathbf{A}_{k\bullet} \mathbf{S}_{\bullet n} - \eta \mathbf{S}_{k\bullet} \mathbf{C}_{\bullet n}, -\alpha\}\} / (1 + \eta \mathbf{G}_{nn})$
- 10: **end for**
- 11: **end for**
- 12: **Update B:**
- 13: Compute $\mathbf{F} = (\mathbf{S}\mathbf{S}^T) \odot (\mathbf{1} - \mathbf{I}), \mathbf{W} = \mathbf{X}\mathbf{S}^T$
- 14: **for** $k = 1; k \leq K; k++$ **do**
- 15: $\mathbf{B}_{\bullet k} = \frac{\mathbf{W}_{\bullet k} - \mathbf{B}\mathbf{F}_{\bullet k}}{\|\mathbf{W}_{\bullet k} - \mathbf{B}\mathbf{F}_{\bullet k}\|_2}$
- 16: **end for**
- 17: **Update the objective function:**
- 18: $f = \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 + 2\alpha\|\mathbf{S}\|_1 + \eta \operatorname{tr}\{\mathbf{S}\mathbf{G}\mathbf{S}^T\}$
- 19: **end while**
- 20: **return** \mathbf{B} , and \mathbf{S}

Here, $\mathbf{1} \in \mathbb{R}^{K \times K}$ is a square matrix with all elements 1, $\mathbf{I} \in \mathbb{R}^{K \times K}$ is the identity matrix, and \odot indicates element dot product. By iterating \mathbf{S} and \mathbf{B} alternately, the sparse codes are obtained, and the corresponding bases are learned.

5. Implementation for image coding and spatial pooling

After learning the dictionary \mathbf{B} , we implement our algorithm on image coding and spatial pooling. First, we find the relationship between the features $\mathbf{Y} \in \mathbb{R}^{D \times M}$ in image \mathbf{J} and \mathbf{X} (features for learning dictionary), let $\mathbf{P} \in \mathbb{R}^{N \times M}$ represent this relationship. \mathbf{P} can be obtained according LLE algorithm [7] for DLSM and feature-sign search algorithm [22] for DLMM. Then, the objective function in the image coding step with dictionary fixed can be formulated as

$$\min f(\mathbf{V}) = \|\mathbf{Y} - \mathbf{B}\mathbf{V}\|_F^2 + 2\alpha\|\mathbf{V}\|_1 + \eta\|\mathbf{V} - \mathbf{S}\mathbf{P}\|_F^2 \quad (26)$$

\mathbf{V} is the sparse codes for image features \mathbf{Y} . Similar to solve (17), the solution to the minimization of (26) is as follows:

$$\mathbf{V}_{km} = (\max\{\mathbf{H}_{km}, \alpha\} + \min\{\mathbf{H}_{km}, -\alpha\}) / (1 + \eta) \quad (27)$$

where $\mathbf{H}_{km} = [\mathbf{B}^T \mathbf{Y}]_{km} + \eta [\mathbf{S} \mathbf{P}]_{km} - \sum_{l=1, l \neq k}^K [\mathbf{B}^T \mathbf{B}]_{kl} \mathbf{V}_{lm}$.

After coding each image, we split it into three levels (i.e. a three level spatial pyramid) with $L(1+4+16)$ regions. Then each region is pooled to a vector, and all L vectors are concatenated together to form the representation of the target image. The pooling strategy adopted is max. An image can be represented as follows:

$$\mathbf{Z} = \{\mathit{hist}_1^{\max T}, \mathit{hist}_2^{\max T}, \dots, \mathit{hist}_L^{\max T}\} \quad (28)$$

where $\mathit{hist}_l^{\max} = \max\{\mathbf{V}_{\bullet i}\}_{i \in O^l}$, O^l represents the features distributing in the l th region. Image coding and SPM pooling is shown in Algorithm 2.

Algorithm 2. Image coding and SPM pooling.

```

Require  $K, L, \alpha, \eta, \mathbf{Y} \in \mathbb{R}^{D \times M}, \mathbf{P} \in \mathbb{R}^{N \times M}, \mathbf{B} \in \mathbb{R}^{D \times K}, \mathbf{S} \in \mathbb{R}^{K \times N}$ 
 $\mathbf{V} \leftarrow \mathit{zeros}(K, M), \mathbf{Z} \leftarrow \mathit{zeros}(K \times L, 1)$ 
2: Image Coding:
 $iter = 0$ 
4: while  $(f(iter) - f(iter + 1)) / f(iter) > 1e-5$  do
 $iter \leftarrow iter + 1$ 
6: Update V:
Compute  $\mathbf{A} = (\mathbf{B}^T \mathbf{B}) \odot (\mathbf{1} - \mathbf{I})$  and  $\mathbf{E} = \mathbf{B}^T \mathbf{Y} + \eta \mathbf{S} \mathbf{P}$ 
8: for  $k = 1; k \leq K; k++$  do
 $\mathbf{V}_{k\bullet} = \{\max\{\mathbf{E}_{k\bullet}, -\mathbf{A}_{k\bullet} \mathbf{V}, \alpha\} + \min\{\mathbf{E}_{k\bullet}, -\mathbf{A}_{k\bullet} \mathbf{V}, -\alpha\}\} / (1 + \eta)$ 
10: end for
Update the objective function:
12:  $f = \|\mathbf{Y} - \mathbf{B} \mathbf{V}\|_F^2 + 2\alpha \|\mathbf{V}\|_1 + \eta \|\mathbf{Y} - \mathbf{S} \mathbf{P}\|_F^2$ 
end while
14: SPM Pooling:
for  $l = 1; l \leq L; l++$  do
16:  $\mathit{hist}_l^{\max} = \max\{\mathbf{V}_{\bullet i}\}_{i \in O^l}, l = 1, \dots, L$ 
end for
18:  $\mathbf{Z} = \{\mathit{hist}_1^{\max T}, \mathit{hist}_2^{\max T}, \dots, \mathit{hist}_L^{\max T}\}$ 
return Z

```

6. Experimental results

In this section, our DLSM and DLMM algorithms for image classification are evaluated on four benchmark data sets, including UIUC-Sports data set [26], Scene 15 data set [1,27,28], Caltech-101 data set [29] and Caltech-256 data set [30]. Parameter settings are given first. Then experimental results and analysis are demonstrated, and some discussions are listed finally.

6.1. Parameter settings

For each data set of image classification, the data are randomly split into the training set and the testing set based on published protocols. To make the results more convincing, the experimental process is repeated eight times, and the mean and standard deviation of the classification accuracy are recorded. Each image is resized with maximum side 300 pixels first.¹ As for the image features, densely sampling patches are extracted with the patch size and step size 16×16 and 8 pixels, respectively, and 128 dimensional SIFT descriptors [31] are obtained with grid size 4×4 .

The number of samples used for learning dictionary is about 120,000 and the dictionary size is 1024. Spatial pyramid matching kernel is embedded in the pooling step (the image is split into three layers, each of which has 1, 4, and 16 segments, respectively). The pooling strategy is max [2]. An image is represented as the concatenation of each segment with length 21,504 and normalized to 1 with l_2 -norm. Linear kernel SVM classifier and one against all multi-classification strategy are adopted, and LIBSVM [23] package is used.

There are two parameters: α and η . The parameter α is used for adjusting the sparsity of the codes. The bigger α is, the sparser the codes are. Yang et al. [2] have obtained the empirical value 0.15 for α . The parameter η is used for balancing the weight between the fitting goodness term and the geometrical structure preserving term. In this paper, α is set to 0.15 and η is set to 0.2 (for details, see Section 6.6.1).

For experiment of comparison on strong manifolds structures, the procedure is similar to the image classification in addition to the lack of pooling step. The dimension of feature is $112 \times 92 = 10,304$, and the dictionary size is 200. The parameter α is set to 0.15 and η is set to 0.8. The first five faces of each class are selected as the training set and the rest are the testing set. Linear kernel SVM classifier and one against all multi-classification strategy are adopted, and we use LIBSVM [23] package for the implementation.

6.2. UIUC-Sports data set

For UIUC-Sports data set [26], there are eight classes with totally 1579 images: rowing (250 images), badminton (200 images), polo (182 images), bocce (137 images), snow boarding (190 images), croquet (236 images), sailing (190 images), and rock climbing (194 images). For each class, the sizes of the instances in the same scene are very different, and the poses of the objects vary a lot. In addition, the background of each image is highly cluttered and discrepant. Even the number of the instances in the same category changes greatly. Some images from different classes have similar background (see in Fig. 5). We follow the common setup: 70 images per class are randomly selected as the training data, and 60 images per class for testing. Figs. 6 and 7 show the confusion matrices of our DLSM and DLMM method for image classification. Table 2² shows the performance of different methods. We notice that the classes of croquet and bocce have a high probability of being classified mistakenly, because these two classes are visually similar to each other. Furthermore, the classification rate by DLSM and DLMM is almost the same.

6.3. Scene 15 data set

For Scene 15 data set, there are 15 classes, with totally 4485 images. Each class varies from 200 to 400 images. The images contain not only indoor scenes, such as bedroom, living room, PARoffice, kitchen, and store, but also outdoor scenes, such as industrial, forest, mountain, tallbuilding, highway, street, open-country, and so on (see in Fig. 8). We use an identical experimental setup as [1], where 100 images per class are randomly selected as the training data, and the rest for testing. Figs. 9 and 10 show the confusion matrices of our DLMM and DLSM method for image classification. Table 3 lists the comparisons of our two methods with previous work for image classification.

¹ For UIUC-Sports data set, we resize the maximum side to 400 due to the high resolution of original image.

² All the results of OCSVM and HIKVQ are based on step size 8 and without concatenated Sobel images.

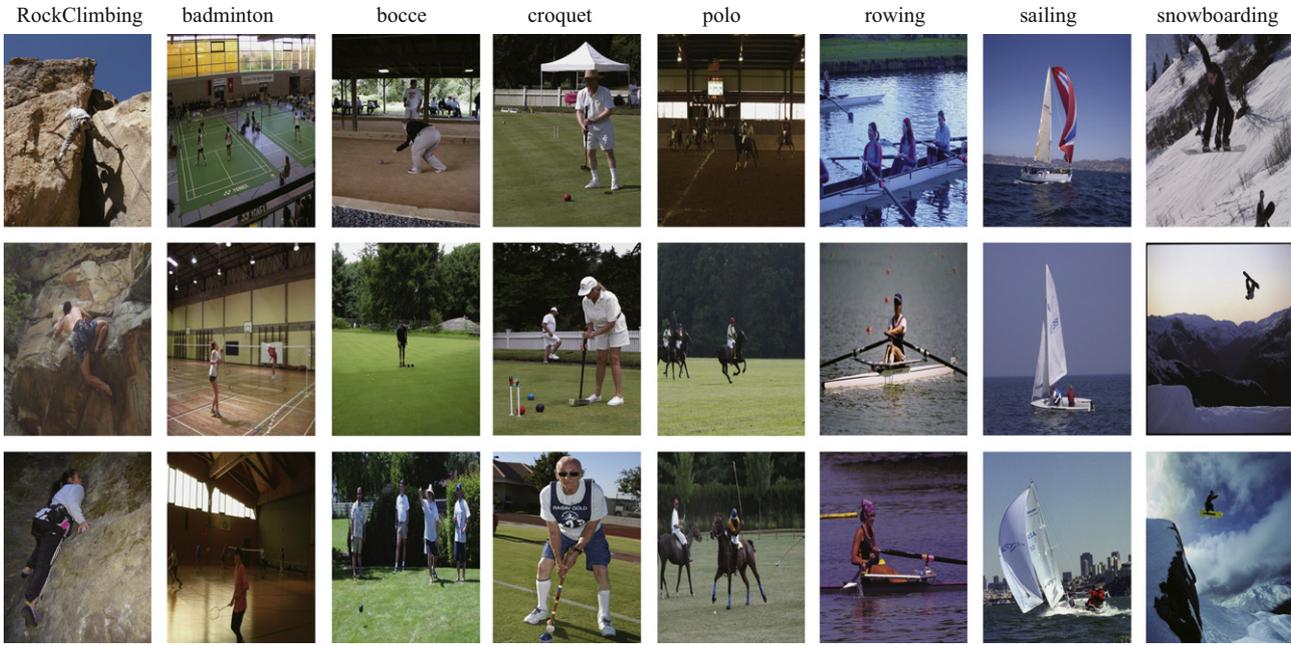


Fig. 5. Example images from the UIUC-Sports data set. For each class, the instances in the same category are very different, the pose of the objects vary a lot, and the background of each is highly clutter and discrepancy. Some images from different classes have similar background.

RockClimbing	93.54	0.00	3.33	2.92	1.25	0.83	0.00	3.75
badminton	0.00	95.21	1.67	1.04	1.46	0.42	0.00	1.04
bocce	0.63	1.67	71.46	12.92	2.71	2.50	0.63	5.00
croquet	0.21	0.83	14.79	77.29	1.88	0.83	3.33	1.04
polo	0.21	0.83	3.33	2.92	89.17	2.08	0.00	1.25
rowing	2.29	0.42	2.92	1.46	1.88	88.75	2.50	1.46
sailing	0.00	0.63	0.21	0.83	0.00	1.25	92.92	0.21
snowboarding	3.13	0.42	2.29	0.63	1.67	3.33	0.63	86.25
	RockClimbing	badminton	bocce	croquet	polo	rowing	sailing	snowboarding

Fig. 6. Confusion matrix on UIUC-Sport data set (%) by DLSSM.

RockClimbing	93.13	0.00	3.54	2.92	1.25	0.83	0.00	3.33
badminton	0.00	95.42	2.08	0.83	1.04	0.42	0.00	1.25
bocce	0.63	1.46	71.46	13.96	2.50	2.50	0.63	5.42
croquet	0.21	0.83	14.37	76.88	1.88	0.63	3.33	1.04
polo	0.42	0.83	3.12	2.50	90.00	2.08	0.00	1.04
rowing	2.29	0.42	2.92	1.46	1.88	89.17	2.50	1.46
sailing	0.00	0.63	0.21	0.83	0.00	1.25	92.92	0.21
snowboarding	3.33	0.42	2.29	0.63	1.46	3.13	0.63	86.25
	RockClimbing	badminton	bocce	croquet	polo	rowing	sailing	snowboarding

Fig. 7. Confusion matrix on UIUC-Sport data set (%) by DLMM.

Table 2
Performance comparison on UIUC-Sports data set.

Methods	Average classification rate (%)
HIKVQ [3]	81.87 ± 1.14
OCSVM [3]	81.33 ± 1.56
ScSPM [2]	82.74 ± 1.46
KScSPM [32]	84.92 ± 0.78
LScSPM [14]	85.31 ± 0.51
DLSM	86.82 ± 1.04
DLMM	86.93 ± 0.99

6.4. Caltech-101

Caltech-101 data set introduced in [29] contains 102 classes, one of which is the background. After removing the background class, the rest 101 classes with totally 8677 images are used for classification, with each class varying from 31 to 800 images. We follow the common experimental setup for this data set, where 15 and 30 images per category are selected as the training set, and the rest for the testing set (for training on 15 images per category, the maximum is 20 images per category for testing, and for training on 30 images per category, the maximum is 50 images



Fig. 8. Example images from the Scene 15 data set. The images contains not only indoor scene, but also outdoor scene.

suburb	98.76	0.00	0.33	0.08	0.42	0.00	0.81	0.00	0.00	0.00	0.00	0.06	0.00	0.07	0.00
coast	0.00	90.19	0.11	3.44	0.30	1.73	16.98	0.00	0.29	0.00	0.00	0.18	0.00	0.00	0.00
forest	0.00	0.53	92.05	0.00	0.00	1.14	3.79	0.00	0.15	0.00	0.00	0.06	0.00	0.00	0.23
highway	0.00	2.31	0.00	89.77	0.12	0.68	3.15	2.34	0.00	0.00	0.00	0.36	0.00	0.00	0.00
insidicity	0.00	0.00	0.00	1.72	85.82	0.05	0.00	3.26	3.52	0.76	0.43	1.54	0.34	0.00	3.20
mountain	0.00	0.82	4.33	0.63	0.00	90.97	4.40	0.26	0.93	0.00	0.11	0.53	0.57	0.40	2.67
opencountry	0.00	5.14	1.04	1.56	0.00	3.33	68.15	0.00	0.15	0.00	0.00	0.06	0.00	0.00	0.00
street	0.00	0.05	0.77	0.47	3.31	0.23	0.69	90.63	0.05	0.00	0.11	0.77	0.00	0.13	1.05
tallbuilding	0.00	0.24	0.27	0.70	4.69	1.09	0.04	1.69	91.55	0.00	0.11	2.19	0.00	0.00	1.22
office	0.27	0.19	0.00	0.00	0.00	0.00	0.00	0.00	0.00	92.17	2.59	0.71	3.52	1.92	0.99
bedroom	0.09	0.05	0.00	0.08	0.48	0.09	0.12	0.13	0.05	0.98	76.83	1.66	5.68	18.06	1.69
industrial	0.18	0.48	0.55	0.86	2.04	0.55	1.01	1.04	2.05	1.41	3.02	77.67	2.16	4.17	8.02
kitchen	0.00	0.00	0.00	0.00	0.96	0.00	0.16	0.00	0.05	2.61	6.57	2.31	74.43	9.72	5.47
livingroom	0.53	0.00	0.00	0.08	0.24	0.00	0.20	0.13	0.20	1.41	9.27	1.18	10.00	61.38	4.88
store	0.18	0.00	0.55	0.63	1.62	0.14	0.52	0.52	1.03	0.65	0.97	10.72	3.30	4.17	70.58
	suburb	coast	forest	highway	insidicity	mountain	opencountry	street	tallbuilding	office	bedroom	industrial	kitchen	livingroom	store

Fig. 9. Confusion matrix on Scene 15 data set (%) by DLSM.

suburb	98.76	0.00	0.33	0.08	0.48	0.00	0.73	0.00	0.00	0.00	0.00	0.06	0.00	0.07	0.00
coast	0.00	89.86	0.11	3.28	0.24	1.73	16.69	0.00	0.29	0.00	0.00	0.24	0.00	0.00	0.00
forest	0.00	0.53	92.05	0.08	0.00	1.14	3.55	0.00	0.15	0.00	0.00	0.06	0.00	0.00	0.17
highway	0.00	2.21	0.00	90.31	0.06	0.68	3.02	2.47	0.00	0.00	0.00	0.36	0.00	0.00	0.00
insidicity	0.00	0.00	0.00	1.80	86.72	0.05	0.00	3.26	3.52	0.87	0.43	1.72	0.34	0.00	3.37
mountain	0.00	0.82	4.17	0.55	0.00	91.01	4.72	0.33	0.93	0.00	0.22	0.59	0.57	0.40	2.62
opencountry	0.00	5.58	1.15	1.33	0.00	3.33	68.51	0.00	0.05	0.00	0.11	0.12	0.00	0.00	0.00
street	0.00	0.05	0.82	0.39	3.25	0.23	0.73	90.49	0.10	0.00	0.11	0.89	0.00	0.26	0.93
tallbuilding	0.00	0.24	0.22	0.70	4.63	1.00	0.04	1.69	91.65	0.00	0.22	2.25	0.00	0.00	1.16
office	0.27	0.19	0.00	0.00	0.00	0.00	0.00	0.00	0.00	93.15	2.80	0.65	3.52	2.12	0.93
bedroom	0.09	0.05	0.00	0.08	0.42	0.09	0.12	0.13	0.10	0.87	77.16	1.72	5.11	18.78	1.51
industrial	0.18	0.48	0.49	0.70	1.62	0.59	1.01	0.98	2.00	1.30	2.80	77.37	2.16	3.64	8.02
kitchen	0.00	0.00	0.00	0.00	0.84	0.00	0.16	0.00	0.05	1.96	6.36	2.19	75.34	9.06	5.29
livingroom	0.53	0.00	0.00	0.08	0.24	0.00	0.24	0.13	0.20	1.30	8.84	1.13	9.89	61.57	4.83
store	0.18	0.00	0.66	0.63	1.50	0.14	0.48	0.52	0.98	0.54	0.97	10.66	3.07	4.10	71.16

Fig. 10. Confusion matrix on Scene 15 data set (%) by DLMM.

Table 3

Performance comparison on scene 15 data set.

Methods	Average classification rate (%)
KSPM [1]	81.4 ± 0.5
KC [11]	76.7 ± 0.4
HIKQ [3]	81.77 ± 0.49
OCSVM [3]	82.02 ± 0.54
ScSPM [2]	80.28 ± 0.93
KScSPM [32]	83.68 ± 0.61
LScSPM [14]	89.75 ± 0.50
DLSM	83.40 ± 0.44
DLMM	83.67 ± 0.49

Table 5

Performance comparison on Caltech-256 data set.

Methods	15 training	30 training	45 training	60 training
LLC ^a [12]	28.00 ± 0.36	33.34 ± 0.57	36.24 ± 0.37	38.08 ± 0.39
ScSPM [2]	27.73 ± 0.51	34.02 ± 0.35	37.46 ± 0.55	40.14 ± 0.91
KScSPM [32]	29.77 ± 0.14	35.67 ± 0.10	38.61 ± 0.19	40.30 ± 0.22
LScSPM [14]	30.00 ± 0.14	35.74 ± 0.10	38.54 ± 0.36	40.43 ± 0.38
DLSM	29.31 ± 0.58	35.12 ± 0.34	37.62 ± 0.57	39.96 ± 0.62
DLMM	30.35 ± 0.42	36.22 ± 0.33	38.97 ± 0.56	41.09 ± 0.44

^a For LLC, we adopt the code of local feature coding provided by [12] and do experiment on our data set with single scale features and the size of dictionary 1024.

Table 4

Performance comparison on Caltech-101 data set.

Methods	15 training	30 training
KSPM [1]	–	64.6 ± 0.8
KC [11]	–	64.1 ± 1.2
LLC ^a [12]	63.92 ± 0.46	70.63 ± 0.99
ScSPM [2]	67.0 ± 0.45	73.2 ± 0.54
DLSM	66.88 ± 0.53	74.39 ± 0.82
DLMM	67.54 ± 0.41	74.87 ± 0.67

^a For LLC, we adopt the code of local feature coding provided by [12] and do experiment on our data set with single scale features and the size of dictionary 1024.

per category for testing). Table 4 shows the performance of different methods.

6.5. Caltech-256

Caltech-256 data set introduced in [30] contains 257 classes, one of which is the background. After removing the background class, the rest 256 classes with totally 29,780 images are used for classification. Due to much higher intra-class variability and higher object location variability compared with Caltech-101, Caltech-256 is a very challenging data set so far for image classification. We follow the common experimental setup for this data set: 15, 30, 45 and 60 training images per category and 15

testing images per category. Table 5 shows the performance of different methods.

6.6. Experiments revisit

In this section, we revisit the experiments, and give the best parameter selection for balancing the weight among the fitting goodness, sparsity and the geometrical structure at first. Then, the comparison of DLSM and DLMM is also presented.

6.6.1. Parameter selection

The parameter α is the regularization parameter to control the tradeoff between the fitting goodness and the sparseness. With the increasing value of α , the codes become sparser and sparser, more and more salient, easy to distinguish; on the contrary, the reconstruction error becomes larger and larger, which will lead to inaccurate description of the codes. The parameter η is the regularization parameter to balance between the fitting goodness and the geometrical structure. With the increasing value of η , the geometrical structure becomes more and more reliable; the codes becomes more and more reasonable, and the reconstruction error becomes larger and larger. These two parameters are the competitors of reconstruction error.

First, we studied the effect of different η for DLMM and DLSM. Fig. 11(a) lists the performance when $\eta = \{0.0, 0.1, 0.2, 0.3, 0.4, 0.6, 0.8, 1.0\}$ with $\alpha = 0.15$. As can be seen, when $\eta = 0.0$, DLSM

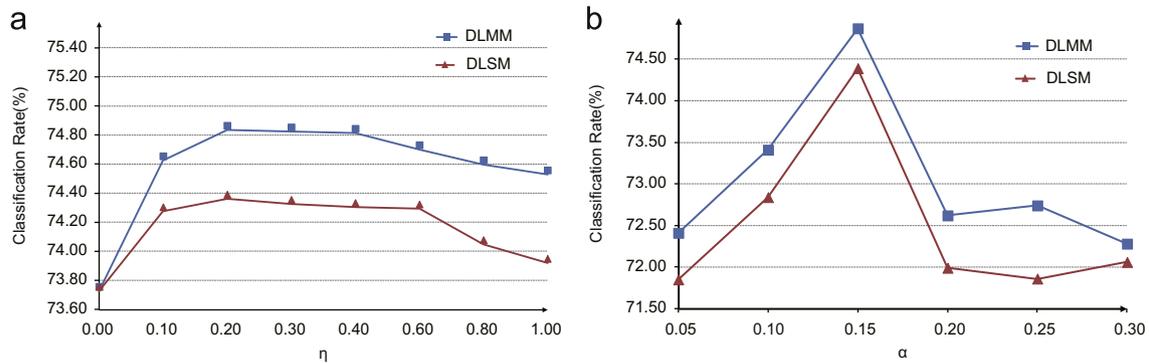


Fig. 11. Selection of parameters. (a) The classification rate under different η (Caltech-101) with $\alpha = 0.15$ and (b) the classification rate under different α (Caltech-101) with $\eta = 0.2$.

and DLMM are equivalent to sparse coding based dictionary learning method. With η growing, the classification rate increases. The best classification accuracy can be obtained when $\eta = 0.2$. After that, the performance starts to degenerate. When applying $\eta = 0.2$ to other data sets, we can also achieve the optimal results. Then, following the steps to obtain the optimal η , we fix $\alpha = \{0.05, 0.1, 0.15, 0.2, 0.25, 0.3\}$ with $\eta = 0.2$, and obtain the optimal $\alpha = 0.15$ (see Fig. 11(b)).

6.6.2. Comparison of DLSM and DLMM

According to the experimental results, we can see that, for the small size of data set, both DLSM and DLMM can achieve satisfying results, such as UIUC-sports data set and Scene 15 data set. However, with the increase of the size of the data and the number of the category, DLMM method becomes superior to DLSM. This is because, with the increase of the size of the data and the number of the category, the complexity of image patches is raising, so that single manifold cannot fit the cases well.

With the increasing number of the training data, the improvement of DLSM decreases faster than DLMM compared with ScSPM. When 60 training data per category is adopted, the classification accuracy is lower than ScSPM, which means that the single manifold cannot represent the real image features. Co-incidentally, in [14], the author also concluded that as the number of training data increases, the improvement of LScSPM also decreases. When 60 training data per category is adopted, the classification rate merely exceeded 0.29%. However, our DLMM method exceeded 0.95%.

6.6.3. Comparisons on strong manifolds structures

To evaluate the effect of k -neighbors, we carry out an experiment on ORL face databases with strong manifolds structures. The ORL face data set contains 10 images for each of 40 human faces with each lying on the same manifold. The mean classification accuracy for each class is recorded (see in Table 6). From Table 6, the classification accuracy for the traditional sparse coding is 93%. The best classification accuracy for DLMM is 97%. And the best classification accuracy for DLSM is 95.5%. The classification accuracy for DLSM is high, because all the human faces are very similar, and single manifold may be approaching smooth. Obviously, the multiple manifolds structure is more reasonable than the single manifold structure.

From the comparison of k -neighbors, we can see that for DLMM, the classification accuracy is stable except when the number of neighbors is set to 5, because 5-neighborhood is too small to construct the multiple manifolds structures. For DLSM, the classification accuracy is unstable. When the number of neighbors is set to a large number, such as 150 or 199, the classification rate is less than the traditional sparse coding algorithm.

Table 6
Comparison on ORL face data set among different k -neighbors.

k -Neighbors	5	10	15	20	30	40	50
DLMM (%)	92.5	94	94	95.5	97	96	95.5
DLSM (%)	94	94	93.5	94	92.5	94.5	95
SC (%)	93						
k -neighbors	60	70	80	90	100	150	199
DLMM (%)	94.5	95.5	95	95	95	95	95.5
DLSM (%)	95.5	95	93.5	93.5	93.5	91.5	91.5

7. Conclusion

In this paper, we have proposed dictionary learning methods on manifolds for image classification. The methods consider the intrinsic geometrical structures identified by locally linear embedding on single manifold and sparse representation on multiple manifolds, respectively. Then the geometrical structures are embedded into sparse coding algorithm in a unified framework so as to preserve them during the sparse coding procedure. After that, a coordinate descent scheme is proposed to solve the optimization subproblems. The uncovered manifold structure makes it more congruent with image classification task. Experimental results on four benchmark data sets demonstrate that our proposed algorithm leads to more effective image representation and gains better classification performance. For the future, learning kernel based dictionary on manifolds will be carried out. And more rational topology will also be explored to improve the performance of image classification.

Acknowledgments

This work has been supported by the National Natural Science Foundation under Grants NNSF-61171118 and the Ministry of Education under Grants SRFDP-20110002110057.

We would like to thank Qianhaoze You for the suggestion on revision of this manuscript.

References

- [1] A. Goldberg, X. Zhu, A. Singh, Z. Xu, R. Nowak, Multi-manifold semi-supervised learning, in: Proceedings of 12th International Conference on Artificial Intelligence and Statistics, 2009.
- [2] S. Roweis, L. Saul, Nonlinear dimensionality reduction by locally linear embedding, Science 290 (5500) (2000) 2323.
- [3] C. Chang, C. Lin, LIBSVM: a library for support vector machines, ACM Transactions on Intelligent Systems and Technology 2 (2011) 27:1–27:27.

- [4] J. Sivic, A. Zisserman, Video Google: a text retrieval approach to object matching in videos, in: Proceedings of the Ninth International Conference on Computer Vision, IEEE, vol. 2, 2003, pp. 1470–1477.
- [5] J. Yang, S. Yang, Y. Fu, X. Li, T. Huang, Non-negative graph embedding, in: Proceedings of the 21st International Conference on Computer Vision and Pattern Recognition, IEEE, 2008, pp. 1–8.
- [6] J. Wu, J. Rehg, Beyond the euclidean distance: creating effective visual codebooks using the histogram intersection kernel, in: Proceedings of the 12th International Conference on Computer Vision, IEEE, 2009, pp. 630–637.
- [7] M. Zheng, J. Bu, C. Chen, C. Wang, L. Zhang, G. Qiu, D. Cai, Graph regularized sparse coding for image representation, IEEE Transactions on Image Processing 20 (5) (2011) 1327–1336.
- [8] J. Lee, Introduction to Smooth Manifolds, vol. 218, Springer Verlag, 2003.
- [9] J. Tenenbaum, V. De Silva, J. Langford, A global geometric framework for nonlinear dimensionality reduction, Science 290 (5500) (2000) 2319–2323.
- [10] J. Yang, K. Yu, Y. Gong, T. Huang, Linear spatial pyramid matching using sparse coding for image classification, in: Proceedings of the 22nd International Conference on Computer Vision and Pattern Recognition, IEEE, 2009, pp. 1794–1801.
- [11] L. Fei-Fei, P. Perona, A bayesian hierarchical model for learning natural scene categories, in: Proceedings of the 18th International Conference on Computer Vision and Pattern Recognition, IEEE, vol. 2, 2005, pp. 524–531.
- [12] M. Belkin, P. Niyogi, Laplacian eigenmaps and spectral techniques for embedding and clustering, Advances in Neural Information Processing Systems 14 (2001) 585–591.
- [13] B. Shen, L. Si, Nonnegative matrix factorization clustering on multiple manifolds, in: Proceedings of the 24th AAAI Conference on Artificial Intelligence, 2010, pp. 575–580.
- [14] J. van Gemert, C. Veenman, A. Smeulders, J. Geusebroek, Visual word ambiguity, IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (7) (2010) 1271–1283.
- [15] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, Y. Gong, Locality-constrained linear coding for image classification, in: Proceedings of the 23rd International Conference on Computer Vision and Pattern Recognition, IEEE, 2010, pp. 3360–3367.
- [16] K. Ramamurthy, J. Thiagarajan, A. Spanias, Improved sparse coding using manifold projections, in: Proceedings of the 18th IEEE International Conference on Image Processing, IEEE, 2011, pp. 1237–1240.
- [17] K. Yu, T. Zhang, Y. Gong, Nonlinear learning using local coordinate coding, Advances in Neural Information Processing Systems 22 (2009) 2223–2231.
- [18] D. Lowe, Distinctive image features from scale-invariant keypoints, International Journal of Computer Vision 60 (2) (2004) 91–110.
- [19] S. Yan, D. Xu, B. Zhang, H. Zhang, Q. Yang, S. Lin, Graph embedding and extensions: a general framework for dimensionality reduction, IEEE Transactions on Pattern Analysis and Machine Intelligence 29 (1) (2007) 40–51.
- [20] Z. Lu, Y. Peng, Latent semantic learning by efficient sparse coding with hypergraph regularization, in: Proceedings of the 25th AAAI Conference on Artificial Intelligence, 2011.
- [21] S. Gao, I. Tsang, L. Chia, Kernel sparse representation for image classification and face recognition, in: Proceedings of the 11th European Conference on Computer Vision, Springer, 2010, pp. 1–14.
- [22] A. Oliva, A. Torralba, Modeling the shape of the scene: a holistic representation of the spatial envelope, International Journal of Computer Vision 42 (3) (2001) 145–175.
- [23] L. Fei-Fei, R. Fergus, P. Perona, Learning generative visual models from few training examples: an incremental bayesian approach tested on 101 object categories, in: Workshop of the 17th International Conference on Computer Vision and Pattern Recognition, vol. 12, 2004, p. 178.
- [24] S. Gao, I. Tsang, L. Chia, P. Zhao, Local features are not lonely—Laplacian sparse coding for image classification, in: Proceedings of the 23rd International Conference on Computer Vision and Pattern Recognition, IEEE, 2010, pp. 3555–3561.
- [25] F. Jurie, B. Triggs, Creating efficient codebooks for visual recognition, in: Proceedings of the 10th International Conference on Computer Vision, 2005, pp. 604–610.
- [26] S. Lazebnik, C. Schmid, J. Ponce, Beyond bags of features: spatial pyramid matching for recognizing natural scene categories, in: Proceedings of the 19th International Conference on Computer Vision and Pattern Recognition, IEEE, vol. 2, 2006, pp. 2169–2178.
- [27] Y. Boureau, F. Bach, Y. LeCun, J. Ponce, Learning mid-level features for recognition, in: Proceedings of the 23rd International Conference on Computer Vision and Pattern Recognition, IEEE, 2010, pp. 2559–2566.
- [28] H. Lee, A. Battle, R. Raina, A.Y. Ng, Efficient sparse coding algorithms, Advances in Neural Information Processing Systems 19 (2007) 801–808.
- [29] L. Li, L. Fei-Fei, What, where and who? Classifying events by scene and object recognition, in: Proceedings of the 11th International Conference on Computer Vision, IEEE, 2007, pp. 1–8.
- [30] D.P. Bertsekas, Nonlinear Programming, Athena Scientific, Belmont, MA, 1999.
- [31] G. Griffin, A. Holub, P. Perona, Caltech-256 Object Category Dataset, Technical Report 7694, California Institute of Technology, 2007.
- [32] B. Liu, Y. Wang, Y. Zhang, Y. Zheng, Discriminant sparse coding for image classification, in: Proceedings of the 37th International Conference on Acoustics, Speech and Signal Processing, IEEE, 2012, pp. 2193–2196.

Bao-Di Liu was born in Shandong, China. He received his bachelor degree and master degree in China University of Petroleum in 2004 and 2007 respectively. Currently, he is a PhD student from the Department of Electronic Engineering, Tsinghua University, China. His research interests include computer vision and machine learning.

Yu-Xiong Wang received the BS degree in Electronic Engineering from Beijing Institute of Technology (BIT), Beijing, China, in 2009. Currently, he is a master graduate student in the Department of Electronic Engineering at Tsinghua University, Beijing, China. His research interests include image processing, computer vision and machine learning.

Yu-Jin Zhang received the PhD degree in Applied Science from the State University of Liège, Liège, Belgium, in 1989. From 1989 to 1993, he was post-doc fellow and research fellow with the Department of Applied Physics and Department of Electrical Engineering at the Delft University of Technology, Delft, the Netherlands. In 1993, he joined the Department of Electronic Engineering at Tsinghua University, Beijing, China, where he is a professor of Image Engineering (since 1997). He is an active researcher in Image Engineering, with current interests on object segmentation from images and video, segmentation evaluation and comparison, moving object detection and tracking, face recognition, facial expression detection/classification, content-based image and video retrieval, information fusion for high-level image understanding, etc. He has authored more than 20 books and published more than 400 papers in the areas of image processing, image analysis, and image understanding. Professor Zhang is vice president of China Society of Image and Graphics and director of academic committee of the Society, and is a Fellow of SPIE.

Bin Shen is a PhD student in Computer Science, Purdue University. His research interests include machine learning and its applications to data mining, computer vision. Before that, he got his BS and MS in 2007 and 2009 respectively, both from EE, Tsinghua University, China.