

# Scene Recognition by Manifold Regularized Deep Learning Architecture

Yuan Yuan, *Senior Member, IEEE*, Lichao Mou, and Xiaoqiang Lu

**Abstract**—Scene recognition is an important problem in the field of computer vision, because it helps to narrow the gap between the computer and the human beings on scene understanding. Semantic modeling is a popular technique used to fill the semantic gap in scene recognition. However, most of the semantic modeling approaches learn shallow, one-layer representations for scene recognition, while ignoring the structural information related between images, often resulting in poor performance. Modeled after our own human visual system, as it is intended to inherit humanlike judgment, a manifold regularized deep architecture is proposed for scene recognition. The proposed deep architecture exploits the structural information of the data, making for a mapping between visible layer and hidden layer. By the proposed approach, a deep architecture could be designed to learn the high-level features for scene recognition in an unsupervised fashion. Experiments on standard data sets show that our method outperforms the state-of-the-art used for scene recognition.

**Index Terms**—Deep architecture, machine learning, manifold kernel, manifold regularization, scene recognition.

## I. INTRODUCTION

**S**CENE recognition, the process of categorizing images into different bins (e.g., coast, highway, street, bedroom, and store), is a challenging problem that is of importance in the field of computer vision. It is helpful to reduce the gap between computers and humans when acquiring an understanding for a scene. Thus, scene recognition has an abundance of applications in fields of computer vision and multimedia. Majority of the early approaches [1]–[3] focused on finding a mapping relating a set of low-level features (e.g., color histogram, Local Binary Pattern, and Scale-invariant Feature Transform (SIFT)) to meaningful semantic categories. However, a large semantic gap between the low-level vision features and high-level semantic categories still exist (Fig. 1).

To overcome the gap, semantic modeling was proposed for scene recognition. The semantic modeling focuses on

Manuscript received August 30, 2013; revised March 31, 2014 and June 26, 2014; accepted September 11, 2014. Date of publication January 22, 2015; date of current version September 16, 2015. This work was supported in part by the National Basic Research Program (973 Program) of China under Grant 2012CB719905 and in part by the National Natural Science Foundation of China under Grant 61172143 and Grant 61100079.

The authors are with the Center for OPTical IMagery Analysis and Learning (OPTIMAL), State Key Laboratory of Transient Optics and Photonics, Xi'an Institute of Optics and Precision Mechanics, Chinese Academy of Sciences, Xi'an 710119, Shaanxi, P. R. China (e-mails: yuany@opt.ac.cn; moulitchao@opt.ac.cn; luxq66666@opt.ac.cn).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2014.2359471

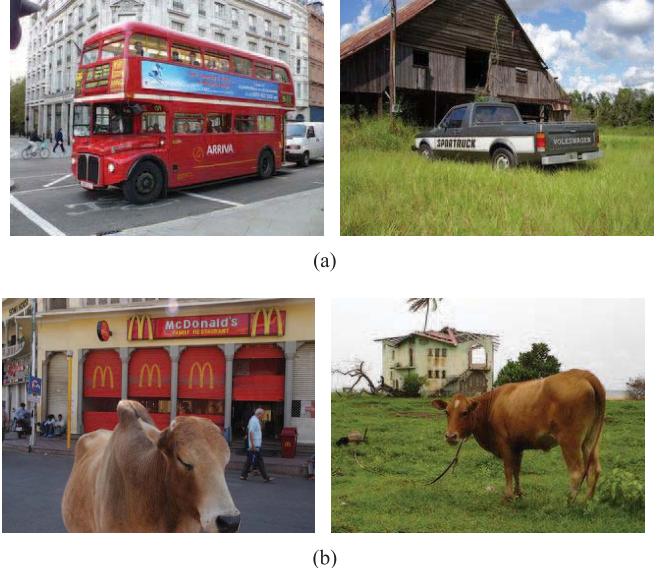


Fig. 1. Semantic gap. (a) City and countryside scenes both contain *vehicle* and *building*. (b) City and countryside scenes both contain *cow* and *building*. Each pair of images in (a) and (b) share similar objects, but their scene categories (labels) are totally different.

finding an intermediate semantic representation that bridges this gap. The semantic modeling can be roughly categorized into three categories: 1) object-based approaches; 2) bag-of-visual-features (BoF) approaches; and 3) attribute-based approaches. Object-based approaches [4]–[8] exploit a set of materials or objects that appear in the images as intermediate semantic representation, therefore, revealing the semantic gap. In general, object-based approaches can be described to the following steps. First, different regions of images can be segmented. Second, local classifiers are used to label each region according to an object it belongs to. Lastly, the global scene is classified using the object information. The *BoF* approaches [9]–[12] stem from the *bag-of-words* (BoW), having been proved successful with text analysis. BoF-based approaches construct an orderless collection of visual words extracted from images at nodes evenly spaced in a gridlike manner. The image can then be characterized by the frequency of each visual word. BoF-based approaches can learn an intermediate level representation by describing images as abstract, high semantic level. For BoF-based approaches [9]–[12], although the modeling power of local descriptors, such as SIFT, can be exploited, the feature coding step includes a fixed and flat single-layer operation. The

flat structure limits the representation power of data, and fails to adapt to different data because of the lack of learning in BoF-based approaches. In addition, the process of characterizing an image with a histogram of visual words is suboptimal [9]–[12]. The discriminative power of the local descriptors gets greatly reduced by the coarse quantization introduced as a predefined visual vocabulary [13]. Recently, attribute-based approaches have achieved some success with scene recognition. This is due to the interesting properties found in [14]–[17] (e.g., carrying abundant semantic meanings [16], powerful representation for visual recognition task [17], strong visuality [14], among other similar approaches). Attribute-based approaches set out to describe an image with a set of meaningful, visual attributes (e.g., arm moving forward, animal, etc.). For example, as done in [16], an image is reshaped to a 2659-dimension vector, and is represented as a set of visual attributes. In [17], both discriminative and nameable visual attributes are learnable in a semisupervised fashion. However, the attribute-based approaches require large amounts of manned effort to define visual attributes. In addition, the classification is out performed by the simple BoF histogram approach used on the same low-level features. Although object-based approaches and BoF-based approaches are popular for scene recognition, there are some problems with their learning process. Either approach needs amount of labeled images. It is a tough task to label efficiently kinds of images. Therefore, it is important to design an unsupervised feature learning approach for scene recognition. Moreover, benefitted from recent advances in machine learning, there are many techniques such as feature construction [18], feature learning [19], [20], subspace learning [21]–[23], and manifold learning [24], that provide potential for visual recognition task, while helping to improve performance.

Recently, computer architects have been exploited in the field of computer vision because of the success in simulating the human brain. Amongst these models, deep architectures and the related deep learning approaches attracted attention in many fields with good performance, such as, [25]–[30], etc. Deep architectures try to learn a hierarchical structure, as simple features are learned for the lower layers and more complex features for the higher layers. In other words, deep architectures attempt to learn simple concepts first, to then generate more complex concepts. Moreover, lots of work recently in machine learning focused on learning good feature representations from unlabeled input data used for higher-level tasks, such as classification. Current solutions generally learn multilevel representations by deep learning [26], [27], or other unsupervised learning algorithm [18].

Despite some exciting results with image recognition from the related deep learning approaches, the performance of recognition can be further improved with the structural information of data considered. In this paper, a new deep architecture framework is presented and preserved the structural information between images. The main objective of this paper is to encode powerful local descriptors, such as SIFT, using a manifold regularized deep architecture for scene recognition. This is contrary to other deep networks that

learn representations from pixels. Our approach is to study the construction of a manifold regularized deep architecture to learn from local descriptors as a starting point with a greater representational power than raw pixels. Unlike the previous deep learning methods, the proposed framework aims to capture the nonlinear structure of the data by incorporating the sparse regularizer into the kernel manifold space. First, the structure in a set of images cannot be well represented by exploiting a single-layer model [31]. And, layer-by-layer learning model was introduced to improve the learning performance in scene recognition. We consider the previous layer as the basic unit for learning the next hidden space layers. Second, in the basic unit of each layer, the sparse weight matrix can be learned to improve the learning performance in the manifold kernel space, which is generated by incorporating the manifold structure into the kernel space. Finally, The processing hidden space layers that follow can be learned from the proceeding activities, which can be steered according to the data of the previous layer. In this case, layer-by-layer deep learning architecture can progressively reveal low-dimensional nonlinear structure. The contributions of this paper are as follows.

- 1) We propose a novel, improved multilayer learning model to overcome the limitation pointed out in [31], which limits the model of the structure to single-layer learning. This is done by modeling the higher order correlations between data points in the layer below.
- 2) We propose deep architecture as a novel way of bettering the original deep learning by incorporating the geometrical structure of the data into the kernel space.
- 3) We learn the next hidden space layers from the previous layer's activities in this paper, this is defined as a basic unit. In each basic unit, a sparse weighted graph can be adaptively learned in the manifold kernel space. Whereas, the hidden space layer can preserve the clear structure of data, and further improve the embedding. Hence, layer-by-layer deep learning architecture is an effective way to progressively reveal low-dimensional, nonlinear structures.

The rest of this paper is organized as follows. Section II briefly reviews related work in deep learning. Section III presents our deep architecture for scene recognition. To validate the proposed approach, the experimental results for various image sets are reported in Section IV. Section V concludes this paper.

## II. RELATED WORK ON DEEP LEARNING

In this section, we briefly review some classic deep architectures and related deep learning approaches, as we later compare their architectures with the proposed model.

The performance of deep architecture has been notable, and especially after the appearance of *deep belief networks* (DBNs). The learning process for DBNs is made up of two stages: 1) extracting feature layer by layer and 2) fine-tuning the model as a whole, to improve performance. After these two stages, DBNs can combine the low-level features to form more efficient high-level features by unsupervised learning. In the first stage, DBNs implement a family

of RBMs that attempt to extract feature layer by layer. RBMs are two-layer, undirected graphical model. This model has symmetric connections between the hidden layer and the visible layer, no connections exist within the hidden layer or the visible layer. The energy function of RBMs can be defined as follows:

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}) &= \frac{1}{2} \sum_i v_i^2 - \left( \sum_{i,j} v_i w_{i,j} h_j + \sum_i c_i v_i + \sum_j b_j h_j \right) \\ &= \frac{1}{2} \mathbf{v}^T \mathbf{v} - (\mathbf{h}^T \mathbf{W} \mathbf{v} + \mathbf{c}^T \mathbf{v} + \mathbf{b}^T \mathbf{h}) \end{aligned} \quad (1)$$

where  $\mathbf{v} = (v_1, v_2, \dots, v_N)$  is visible random variable,  $\mathbf{h} = (h_1, h_2, \dots, h_M)$  is the binary hidden random variable. The joint probability distribution of the RBMs is defined as follows:

$$P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})) \quad (2)$$

where  $Z$  is a normalization constant. The form of Restricted Boltzmann Machine (RBM) makes the conditional probability distributions a simple calculating, when  $\mathbf{v}$  or  $\mathbf{h}$  are fixed.

The power of the RBM is limited to the extent that itself can be represented by. However, its real power emerges when the RBMs are stacked, forming a DBN. Hinton *et al.* [26] proposed a greedy approach that trains RBM in each layer to train efficient DBN. In general, once a layer of DBN is trained, the parameters of this layer are frozen and the values of hidden layer are inferred. The inferred values are used to train the proceeding layer that is one higher than the respective DBN. The greedy layerwise training approach has been proved effective when training deep architectures, like DBNs. Recently, most deep learning approaches are unsupervised feature learning, such as [32]–[34], which have seen an increase of interest as a result.

Although DBNs were a success, both RBMs and DBNs ignore an image's 2-D structural information. Recently, deep convolutional architectures [27], [35] have attracted amount attention owing to preserving the space structure and having the tendency insensitive to small variations in the image. All of the layers in *deep convolutional neural network* (DCNN) are trained from data in an integrated fashion. Currently, DCNN has been successfully used to extract features in different applications, such as image classification [36], and scene labeling [37], [38]. Donahue *et al.* [36] focus on learning features that have sufficient representational power and generalization ability by leveraging an auxiliary large labeled object database to train a deep convolutional architecture for classification task. Farabet *et al.* [37] propose an approach for scene labeling that uses a multiscale convolutional network trained from raw pixels to extract dense feature vectors that encode regions of multiple sizes centered on each pixel, while Pinheiro and Collobert [38] propose to use a *recurrent convolutional neural network* (RCNN) for scene labeling, and report promising performances. However, these deep models ignore the structural information between data.

### III. MANIFOLD REGULARIZED DEEP ARCHITECTURE

In this section, a novel kernel embedded deep architecture and its corresponding deep learning approach are proposed. Here, structural information is emphasized. All features used in scene recognition are learned by the proposed deep architecture in an unsupervised manner. The motivation of our approach is to design a deep architecture which can learn more useful high-level features from low-level features for scene recognition task, are without supervised. In this paper, we focus on the unsupervised feature learning.

#### A. Deep Architecture Design

Suppose there are  $n$  sample images with low-level features  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ . Given  $\mathbf{X}$ , the weight matrix  $\mathbf{W}$ , representing the local structural information of  $\mathbf{X}$ , which can be learned by the proposed deep architecture. And then, the hidden space  $\mathbf{Y}$  is learned according to the structure information represented by a learned weight matrix  $\mathbf{W}$ .  $\mathbf{Y}$  is defined as a set of activities  $\mathbf{y}_j$ , that is,  $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$ . Fig. 2 shows the structure of one base unit and the proposed deep architecture with one visible input layer and  $L$  hidden layers.  $\mathbf{X}$  is the original data, and  $\mathbf{Y}_k(k = 1, 2, \dots, L)$  is the output of  $k$ th hidden space layer.  $\mathbf{Y}_k$  then serves as the input data for another base unit that generates the new hidden space  $\mathbf{Y}_{k+1}$ . This process is repeated until the response of the final hidden layer  $\mathbf{Y}_L$  is computed, now  $L$  as the number of layers. To obtain more clear structure between data, which is helpful doing classification, the proposed approach tries to discover a model that uses several hidden layers. The next hidden space layers can be learned from the previous activities, which can be driven by the learned data from the previous layer. This layer-by-layer learning should be a performance improvement on the model assigns to the training data at each iteration.

As shown in Fig. 2, a family of base units are stacked layer by layer to form the proposed deep architecture. Two questions to be addressed: 1) how to describe the structural information between data in this deep architecture and 2) how to map the data into a hidden space based on the learned graph.

To learn the structural information between data, in a natural way, local geometry information learned by using *locally linear embedding* (LLE). LLE is a nonlinear dimension reduction approach and it supposes each data point and its neighbors lying on or close to a locally linear manifold, to determine the weight coefficients matrix  $\mathbf{W}$  constructed. The structure information can be described using LLE, which is as follows:

$$\min_{\mathbf{W}} \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} W_{ij} \mathbf{x}_j \right\|^2 \quad (3)$$

where the weight coefficients matrix  $\mathbf{W}$  is composed of  $W_{ij}$  and  $\sum_j W_{ij} = 1$ . This indicates the contribution of the  $j$ th data point in the  $K$  nearest neighborhood of  $i$ th data point to the construction of  $i$ th data point.  $\mathcal{N}_i$  is the  $K$  nearest neighbor (kNN) of  $x_i$ , that is,  $W_{ij}$  is zero, if  $j \notin \mathcal{N}_i$ . During the experiments, it can be found that most, if not all scene categories of kNN of a given image may differ. This is due to the

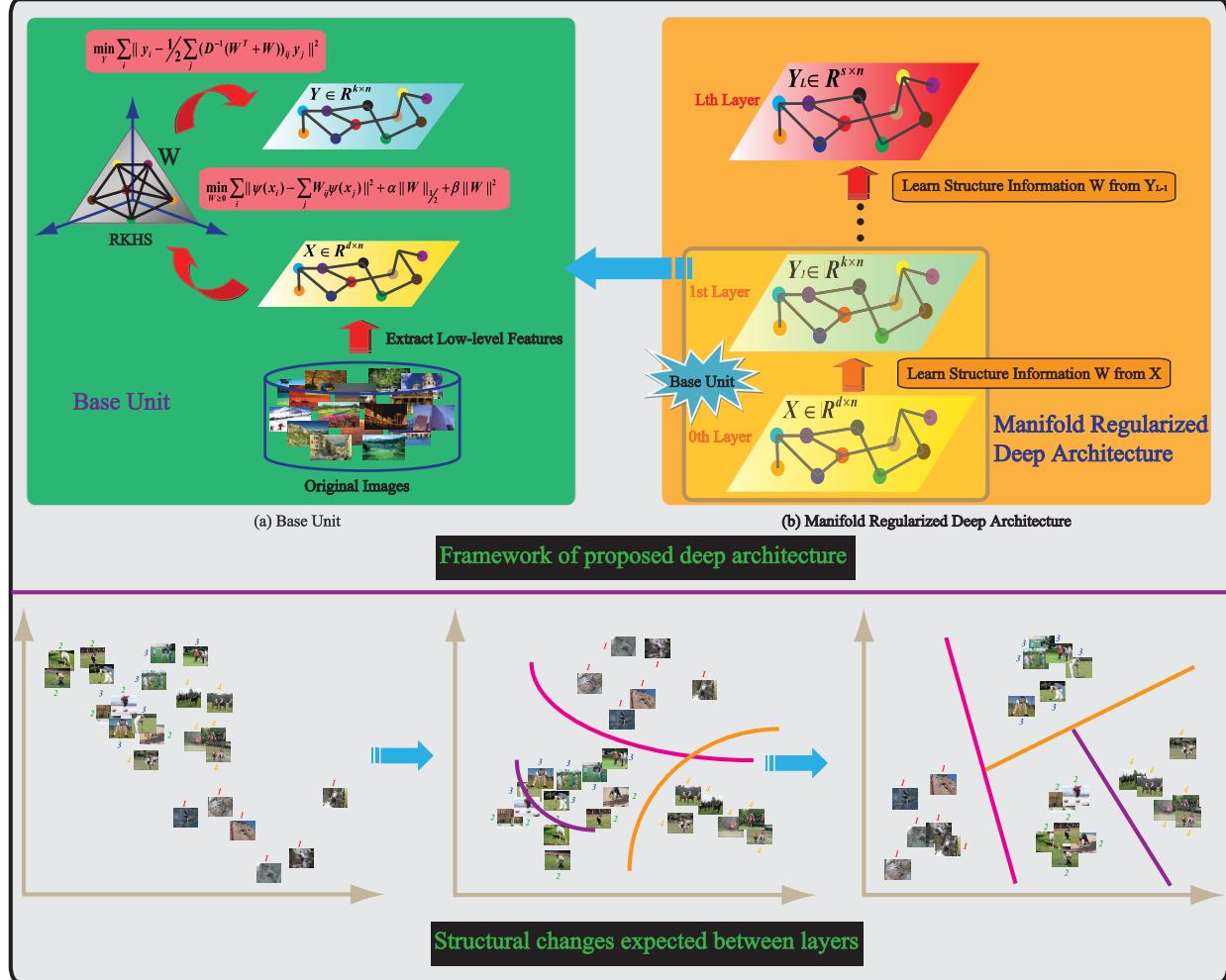


Fig. 2. Framework overview. (a) Base unit. (b) Manifold regularized deep architecture.

Euclidean distance not entirely reflecting the relation among different high-level semantic categories. In such cases, the restriction of  $K$  nearest neighbors is not able to provide high correlated description on the level of the semantic category. We hope that a given image can be well represented by other images, having the same scene category as much as possible. That is, a given point can be well, sparsely represented by all the points. In other words, we bypass kNN entirely. In fact, a given data point is expected to be represented by the strong, high correlation data points in data space on the semantic level. In this case, the nonzero sparse coding coefficients in the weight coefficients matrix  $\mathbf{W}$  corresponding to the strong, high correction data points can best recover the given data point from the other data points. Hence, we relax the constraint to let  $W_{ij}$  be nonzero even if  $j \notin \mathcal{N}_i$ , where constructing the weight coefficients matrix  $\mathbf{W}$  to measure the similarity. In other words, a sparse coding matrix  $\mathbf{W}$  is constructed in this paper, which can describe the structure information better. To achieve better performance for scene recognition, the sparse weight coefficients matrix  $\mathbf{W}$  is calculated in a high-dimensional data space. In this paper,  $L_{1/2}$  regularizer, an unbiased estimator [39] that imposes strong sparsity upon the minimization problem, is introduced to enforce the sparsity of

matrix  $\mathbf{W}$ . In addition, constraint  $\mathbf{W} \geq 0$  ( $W_{ij} \geq 0 \quad \forall i, j$ ) is added to (3) to make a connection to graph embedding. The aforementioned process can be written as

$$\min_{\mathbf{W} \geq 0} \sum_i \left\| \mathbf{x}_i - \sum_j W_{ij} \mathbf{x}_j \right\|^2 + \alpha \|\mathbf{W}\|_{1/2} + \beta \|\mathbf{W}\|^2 \quad (4)$$

where  $\alpha$  and  $\beta$  are the regularization parameters, and  $\|\mathbf{W}\|_{1/2}$  is defined as

$$\|\mathbf{W}\|_{1/2} = \sum_i \sum_j W_{ij}^{1/2}. \quad (5)$$

The local structural information learned by (4) is used to define the hidden space. It results in the following optimization problem:

$$\min_{\mathbf{Y}} \sum_i \left\| \mathbf{y}_i - \frac{1}{2} \sum_j (\mathbf{D}^{-1}(\mathbf{W}^T + \mathbf{W}))_{ij} \mathbf{y}_j \right\|^2 \quad (6)$$

where  $\mathbf{D}$  is a diagonal matrix containing node degree of the learned graph from (4) and satisfies  $\mathbf{Y}\mathbf{D}\mathbf{Y}^T = \mathbf{I}$ . More details about the solver are provided in Section III-C.

One base unit of each layer in the proposed deep architecture is obtainable with (4) and (6). Each basic unit consists

of two steps. First, weight matrix  $\mathbf{W}$  can be learned over data points. Then, data from the previous layer gets mapped to the next layer by exploiting the learned weight matrix  $\mathbf{W}$ . Each hidden space layer can capture the correlations between the activities of units in the previous layer. Compared with the single space layer model, layer-by-layer deep learning architecture learns the structure of the data, and further improves the embedding. Layer-by-layer deep learning architecture is an effective way to progressively reveal low-dimensional nonlinear structure.

### B. Weight Matrix $\mathbf{W}$ Modeling

In this section, the kernel trick aids the description of  $\mathbf{W}$ . The kernel trick discovers the nonlinear structure in the data by mapping the original nonlinear observations to a linear space of a higher-dimension. There are many ways to construct a kernel, the most common approaches use classic kernels, such as Gaussian kernel, linear kernel, polynomial kernel, and so on. However, the nonlinear structure learned by the classic kernels may be inconsistent, with the inherent manifold structure of data, because of classic kernels are independent of the data. In this paper, the manifold structure is introduced into the kernel space to generate the manifold kernel space that can then capture the underlying structure of the data. The weight matrix  $\mathbf{W}$  can be learned to improve the performance in the manifold kernel space.

The structure information between data is described by a sparse weight matrix  $\mathbf{W}$  in (4). In fact, the weight matrix  $\mathbf{W}$  represents the distribution of data on the manifold implicitly. Therefore, a manifold kernel is used to describe the structure information between data. Let  $\mathcal{X}$  be a compact domain in Euclidean space or a manifold. Denote  $\mathcal{H}$  be a complete Hilbert space of function  $\mathcal{X} \rightarrow \mathbb{R}$  and  $\mathcal{Q}$  be a linear space with a positive semidefinite inner product (quadratic form). Define a bounded linear operator  $S : \mathcal{H} \rightarrow \mathcal{Q}$ . Then a space  $\tilde{\mathcal{H}}$  of functions from Hilbert space  $H$  can be defined with the modified inner product

$$\langle f, g \rangle_{\tilde{\mathcal{H}}} = \langle f, g \rangle_{\mathcal{H}} + \langle S(f), S(g) \rangle_{\mathcal{Q}}.$$

Sindhwani *et al.* [40] have shown that  $\tilde{\mathcal{H}}$  is still a *Reproducing Kernel Hilbert Space* (RKHS).

Given the set of data points  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , let  $S : \mathcal{H} \rightarrow \mathbb{R}^n$  be the sampling operator

$$S(f) = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^T.$$

Denote  $\mathbf{f} = S(f) = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^T$ ,  $\mathbf{g} = (g(\mathbf{x}_1), g(\mathbf{x}_2), \dots, g(\mathbf{x}_n))^T$ . The inner product on  $\mathbb{R}^n$  can be obtained by matrix  $\mathbf{M}$ , which is a symmetric positive semidefinite matrix

$$\langle \mathbf{f}, \mathbf{g} \rangle = \langle S(f), S(g) \rangle_{\mathcal{Q}} = \mathbf{f}^T \mathbf{M} \mathbf{g}.$$

The norm on  $\mathbb{R}^n$  can be defined as follows:

$$\|S(f)\|_{\mathcal{Q}}^2 = \mathbf{f}^T \mathbf{M} \mathbf{f}.$$

Define

$$\mathbf{k}_{\mathbf{x}_i} = (\mathcal{K}(\mathbf{x}_i, \mathbf{x}_1), \mathcal{K}(\mathbf{x}_i, \mathbf{x}_2), \dots, \mathcal{K}(\mathbf{x}_i, \mathbf{x}_n))$$

where  $\mathcal{K}$  is a kernel which naturally defines a unique RKHS  $\tilde{\mathcal{H}}$ . The reproducing kernel induced the RKHS  $\tilde{\mathcal{H}}$  is shown as follows:

$$\tilde{\mathcal{K}}_{ij} = \tilde{\mathcal{K}}(\mathbf{x}_i, \mathbf{x}_j) = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) - \theta \mathbf{k}_{\mathbf{x}_i}^T (\mathbf{I} + \mathbf{MK})^{-1} \mathbf{M} \mathbf{k}_{\mathbf{x}_j} \quad (7)$$

where  $\mathbf{I}$  is an unit matrix,  $\mathbf{K}$  is the kernel matrix, and  $\tilde{\mathcal{K}}_{ij}$  is the element of matrix  $\tilde{\mathcal{K}}$ .  $\theta$  is a parameter which can adjust the smoothness of the functions and it satisfies  $\theta \geq 0$ . The choice of matrix  $\mathbf{M}$  is a key issue for the construction of the reproducing kernel. In this paper, the deformation of the kernel can be induced by the data-dependent norm, which is motivated according to the inherent manifold of the data.

To describe the structure information between data better by the kernel trick, the aforementioned manifold kernel is used. As mentioned previously, the choice of  $\mathbf{M}$  is critical, we will introduce how to construct the matrix  $\mathbf{M}$  in the following content. According to LLE criterion, the reconstruction weight  $\lambda_{ji}$  of the neighbors in  $\mathbf{x}_i$  can be computed by minimizing the local reconstruction error

$$\hat{\lambda}_{ji} = \operatorname{argmin}_{\lambda_{ji}} \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} \lambda_{ji} \mathbf{x}_j \right\|_2^2 \quad (8)$$

where  $\mathcal{N}_i$  denotes the  $K$  nearest neighbor of  $\mathbf{x}_i$  (here Euclidean distance is used to define neighborhood);  $\lambda_{ji}$  is the reconstruction weight, which is subjected to the constraints  $\sum_{j \in \mathcal{N}_i} \lambda_{ji} = 1$ , and  $\lambda_{ji} = 0$  for all  $\mathbf{x}_j \notin \mathcal{N}_i$ . As is described in [41], the minimization of (8) can be done with the help of a Gram matrix

$$G_i = (\mathbf{x}_i \mathbf{1}^T - \mathbf{X})^T (\mathbf{x}_i \mathbf{1}^T - \mathbf{X}) \quad (9)$$

where each column of  $\mathbf{X}$  is one neighbor of  $\mathbf{x}_i$  and  $\mathbf{1}$  is a column vector whose elements are all ones. Let  $\lambda_i$  be the concatenation of  $\lambda_{ji}$ . The weight estimator  $\lambda_i$  is then obtained as follows:

$$G_i \lambda_i = 1. \quad (10)$$

Then  $\lambda_i$  is normalized so that  $\sum_{j \in \mathcal{N}_i} \lambda_{ji} = 1$ . Thus the reconstruction weight  $\hat{\lambda}$  is solved.

With the obtained reconstruction weight  $\hat{\lambda}_{ji}$ , the local structure information can be described as follows:

$$\begin{aligned} \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} \hat{\lambda}_{ji} \mathbf{x}_j \right\|^2 &= \sum_i \left\| \mathbf{x}_i - \sum_j \hat{\lambda}_{ji} \mathbf{x}_j \right\|^2 \\ &= \|\mathbf{X} - \mathbf{X} \hat{\Lambda}\|^2 \\ &= \operatorname{Tr}(\mathbf{X}(\mathbf{I} - \hat{\Lambda})(\mathbf{I} - \hat{\Lambda})^T \mathbf{X}^T) \\ &= \operatorname{Tr}(\mathbf{X} \Gamma \mathbf{X}^T) \end{aligned}$$

where  $\operatorname{Tr}$  is the trace of a matrix,  $\mathbf{I}$  is the unit matrix, and  $\Gamma = (\mathbf{I} - \hat{\Lambda})(\mathbf{I} - \hat{\Lambda})^T$ .

Setting  $\mathbf{M} = \Gamma$ , we get the following manifold kernel:

$$\tilde{\mathcal{K}}_{ij} = \tilde{\mathcal{K}}(\mathbf{x}_i, \mathbf{x}_j) = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) - \theta \mathbf{k}_{\mathbf{x}_i}^T (\mathbf{I} + \Gamma \mathbf{K})^{-1} \Gamma \mathbf{k}_{\mathbf{x}_j}. \quad (11)$$

Equation (4) is solved with the obtained manifold kernel  $\tilde{\mathcal{K}}$ . In (4), the mapping  $\psi : \mathbf{X} \rightarrow \psi(\mathbf{X})$  is exploited to map the original data  $\mathbf{X}$  into a high-dimension space in kernel mechanism. It is important to note that only the inner product

$\mathcal{K}_{ij} = \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle$  is needed. Previous literatures [42], [43] have shown that linear models are the best way to construct manifold adaptive kernels, and [44] gives a stronger conclusion. That is, only a linear mapping can unfold any manifold structure in the data. Therefore, to construct the manifold adaptive kernel, linear kernel is regarded as the most reasonable choice. Equation (4) is rewritten as follows:

$$\begin{aligned} & \min_{\mathbf{W} \geq 0} \sum_i \left\| \psi(\mathbf{x}_i) - \sum_j W_{i,j} \psi(\mathbf{x}_i) \right\|^2 + \alpha \|\mathbf{W}\|_{\frac{1}{2}} + \beta \|\mathbf{W}\|^2 \\ &= \min_{\mathbf{W} \geq 0} \|\Psi(\mathbf{X}) - \Psi(\mathbf{X})\mathbf{W}\|^2 + \alpha \|\mathbf{W}\|_{\frac{1}{2}} + \beta \|\mathbf{W}\|^2 \\ &= \min_{\mathbf{W} \geq 0} \text{Tr}(\tilde{\mathcal{K}} - 2\tilde{\mathcal{K}}\mathbf{W} + \mathbf{W}^T \tilde{\mathcal{K}}\mathbf{W}) + \alpha \|\mathbf{W}\|_{\frac{1}{2}} \\ &\quad + \beta \text{Tr}(\mathbf{W}^T \mathbf{W}). \end{aligned}$$

Therefore, the optimization formula is approximated to

$$\min_{\mathbf{W} \geq 0} \text{Tr}[(\tilde{\mathcal{K}} - 2\tilde{\mathcal{K}}\mathbf{W} + \mathbf{W}^T \tilde{\mathcal{K}}\mathbf{W}) + \beta \mathbf{W}^T \mathbf{W}] + \alpha \|\mathbf{W}\|_{\frac{1}{2}}. \quad (12)$$

Furthermore, (12) is equivalent to

$$\min_{\mathbf{W} \geq 0} \text{Tr}[\mathbf{W}^T (\beta \mathbf{I} + \tilde{\mathcal{K}})\mathbf{W} - 2\tilde{\mathcal{K}}\mathbf{W} + \tilde{\mathcal{K}}] + \alpha \|\mathbf{W}\|_{\frac{1}{2}}. \quad (13)$$

It is worth noting that minimizing (12) is subjected to  $\mathbf{W} \geq 0$ . Let  $\zeta \geq 0$  be the corresponding Lagrange multipliers. Consider the Lagrange function  $F(\mathbf{W})$  as

$$\begin{aligned} F(\mathbf{W}) &= \text{Tr}[\mathbf{W}^T (\beta \mathbf{I} + \tilde{\mathcal{K}})\mathbf{W} - 2\tilde{\mathcal{K}}\mathbf{W} + \tilde{\mathcal{K}}] \\ &\quad + \alpha \|\mathbf{W}\|_{\frac{1}{2}} + \text{Tr}(\zeta \mathbf{W}^T). \end{aligned} \quad (14)$$

Then, taking partial derivative with respect to  $\mathbf{W}$  on both sides leads to

$$\frac{\partial F(\mathbf{W})}{\partial W_{ij}} = \left( -2\tilde{\mathcal{K}} + 2\tilde{\mathcal{K}}\mathbf{W} + 2\beta\mathbf{W} + \frac{1}{2}\alpha\mathbf{W}^{-\frac{1}{2}} + \zeta \right)_{ij} \quad (15)$$

where  $\mathbf{W}^{-\frac{1}{2}}$  is equivalent to the inverse matrix of principal square-rooting matrix  $\mathbf{W}^{\frac{1}{2}}$ .

Then the *Karush–Kuhn–Tucker* (KKT) condition  $\zeta \mathbf{W} = 0$  for  $\mathbf{W}$  is

$$\left( -2\tilde{\mathcal{K}} + 2\tilde{\mathcal{K}}\mathbf{W} + 2\beta\mathbf{W} + \frac{1}{2}\alpha\mathbf{W}^{-\frac{1}{2}} \right)_{ij} W_{ij} = 0 \quad \forall i, j. \quad (16)$$

Equation (16) can be rewritten as

$$\left[ -\tilde{\mathcal{K}}_{ij} + \left( \tilde{\mathcal{K}}\mathbf{W} + \beta\mathbf{W} + \frac{1}{4}\alpha\mathbf{W}^{-\frac{1}{2}} \right)_{ij} \right] W_{ij} = 0. \quad (17)$$

An iteratively updating rule on  $W$  is designed as

$$W_{ij} \leftarrow \frac{\tilde{\mathcal{K}}_{ij}}{(\tilde{\mathcal{K}}\mathbf{W} + \beta\mathbf{W} + \frac{1}{4}\alpha\mathbf{W}^{-\frac{1}{2}})_{ij}} \cdot W_{ij}. \quad (18)$$

### C. Hidden Layer Modeling $\mathbf{Y}$

To obtain the solution of the hidden space  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n) \in \mathbb{R}^{k \times n}$ , we minimize the objective in (6). Kong and Ding [45] proposed an approach to solve the problem by normalized cut spectral clustering. To solve the

---

**Algorithm 1** Manifold Regularized Deep Learning Algorithm for Scene Recognition

---

**Require:** Original low-level feature data  $\mathbf{X}$ ; The number of layers  $L$ ;  
**Ensure:** Hidden space output  $\mathbf{Y}_L$ ;  
1: initial  $l = 0$  and  $\mathbf{Y}_0 \leftarrow \mathbf{X}$ ;  
2: **repeat**  
3:      $\mathbf{X} \leftarrow \mathbf{Y}_l$  and then computing kernel matrix  $\tilde{\mathcal{K}}$  by using (7);  
4:     Compute (12), obtain the  $\mathbf{W}$  of  $l$ th layer with updating rule  $W_{ij} \leftarrow \frac{\tilde{\mathcal{K}}_{ij}}{(\tilde{\mathcal{K}}\mathbf{W} + \beta\mathbf{W} + \frac{1}{4}\alpha\mathbf{W}^{-\frac{1}{2}})_{ij}} \cdot W_{ij}$ ;  
5:     Compute  $\tilde{\mathbf{W}} = \frac{1}{2}(\mathbf{W} + \mathbf{W}^T)\mathbf{D}^{-1}$ ;  
6:     Obtaining  $\mathbf{G}$  by computing (21) with normalized cut;  
7:     Solve (19) by (22), obtain the  $\mathbf{Y}_{l+1}$  of  $(l+1)$ th layer;  
8:      $l \leftarrow l + 1$ ;  
9: **until** ( $l > L$ )

---

optimization problem by normalized cut [46], (6) can be transformed as

$$\min_{\mathbf{Y}} \sum_i d_i \left\| \mathbf{y}_i - \frac{1}{2} \sum_j (\mathbf{D}^{-1}(\mathbf{W} + \mathbf{W}^T))_{ij} \mathbf{y}_j \right\|^2 \quad (19)$$

where  $d_i = D_{ii}$  and  $\mathbf{D} = \text{diag}((1/2)(\mathbf{W} + \mathbf{W}^T))$  is a diagonal matrix containing node degrees. Now (19) can be rewritten as

$$\begin{aligned} & \sum_i d_i \left\| \mathbf{y}_i - \frac{1}{2} \sum_j (\mathbf{D}^{-1}(\mathbf{W}^T + \mathbf{W}))_{ij} \mathbf{y}_j \right\|^2 \\ &= \text{Tr}[\mathbf{Y}\mathbf{D}^{\frac{1}{2}}(\mathbf{I} - \tilde{\mathbf{W}})(\mathbf{I} - \tilde{\mathbf{W}})\mathbf{D}^{\frac{1}{2}}\mathbf{Y}^T] \\ &= \text{Tr}[\mathbf{Y}\mathbf{D}^{\frac{1}{2}}(\mathbf{I} - \tilde{\mathbf{W}})^2\mathbf{D}^{\frac{1}{2}}\mathbf{Y}^T] \end{aligned}$$

where  $\tilde{\mathbf{W}} = 1/2(\mathbf{W} + \mathbf{W}^T)\mathbf{D}^{-1}$ . Thus, (19) becomes

$$\min_{\mathbf{Y}} \text{Tr}[\mathbf{Y}\mathbf{D}^{\frac{1}{2}}(\mathbf{I} - \tilde{\mathbf{W}})^2\mathbf{D}^{\frac{1}{2}}\mathbf{Y}^T]. \quad (20)$$

Denote  $\mathbf{G} = \mathbf{D}^{1/2}\mathbf{Y}^T$ . Equation (20) can be rewritten as

$$\min_{\mathbf{G}} \text{Tr}[\mathbf{G}^T(\mathbf{I} - \tilde{\mathbf{W}})^2\mathbf{G}] \quad (21)$$

where  $\mathbf{G}$  needs to satisfy the constraint condition  $\mathbf{G}^T \mathbf{G} = \mathbf{I}$ . In this case, the optimization of (6) is finally converted into a normalized cut problem. Normalized cut is an effective graph partitioning (clustering) and  $\mathbf{G} = (\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_k)$  are relaxed cluster indicators. The optimal solution for  $\mathbf{G}$  is the smallest  $k$  eigenvectors of  $(\mathbf{I} - \tilde{\mathbf{W}})$ . The optimal solution can be represented as

$$\mathbf{Y}^* = \mathbf{G}^T \mathbf{D}^{-\frac{1}{2}}. \quad (22)$$

### D. Procedure of Manifold Regularized Deep Architecture

The detailed procedure of the proposed deep architecture is described in Algorithm 1. One base unit of each layer is achieved from steps 3 to 7 and the output of the base unit is regarded as input of the next base unit in the next layer. Repeat this process until the  $L$  layers deep architecture is achieved.

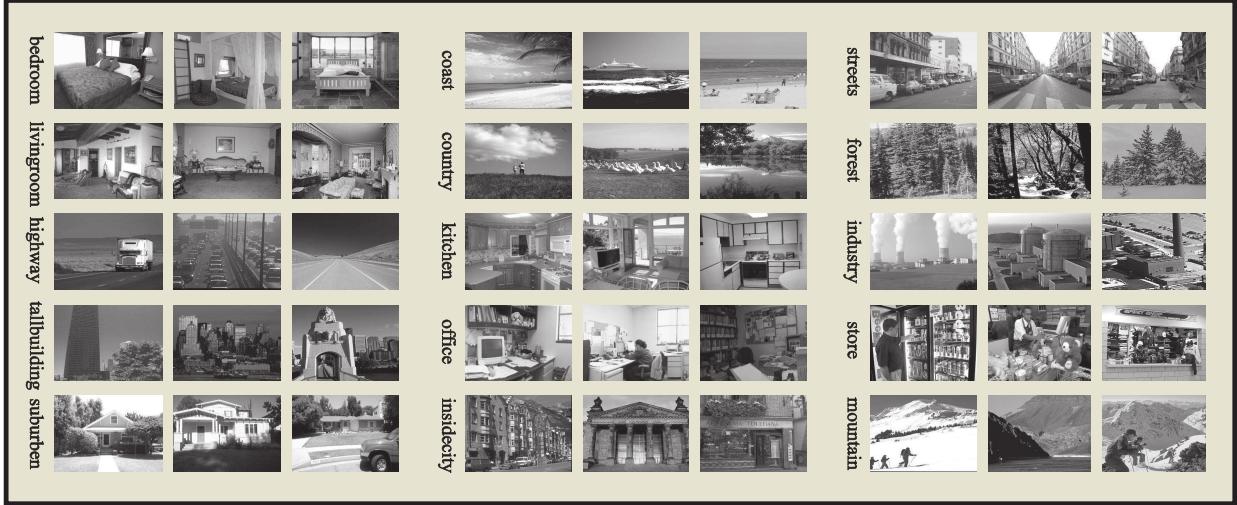


Fig. 3. Sample images of 15 scene categories.

#### IV. EXPERIMENTS

In this section, we evaluate the proposed deep architecture approach on three data sets for scene recognition and compare it with the existing works.

##### A. Data Sets

Three data sets are used to test the proposed deep architecture approach for scene recognition. The first data set is 15 scene categories data set [47], which is the extension of 13 scene categories data set provided by Fei-Fei and Perona [9] (eight of these were originally collected by Oliva and Torralba [48]). This data set contains: coast (360 images), forest (328 images), mountain (274 images), opencountry (410 images), highway (260 images), insidecity (308 images), tallbuilding (365 images), street (292 images), bedroom (216 images), kitchen (210 images), livingroom (289 images), PARoffice (215 images), CALsuburb (241 images), industrial (311 images), and store (315 images). The average resolution of images is  $300 \times 250$ . Fig. 3 shows some sample images.

The second data set is eight sports event categories data set provided by Li and Fei-Fei [49]. The data set contains eight sports event categories: rowing (250 images), badminton (200 images), polo (182 images), bocce (137 images), snowboarding (190 images), croquet (236 images), sailing (190 images), and rock climbing (194 images). These images are high-resolution ones (from  $800 \times 600$  to thousands of pixels per dimension). Sample images are shown in Fig. 4.

Finally, the proposed approach is evaluated on a larger and more challenging data set which is referred to as the SUN data set. So far, the SUN data set is the largest scene recognition data set, and it contains a full variety of 899 scene categories. In this paper, the experiment is conducted on the well-known SUN-397 data set provided by Xiao *et al.* [50], which is a reasonably good subset of the larger SUN data set for scene recognition.

##### B. Implementation

These three popular benchmark data sets are used in our scene recognition experiments. The prevalent training/testing configurations are followed in the literature. For 15 scene categories data set, each category is divided into two separate sets of images, 100 images for training and the rest for testing. For eight sports event categories data set, 70 images per category are used for training and 60 images are used for testing. For each category in SUN data set, 20 images are used to train the classifier and the other 50 images are exploited to evaluate the performance of learned feature. All the compared approaches are conducted on the same partition of the data set. The appearance representation is based on SIFT descriptors which are extracted over  $4 \times 4 \times 8$  bins from  $9 \times 9$  patches. All experiments involving spatial pyramids relied on three pyramid levels aim to consider the spatial structure of images. These low-level features  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  are input to the proposed deep architecture. In this case, the  $\mathbf{Y}_L$  is regarded as high-level features to train classifier. In this paper, for the tests on 15 scene categories data set and eight sports event categories data set, multiclass classification is done with a *support vector machine* (SVM) classifier using the one-versus-all rule. The rule is defined that a classifier is learned to separate each class from the rest, and a test image is assigned the label of the classifier with the highest response. The LIBLINEAR is used to train a linear SVM and switched from one-versus-one to one-versus-all multiclass classification. For the parameter setting of Lagrange multipliers, there are two free parameters:  $\alpha$  and  $\beta$ . It is worth mentioning that the parameters ( $\alpha$  and  $\beta$ ) have few influence on the results from many experiments and we define the parameters  $\alpha = 0.1$  and  $\beta = 0.1$  in our experiments. On the other hand, for the choosing of sparse regularizer of (12), we adopt  $L_{1/2}$  sparse regularizer instead of classical  $L_1$  regularizer. It is worth noting that the recognition accuracy of model using  $L_1$  is less than the recognition accuracy of model using  $L_{1/2}$  regularizer by 0.4% averagely. Moreover, model without manifold kernel has also been experimented

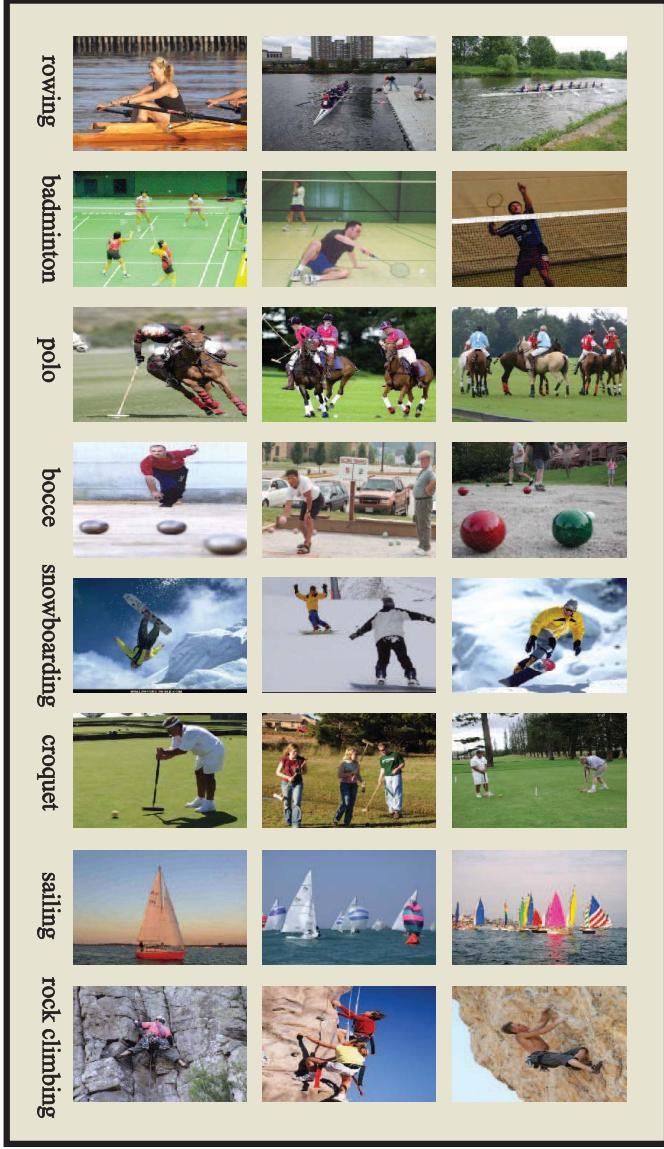


Fig. 4. Sample images of eight event categories.

for testing the performance. The experimental results show the accuracy of model without manifold kernel is poor. Only 54.7% is achieved on 15 scene categories data set and 41.9% is obtained on eight sports event categories data set. The reason is that the model may lose the ability to distinguish the categories of different scene images for lacking kernel.

#### C. Demonstration of Learned Features

The aim of our approach is to use a deep architecture to extract more powerful features for scene recognition task by considering the data structure. We argue that the structure information among data is clearer by the proposed deep architecture, and it is of great benefit for scene recognition task. This is because the structure of the data is helpful with training a classifier. To illustrate this point, we evaluate the experimental results from two-, three-, and four-layer deep architectures. We demonstrate the effectiveness of the proposed deep architecture and learned features by using

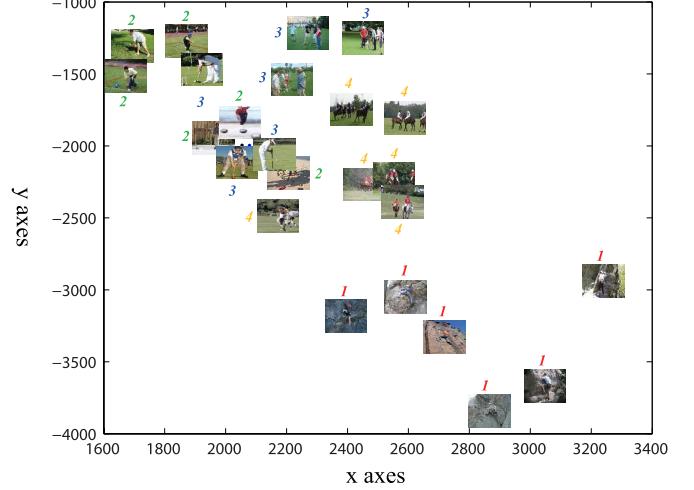


Fig. 5. 2-D visualizations of embedding result from two-layer deep architecture.

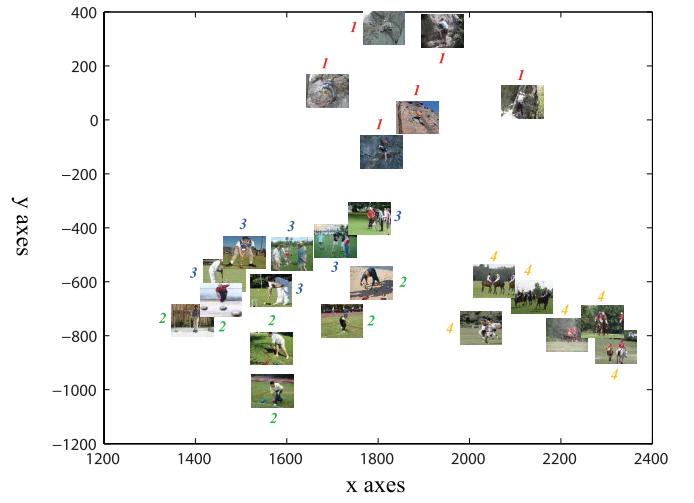


Fig. 6. 2-D visualizations of embedding result from three-layer deep architecture.

2-D visualization. Four scene categories (*rock climbing*, *bocce*, *croquet*, and *polo*) have been selected randomly from eight sport event data sets and each category includes six images which is also sampled randomly. Figs. 5–7 show the 2-D visualizations of embedding results utilizing learned  $\mathbf{W}$  from two-, three-, and four-layer deep architectures, respectively. The numbers 1–4 in figures correspond, respectively, to the sport categories *rock climbing*, *bocce*, *croquet*, and *polo* on eight sports event categories data set. In the result of two-layer deep architecture, all images from different scene categories cluster together. For the results obtained from three-layer and four-layer deep architecture, the classes, however, become more separable. This indicates that more clear data structure and much better recognition performance can be obtained by the proposed deep architecture.

#### D. Comparing to the State-of-the-Art

In this section, we compare the proposed deep architecture with the state-of-the-art in the literature on three standard

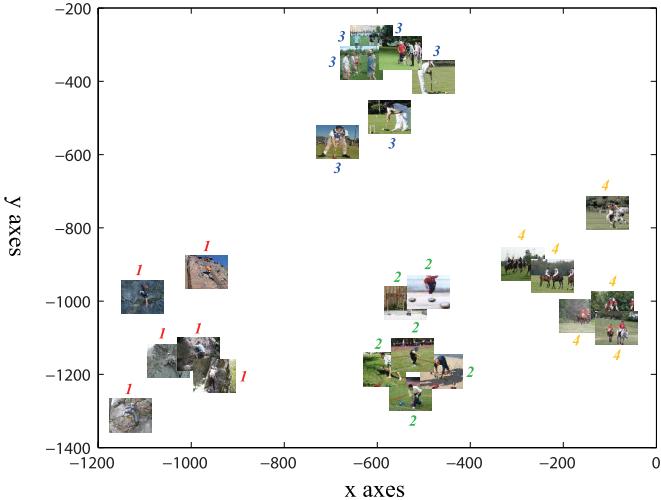


Fig. 7. 2-D visualizations of embedding result from four-layer deep architecture.

data sets. First, the performance of some approaches can be observed on the 15 scene categories data set. The outstanding spatial pyramid matching was provided by Lazebnik *et al.* [47]. In [47], the images were partitioned into the fine subregions and the histograms of local features inside each subregion are computed. In this case, the recognition accuracy of 81.2% can be obtained in [47]. Dixit *et al.* [51] presented a general formulation of Bayesian adaptation, which targets class adaptation and is applicable to both the generative and discriminative strategies for the task of image classification. The proposed approach in [51] made a correct recognition rate of 82.3% on the 15 scene categories data set. Kwitt *et al.* [52] proposed a new architecture, denoted *spatial pyramid matching on the semantic manifold* (SPMSM) for scene recognition. SPMSM established a connection between the semantic simplex and a Riemmanian manifold, and equipped the architecture with a similarity measure which respects the manifold structure of the semantic space. Kwitt *et al.* [52] achieved a recognition accuracy of 85.4% on the 15 scene categories data set by SPMSM. Very recently, Goh *et al.* [53] proposed a deep hierarchical architecture according to the restricted Boltzmann machines to encode SIFT descriptors and provided the vectorial representation for image categorization. The model merges the complementary strengths of the BoF framework and deep architectures. For the 15 scene categories data set, the deep architecture proposed by Goh *et al.* [53] obtained an average recognition accuracy of 85.4%. It can be seen from Table I that the proposed approach and deep learning models of [53] outperform shallow models including Lazebnik *et al.* [47], Dixit *et al.* [51], and Kwitt *et al.* [52]. To verify the performance of the proposed deep architecture, two-layer deep model, three-layer deep model, and four-layers deep model output are adopted, respectively, as the learned features to run the SVM classifier for the task of scene recognition. As shown in Table I, the recognition accuracy of the proposed approach is high compared with the other approaches. Moreover, it is noteworthy that the deep architecture with more layers significantly

TABLE I  
COMPARISON TO THE STATE-OF-THE-ART ON 15 SCENE CATEGORIES DATA SET

Dataset	State-of-the-Art	Recognition Accuracy
15 scene categories dataset	Lazebnik <i>et al.</i> [47]	81.2%
	Dixit <i>et al.</i> [51]	82.3%
	Kwitt <i>et al.</i> [52]	85.4%
	Goh <i>et al.</i> [53]	85.4%
	proposed (2-layers)	83.2%
	proposed (3-layers)	84.7%
	proposed (4-layers)	<b>86.0%</b>
	proposed (5-layers)	<b>86.9%</b>

TABLE II  
COMPARISON TO THE STATE-OF-THE-ART ON EIGHT SPORTS EVENT CATEGORIES DATA SET

Dataset	State-of-the-Art	Recognition Accuracy
8 sports event categories dataset	Li and Fei-Fei [49]	73.4%
	Kwitt <i>et al.</i> [52]	83.0%
	Wu and Rehg [12]	84.3%
	proposed (2-layers)	82.3%
	proposed (3-layers)	<b>84.4%</b>
	proposed (4-layers)	<b>85.6%</b>
	proposed (5-layers)	<b>86.1%</b>

TABLE III  
COMPARISON WITH THE STATE-OF-THE-ART ON SUN DATA SET

Dataset	State-of-the-Art	Recognition Accuracy
SUN dataset	Xiao <i>et al.</i> [50]	27.2%
	Kwitt <i>et al.</i> [52]	28.9%
	Donahue <i>et al.</i> [36]	30.14%
	proposed (2-layers)	<b>29.2%</b>
	proposed (3-layers)	<b>30.1%</b>
	proposed (4-layers)	<b>30.3%</b>
	proposed (5-layers)	<b>30.3%</b>

improves the recognition accuracy compared with the architecture with less layers.

Second, the performance of the proposed deep architecture can be verified by comparing with the state-of-the-art approaches on the eight sports event categories data set. Li and Fei-Fei [49] proposed an integrative model which learns to classify static images into complicated sport games by interpreting the semantic components of the images as detailed as possible. By using the integrative model, they show that their system is capable of recognizing these event categories with 73.4% accuracy. Wu and Rehg [12] showed that the *Histogram Intersection Kernel* (HIK) is either more effective than the Euclidean distance in supervised learning tasks with histogram features or is used in an unsupervised manner to significantly improve the generation of visual codebooks. The HIK-based codebook generation approach consistently achieves higher accuracy than *k*-means codebooks by 2%–4%, and can achieve the 84.3% accuracy on the eight

TABLE IV  
FLOATING-POINT CALCULATION TIMES FOR EACH ITERATION IN THE PROPOSED DEEP ARCHITECTURE

operator category	Update $\mathbf{W}$	Update $\mathbf{Y}$	Total
addition	$O(n^3 + 2n^2)$	$O(kn)$	$O(n^3 + 2n^2 + kn)$
multiplication	$O(n^3 + 2n^2)$	$O(kn)$	$O(n^3 + 2n^2 + kn)$
division	$O(n^2)$	0	$O(n^2)$

sports event categories data set. Kwitt *et al.* [52] also tested their proposed approach on these sports event data set and achieved a recognition accuracy of 83.0%. It can be seen in Table II that the proposed deep architecture outperforms other state-of-the-art approaches. In general, the comparable performance between the proposed approach and other approaches is due to the proposed approach's capability to learn deeply the effective feature layer by layer.

Finally, to evaluate the performance of the feature learned by the proposed deep architecture on the large data set, the proposed approach is compared with the other approaches on the SUN data set. Convolution neural networks have achieved amazing successes on visual classification recently, so we have compared the proposed deep architecture with convolution neural networks on SUN data set which is the largest and most challenging data set for scene recognition. Donahue *et al.* [36] proposed an approach to apply the features extracted from the activation of a deep convolutional network trained in a fully supervised fashion on a large, fixed data set to novel generic tasks. In this paper, we have compared the proposed unsupervised deep architecture with deep convolution networks of [36] over the SUN data set. For SUN data set, Donahue *et al.* [36] achieved recognition accuracy of 30.14%. It can be seen from Table III that the best performance of the proposed approach beats the state-of-the-art on the SUN data set. Moreover, it can be found that the deep learning models perform much more significantly better than shallow models.

Furthermore, we observe a phenomenon generated from the experimental results: with the increase of layers of deep architecture, the growth rate of recognition accuracy will monotonically decrease. For example, in Table I, there is an increase of 1.5% recognition accuracy from two to three layers and an increase of 1.3% from three to four layers. Furthermore, the observed regular pattern is similar with other approaches, such as DBNs [26] and convolutional DBNs [27]. In these deep learning approaches, four layer is often considered owing to the computational cost [26], [27]. Therefore, we can argue that an increase of less than 1.3% will appear when five layers are added. To verify this inference, we conduct the experiments with five-layer model, the recognition accuracies on different data set are shown in Tables I–III. On the other hand, the size and complexity of data set will affect the increase range of recognition accuracy (e.g., an increase of 2.8% from two to four layers appears on the 15 scene categories data set and an increase of only 1.1% is found on the SUN data set).

Here, the computational complexity of the proposed deep architecture will be analyzed. Concentrating on the rules represented by (18) and (22), the cost of computing  $\mathbf{W}^{-1/2}$  is known as  $O(n^3)$ , where  $n$  is the number of data. Except for

this cost, the other three floating-point calculation times for each iteration are listed in Table IV.

Moreover, the training time of the proposed approach is analyzed. The proposed approach is implemented in a MATLAB environment on a modern desktop computer (3.4 GHz Core 2 Quad with 32 GB random access memory) to test the performance using eight sport event data set, 15 scene categories data set, and Sun data set. It takes less than 6 and 13 min to obtain the experimental results over the eight sport events data set and 15 scene categories data set, respectively. Because the SUN data set is the largest and challenging data set for scene recognition task, the performance on this data set is time consuming and it takes about 2 h. In this paper, the proposed approach aims to search the effectiveness of learned high-level feature, and the training process is offline. Therefore, the computational cost of the proposed deep architecture is acceptable.

## V. CONCLUSION

In this paper, a novel manifold regularized deep architecture is proposed for scene recognition. The proposed deep architecture exploits the structure information of data to learn deep, effective features layer by layer. The proposed approach outweighs state-of-the-art performance on three publicly available data sets. We plan to extend the proposed deep architecture to extract the high-level feature in other fields such as image/video retrieval.

## ACKNOWLEDGMENT

The authors would like to thank D. Liu, Editor-in-Chief, and the handling Associate Editor of IEEE TNNLS and four anonymous reviewers, whose insightful comments and constructive suggestions led to significant improvement of the presentation of this paper.

## REFERENCES

- [1] D. Wang, "The time dimension for scene analysis," *IEEE Trans. Neural Netw.*, vol. 16, no. 6, pp. 1401–1426, Nov. 2005.
- [2] A. Vailaya, M. A. T. Figueiredo, A. K. Jain, and H.-J. Zhang, "Image classification for content-based indexing," *IEEE Trans. Image Process.*, vol. 10, no. 1, pp. 117–130, Jan. 2001.
- [3] E. Chang, K. Goh, G. Sychay, and G. Wu, "CBSA: Content-based soft annotation for multimodal image retrieval using Bayes point machines," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 13, no. 1, pp. 26–38, Jan. 2003.
- [4] R. Eckhorn, "Neural mechanisms of scene segmentation: Recordings from the visual cortex suggest basic circuits for linking field models," *IEEE Trans. Neural Netw.*, vol. 10, no. 3, pp. 464–479, May 1999.
- [5] C. Liu, J. Yuen, and A. Torralba, "Nonparametric scene parsing: Label transfer via dense scene alignment," in *Proc. IEEE CVPR*, Jun. 2009, pp. 1972–1979.

- [6] R. Socher, C. C. Lin, A. Y. Ng, and C. Manning, "Parsing natural scenes and natural language with recursive neural networks," in *Proc. 28th ICML*, 2011, pp. 129–136.
- [7] X. Yu, C. Fermüller, C. L. Teo, Y. Yang, and Y. Aloimonos, "Active scene recognition with vision and language," in *Proc. ICCV*, 2011, pp. 810–817.
- [8] J. Yao, S. Fidler, and R. Urtasun, "Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation," in *Proc. IEEE CVPR*, Jun. 2012, pp. 702–709.
- [9] L. Fei-Fei and P. Perona, "A Bayesian hierarchical model for learning natural scene categories," in *Proc. IEEE CVPR*, vol. 2. Jun. 2005, pp. 524–531.
- [10] P. Quelhas, F. Monay, J. Odobez, D. Gatica-Perez, and T. Tuytelaars, "A thousand words in a scene," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 9, pp. 1575–1589, Sep. 2007.
- [11] A. Bosch, A. Zisserman, and X. Muñoz, "Scene classification using a hybrid generative/discriminative approach," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 4, pp. 712–727, Apr. 2008.
- [12] J. Wu and J. M. Rehg, "Beyond the Euclidean distance: Creating effective visual codebooks using the histogram intersection kernel," in *Proc. IEEE 12th ICCV*, Sep./Oct. 2009, pp. 630–637.
- [13] O. Boiman, E. Shechtman, and M. Irani, "In defense of nearest-neighbor based image classification," in *Proc. IEEE CVPR*, Jun. 2008, pp. 1–8.
- [14] L. Li, H. Su, E. Xing, and L. Fei-Fei, "Object bank: A high-level image representation for scene classification: semantic feature sparsification," in *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates, Inc., 2010.
- [15] J. Vogel and B. Schiele, "Semantic modeling of natural scenes for content-based image retrieval," *Int. J. Comput. Vis.*, vol. 72, no. 2, pp. 133–157, 2007.
- [16] L. Torresani, M. Szummer, and A. Fitzgibbon, "Efficient object category recognition using classes," in *Proc. ECCV*, 2010, pp. 776–789.
- [17] D. Parikh and K. Grauman, "Interactively building a discriminative vocabulary of nameable attributes," in *Proc. IEEE CVPR*, Jun. 2011, pp. 1681–1688.
- [18] L. Liu, L. Shao, and X. Li, "Building holistic descriptors for scene recognition: A multi-objective genetic programming approach," in *Proc. 21st ACM Multimedia*, 2013, pp. 997–1006.
- [19] J. Yu, D. Tao, Y. Rui, and J. Cheng, "Pairwise constraints based multiview features fusion for scene classification," *Pattern Recognit.*, vol. 46, no. 2, pp. 483–496, 2013.
- [20] J. Cheng *et al.*, "Peripapillary atrophy detection by sparse biologically inspired feature manifold," *IEEE Trans. Med. Imag.*, vol. 31, no. 12, pp. 2355–2365, Dec. 2012.
- [21] T. Zhou and D. Tao, "Double shrinking sparse dimension reduction," *IEEE Trans. Image Process.*, vol. 22, no. 1, pp. 244–257, Jan. 2013.
- [22] D. Tao, X. Li, X. Wu, and S. J. Maybank, "Geometric mean for subspace selection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 2, pp. 260–274, Feb. 2009.
- [23] D. Tao, X. Li, X. Wu, and S. J. Maybank, "General tensor discriminant analysis and Gabor features for gait recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 10, pp. 1700–1715, Oct. 2007.
- [24] T. Zhou, D. Tao, and X. Wu, "Manifold elastic net: A unified framework for sparse dimension reduction," *Data Mining Knowl. Discovery*, vol. 22, no. 3, pp. 340–371, 2011.
- [25] W. K. Wong and M. Sun, "Deep learning regularized Fisher mappings," *IEEE Trans. Neural Netw.*, vol. 22, no. 10, pp. 1668–1675, Oct. 2011.
- [26] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Comput.*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [27] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *Proc. 26th ICML*, 2009, pp. 609–616.
- [28] K. Sohn, D. Y. Jung, H. Lee, and A. O. Hero, "Efficient learning of sparse, distributed, convolutional feature representations for object recognition," in *Proc. ICCV*, 2011, pp. 2643–2650.
- [29] Q. V. Le, W. Y. Zou, S. Y. Yeung, and A. Y. Ng, "Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis," in *Proc. IEEE CVPR*, Jun. 2011, pp. 3361–3368.
- [30] G. B. Huang, H. Lee, and E. Learned-Miller, "Learning hierarchical representations for face verification with convolutional deep belief networks," in *Proc. IEEE CVPR*, Jun. 2012, pp. 2518–2525.
- [31] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [32] R. Mittelman, H. Lee, B. Kuipers, and S. Savarese, "Weakly supervised learning of mid-level features with Beta-Bernoulli process restricted Boltzmann machines," in *Proc. IEEE CVPR*, Jun. 2013, pp. 476–483.
- [33] Y. Kim, H. Lee, and E. M. Provost, "Deep learning for robust feature generation in audiovisual emotion recognition," in *Proc. IEEE ICASSP*, May 2013, pp. 3687–3691.
- [34] G. Zhou, K. Sohn, and H. Lee, "Online incremental feature learning with denoising autoencoders," in *Proc. AISTATS*, 2012, pp. 1453–1461.
- [35] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional learning of spatio-temporal features," in *Proc. ECCV*, 2010, pp. 140–153.
- [36] J. Donahue *et al.*, "Decaf: A deep convolutional activation feature for generic visual recognition," in *Proc. ICML*, 2014, pp. 647–655.
- [37] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1915–1929, Aug. 2013.
- [38] P. H. O. Pinheiro and R. Collobert, "Recurrent convolutional neural networks for scene labeling," in *Proc. ICML*, 2014, pp. 82–90.
- [39] Z. Xu, H. Zhang, Y. Wang, and Y. Chang, " $L_{1/2}$  regularizer," *Sci. China Inf. Sci.*, vol. 53, no. 6, pp. 1159–1169, 2010.
- [40] V. Sindhwani, P. Niyogi, and M. Belkin, "Beyond the point cloud: From transductive to semi-supervised learning," in *Proc. 22nd ICML*, 2005, pp. 824–831.
- [41] H. Chang, D.-Y. Yeung, and Y. Xiong, "Super-resolution through neighbor embedding," in *Proc. IEEE CVPR*, Jun./Jul. 2004, pp. 275–282.
- [42] T. Joachims, "Text categorization with support vector machines: Learning with many relevant features," in *Proc. ECML*, 1998, pp. 137–142.
- [43] Y. Yang, "An evaluation of statistical approaches to text categorization," *Inf. Retr.*, vol. 1, nos. 1–2, pp. 69–90, 1999.
- [44] D. Cai, *Spectral Regression: A Regression Framework for Efficient Regularized Subspace Learning*. Ann Arbor, MI, USA: ProQuest, 2009.
- [45] D. Kong and C. Ding, "An iterative locally linear embedding algorithm," in *Proc. ICML*, 2012.
- [46] J. Shi and J. Malik, "Normalized cuts and image segmentation," in *Proc. IEEE CVPR*, Jun. 1997, pp. 731–737.
- [47] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE CVPR*, Jun. 2006, pp. 2169–2178.
- [48] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3, pp. 145–175, 2001.
- [49] L.-J. Li and L. Fei-Fei, "What, where and who? Classifying events by scene and object recognition," in *Proc. IEEE ICCV*, Oct. 2007, pp. 1–8.
- [50] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "SUN database: Large-scale scene recognition from abbey to zoo," in *Proc. IEEE CVPR*, Jun. 2010, pp. 3485–3492.
- [51] M. Dixit, N. Rasiwasia, and N. Vasconcelos, "Adapted Gaussian models for image classification," in *Proc. IEEE CVPR*, Jun. 2011, pp. 937–943.
- [52] R. Kwitt, N. Vasconcelos, and N. Rasiwasia, "Scene recognition on the semantic manifold," in *Proc. 12th ECCV*, 2012, pp. 359–372.
- [53] H. Goh, N. Thome, M. Cord, and J.-H. Lim, "Learning deep hierarchical visual feature coding," *IEEE Trans. Neural Netw. Learn. Syst.*, doi: 10.1109/TNNLS.2014.2307532.

**Yuan Yuan** (M'05–SM'09) is currently a Full Professor with the Chinese Academy of Sciences, Beijing, China. She has authored over 100 papers, including over 70 in reputable journals, such as the IEEE TRANSACTIONS and *Pattern Recognition*, and conference papers in the IEEE Conference on Computer Vision and Pattern Recognition, the British Machine Vision Conference, the IEEE International Conference on Image Processing, and the International Conference on Acoustics, Speech and Signal Processing. Her current research interests include visual information processing and image/video content analysis.



**Lichao Mou** is currently pursuing the M.S. degree with the Center for OPTical IMagery Analysis and Learning (OPTIMAL), State Key Laboratory of Transient Optics and Photonics, Xi'an Institute of Optics and Precision Mechanics, Chinese Academy of Sciences, Xi'an, Shaanxi, P. R. China. His current research interests include pattern recognition, machine learning, and computer vision.



**Xiaoqiang Lu** is currently an Associate Professor with the Center for OPTical IMagery Analysis and Learning (OPTIMAL), State Key Laboratory of Transient Optics and Photonics, Xi'an Institute of Optics and Precision Mechanics, Chinese Academy of Sciences, Xi'an, Shaanxi, P. R. China. His current research interests include pattern recognition, machine learning, hyperspectral image analysis, cellular automata, and medical imaging.