

Aggregating Rich Hierarchical Features for Scene Classification in Remote Sensing Imagery

Guoli Wang, Bin Fan, *Senior Member, IEEE*, Shiming Xiang, and Chunhong Pan, *Member, IEEE*

Abstract—Scene classification is one of the most important issues in remote sensing image processing. To obtain a high discriminative feature representation for an image to be classified, traditional methods usually consider to densely accumulate hand-crafted low-level descriptors (e.g., scale-invariant feature transform) by feature encoding techniques. However, the performance is largely limited by the hand-crafted descriptors as they are not capable of describing the rich semantic information contained in various remote sensing images. To alleviate this problem, we propose a novel method to extract discriminative image features from the rich hierarchical information contained in convolutional neural networks (CNNs). Specifically, the low-level and middle-level intermediate convolutional features are, respectively, encoded by vector of locally aggregated descriptors (VLAD) and then reduced by principal component analysis to obtain hierarchical global features; meanwhile, the fully connected features are average pooled and subsequently normalized to form new global features. The proposed encoded mixed-resolution representation (EMR) is the concatenation of all the above-mentioned global features. Due to the usage of encoding strategies (VLAD and average pooling), our method can deal with images of different sizes. In addition, to reduce the computational consumption in the training stage, we directly extract EMR from VGG-VD and ResNet pretrained on the ImageNet dataset. We show in this paper that CNNs pretrained on the natural image dataset are more easily applied to the remote sensing dataset when the local structure similarity between two datasets is higher. Experimental evaluations on the UC-Merced and Brazilian Coffee Scenes datasets demonstrate that our method is superior to the state of the art.

Index Terms—Convolutional neural networks (CNNs), mixed-resolution representation, remote sensing scene classification, vector of locally aggregated descriptors (VLAD).

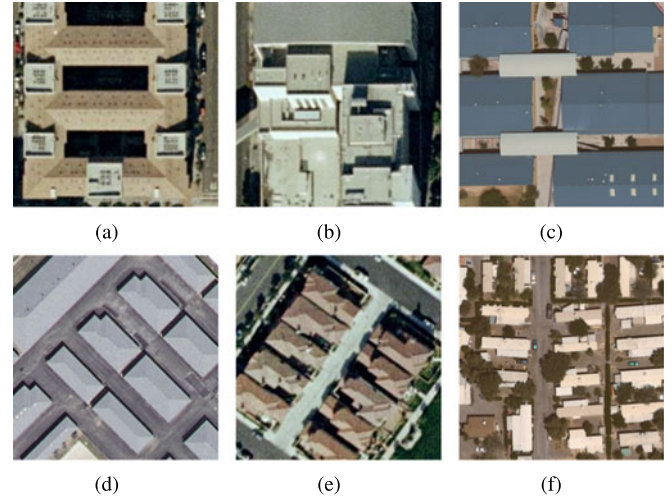


Fig. 1. Examples of different scene image types in the UC-Merced dataset [1]: (a)–(d) buildings, (e) dense residential, (f) mobile home park.

I. INTRODUCTION

SCENE classification of remote sensing imagery refers to the task of classifying an image into different categories of land covers and ground objects. It is an important research topic in the remote sensing community, and plays a significant role in many practical applications, ranging from land management, urban planning to environment prospecting and monitoring.

To achieve a good classification performance, features representing remote sensing images should have small within-class scatter and large between-class scatter. However, due to large variations of land covers and ground objects on the Earth, there are large differences in their sizes, shapes, and structures. As a result, it is very common in remote sensing images that different semantic categories share some similar contents while images from the same category have different appearance and texture. As shown in Fig. 1, although the four images in (a)–(d) belong to a same scene category, they have completely different building structures and shapes. Meanwhile, Fig. 1(d)–(f) shows three images of different categories that contain similar objects. As all of them are constituted by buildings, the different densities and building sizes are the only ways to distinguish them. Owing to these challenges, scene classification in remote sensing images is still an open problem and remains active in the community.

To address this problem, Yang *et al.* [1] considered to use the bag-of-visual-words (BOVW) model combined with a support vector machine (SVM). Specifically, key points are detected in

Manuscript received December 21, 2016; revised March 27, 2017 and May 2, 2017; accepted May 4, 2017. This work was supported in part by the National Natural Science Foundation of China under Grant 61403375, Grant 61472119, Grant 61573352, Grant 61375024, Grant 91338202, and Grant 91646207, in part by the Priority Academic Program Development of Jiangsu Higher Education Institutions, and in part by the Jiangsu Collaborative Innovation Center on Atmospheric Environment and Equipment Technology. (Corresponding author: Bin Fan.)

B. Fan, S. Xiang, and C. Pan are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China (e-mail: bfan@nlpr.ia.ac.cn; smxiang@nlpr.ia.ac.cn; chpan@nlpr.ia.ac.cn).

G. Wang is with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and also with the School of Computer and Control Engineering, University of Chinese Academy of Sciences, Beijing 101408, China (e-mail: glwang@nlpr.ia.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/JSTARS.2017.2705419

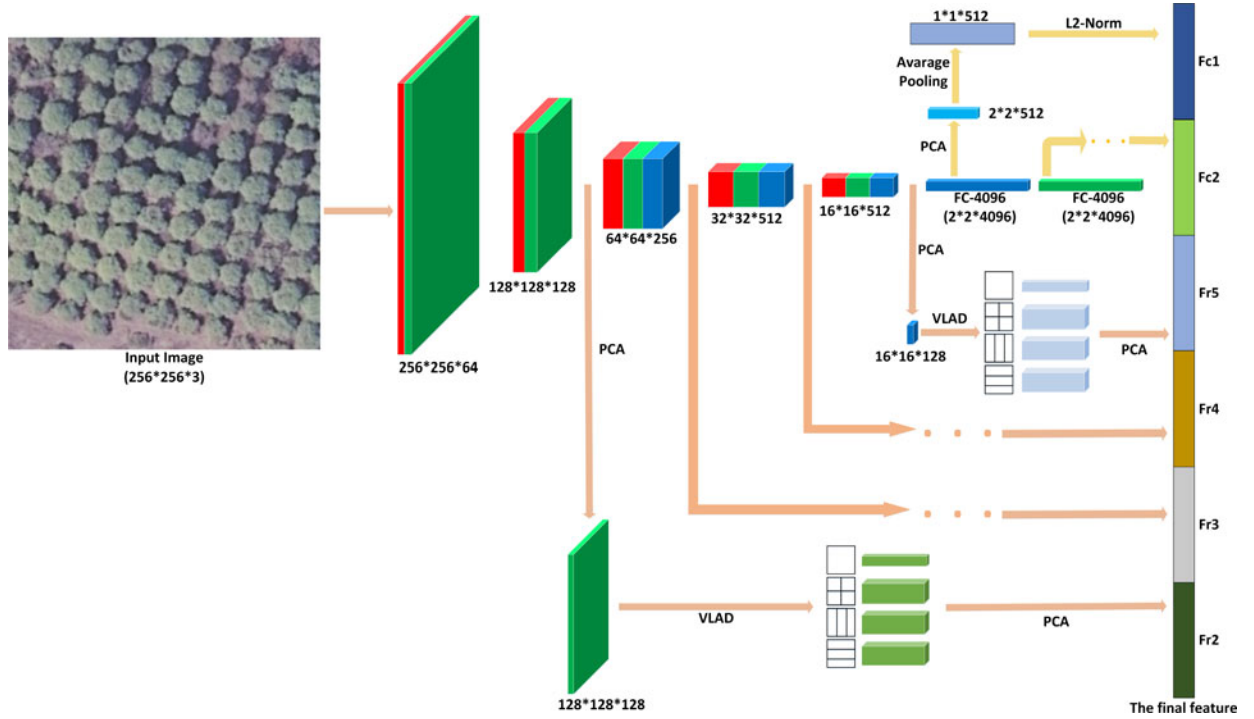


Fig. 2. Proposed framework of EMR by a pretrained VGG-16 network.

a given image by Difference of Gaussian and subsequently represented by scale-invariant feature transform (SIFT) [2]; then, these SIFT descriptors are encoded by hard assignment and spatially pooled together to obtain a single feature vector as the image representation. Since hard assignment only assigns a SIFT descriptor to its closest visual word in the codebook, much useful information is lost. To alleviate this problem, sophisticated encoding methods have been developed, such as sparse coding [3], compressive sensing [4], Fisher kernel coding [5] and so on [6]–[8]. Although these methods can achieve a relative good performance in scene classification, their discriminative abilities are limited by the used local handcraft feature, e.g., SIFT.

Convolutional neural networks (CNNs), a representative data-driven feature learning method, have been gradually considered to be used in remote sensing scene classification due to its large success in natural scene classification [9]–[11] and object recognition [12]. Castelluccio *et al.* [13] studied different learning methods to train CaffeNet [14] and GoogLeNet [15] for land-use classification. Ševo *et al.* [16] used GoogLeNet to automatically detect ground objects in aerial images. Cheng *et al.* [17] learned a rotation-invariant CNN for ground object recognition. Although the above-mentioned works depending on CNNs have achieved amazing performance, two major problems limit the possibility of further improving the performance. The first one is the fixed input image size. For example, the input size of CaffeNet is 227×227 and it has to resize 256×256 images in the UC-Merced dataset to fit the network; some details may be discarded when resizing, hence reducing the classification accuracy. The other problem is that many intermediate convolutional features are discarded in those CNNs. Since CNN is a hierarchical network structure with many intermediate

layers, combining output feature with intermediate features can incorporate richer spatial and semantic information to further improve the discriminative ability.

Besides the above-mentioned two general problems of CNNs, overfitting is a special problem that has to be considered when using CNNs to deal with remote sensing images. Compared with the number of labeled natural scene images, the number of labeled remote sensing scene images is far from sufficient; meanwhile, the manual annotation is too expensive and time consuming. Therefore, fine-tuning a CNN architecture trained on natural images with a small number of labeled remote sensing images becomes a good choice [13]. However, to further adapt the data, thousands of iterations usually are needed to take full advantage of CNNs's potentials [13], [18] and hence fine-tuning CNNs becomes time consuming. Therefore, to reduce the computational consumption in the training stage, we propose to directly extract encoded mixed-resolution representations (EMR) from VGG-VD and ResNet pretrained on the ImageNet dataset. We will show in Section III that CNNs pretrained on the natural image dataset are more easily applied to the remote sensing dataset when the local structure similarity between two datasets is higher.

In this paper, we propose a novel framework to aggregate rich hierarchical features in CNNs for a discriminative image representation. As shown in Fig. 2, a remote sensing image is first fed into a CNN trained on natural images, then both the intermediate layer features in the convolutional layers and the global semantic features in the fully connected layers are encoded and concatenated together to construct the proposed feature. For the intermediate layer features, they are first reduced by principal component analysis (PCA) and then encoded by vector of locally aggregated descriptors (VLAD) [19]; finally,

the VLAD encoded features are subsequently reduced by PCA and concatenated together. For the high-level features in each fully connected layer, they are first reduced by PCA and then average pooled followed by L2 normalization. The final image representation is the concatenation of the obtained vectors from both intermediate layers and fully connected layers. It is worth to note that the dimension of a VLAD feature is independent of the input image size. Meanwhile, the average pooling used in the fully connected layers can deal with input images of different sizes. Therefore, our method can handle images with different sizes, meaning more detail information can be used than the traditional methods fixing the input image size. Our main contributions include three aspects.

- 1) We propose a novel framework of using the rich hierarchical features of a CNN to form a discriminative image representation for scene classification. Our method incorporates features from low-level, middle-level, and high-level simultaneously; hence, we term our method as EMR for its such property.
- 2) Different from the traditional CNN-based feature extraction methods that have to fix the input image size to obtain a fixed dimensional feature, our method is flexible to images of different sizes owing to the used encoding strategies (i.e., VLAD and average pooling).
- 3) We give an experimental analysis about the similarities of local structures between natural images and remote sensing images, which shows that CNNs pretrained on the natural image dataset are more easily applied to the remote sensing dataset when the local structure similarity between two datasets is higher.

The rest of this paper is organized as follows. Section II briefly reviews the related works. Then, Section III elaborates our method, followed by experiments in Section IV. Finally, we conclude this paper in Section V.

II. RELATED WORKS

In recent years, there is a growing concern over scene classification in remote sensing imagery. The BOVW model is perhaps one of the most popular approaches, which encodes dense low-level features into a global semantic representation [1], [20]. Hu *et al.* [21] proposed an unsupervised feature learning framework to learn low-level features for the BOVW model. Zhu *et al.* [22] combined BOVW with global texture features to obtain a discriminative representation. Cheriadat [3] used sparse coding to extend the BOVW model for aerial scene classification. Zhao *et al.* [5] proposed a local Fisher kernel (LFK) framework to incorporate rich spatial information and strong discriminative ability. However, the gap between the low-level features and the semantic meanings limits their descriptive abilities to handle complex scenes [23].

CNN [24] is a trainable multilayer architecture, which incorporates multiple feature extraction stages from low level to high level. It usually consists of several convolutional layers, nonlinear layers, pooling layers, fully connected layers, and a loss layer. The convolutional layer computes filter responses for the input, which are input to the nonlinear layer subsequently.

These generated feature maps incorporate large amounts of low-level or middle-level features to contain abundant local structure and semantic information. The pooling layer reduces the spatial sizes of input feature maps by nonlinear down sampling and provides robustness to translation. After several stacked convolutional, nonlinear and pooling layers, the fully connected layers are commonly used to obtain high-level semantic representations. In the end, the loss layer is used at training time to penalize the deviation between predicted and true labels. The CNN architecture is trained in an end-to-end framework by back propagation, hence it can learn powerful semantic representation about objects from huge amounts of data. It also has good robustness to translation, scale, and distortion. Therefore, CNN has gained extensive attention and achieved best results on a range of vision tasks [12], [25]–[27].

Different CNN architectures have been designed to improve the discriminative ability while reducing the number of model parameters. AlexNet [9] considered to relieve the overfitting problem by some tricks, such as rectified linear units (ReLU) nonlinearity, data augmentation, and dropout. It was the first successfully trained deep CNN architecture and won the competition in the 2012 ImageNet Large Scale Visual Recognition Challenge (ILSVRC-2012). SPP-net [28] used the spatial pyramid pooling in the last feature maps and generated a fixed-length global semantic representation; hence, it had good robustness to different image sizes. Network-in-network (NIN) [29] increased the depth and expressiveness of the CNN architecture by using multilayer perceptrons to convolve the inputs and showed strong discriminative abilities in object recognition. Inspired by NIN, Szegedy *et al.* [15] designed the “inception module” to parallel convolve at different spatial scales and reduce the number of model parameters. Therefore, their final CNN architecture, GoogLeNet, could increase both the depth and width of the network while keeping the computational cost constant. In addition, there were many other CNN architectures with strong discriminative abilities, such as, OverFeat [30], CaffeNet [14], ZF-net [31], VGG-F, VGG-M, VGG-S [32], etc.

Using CNNs in remote sensing images has been reported with high performance. Zhao *et al.* [33] proposed multiscale CNN algorithm to learn spatial-related features for hyperspectral remote imagery classification. In [23], a gradient boosting random convolutional network framework was proposed to combine multiple CNNs for remote sensing scene classification. However, the number of labeled remote sensing scene images is usually far from sufficient for a deep CNN architecture training. This may be caused by three reasons.

- 1) Although it becomes easier to acquire a huge number of remote sensing images with the rapid development of satellite and sensor techniques, images with single purified category labels are difficult to be observed.
- 2) The manual annotation is too expensive and time consuming, which dispels researchers’ enthusiasm of labeling training samples.
- 3) Compared with traditional supervised methods, deep CNN based methods usually need more training samples to overcome the overfitting problem.

Therefore, researchers proposed to use a network architecture pretrained in natural scene images for initialization and then fine-tune the network by remote sensing images [13], [18]. Actually, directly using the features of CNNs pretrained in natural scene images to classify remote sensing scenes could also have excellent results [34]. Napoletano *et al.* [35] represented aerial images by the fully connected features of multiple CNNs pretrained on the ImageNet dataset and achieved better results than traditional methods.

Recently, encoding local features from a pretrained CNN has been widely used to extract a discriminative high-level semantic representation [36]–[38]. In these methods, the last convolutional feature maps are treated as a set of local features; then, these local features are encoded into a global semantic representation by some encoding methods [3], [5], [6]. Compared with traditional handcraft descriptors, such as SIFT, local binary patterns (LBP) [39], Haarlike, etc., local features from pretrained CNN could describe the semantic meanings more precisely; hence, the final image representation becomes more discriminative. In [40], multiscale images were input into a pretrained CNN architecture; then, local features of different scale images were combined together and further encoded into a single representation by Fisher Vector. Although their work successfully employed local features from a pretrained CNN to construct a global image representation and achieved good performance for remote sensing classification, it still ignored huge amounts of semantic information included in rich intermediate layers of CNNs. As described in [41], the low-level features from lower convolutional layers are full of specific structure information whereas the middle-level features from upper convolutional layers incorporate semantic information invariant to illumination, pose, location, etc. Therefore, using features from multiple convolutional layers could incorporate richer semantic information. Long *et al.* [42] proposed fully convolutional network (FCN) to handle arbitrary sizes of images and incorporate multilayer convolutional features to improve the semantic segmentation accuracy. However, FCN cannot be used for scene classification, which requires image level representation.

The spatial information is also an important consideration for constructing discriminative representations. Wang *et al.* [43] proposed the spatial latent Dirichlet allocation model to encode spatial structures among visual words. Yang *et al.* [1] extended the spatial pyramid match kernel to BOVW for remote sensing scene classification. They also proposed a novel spatial co-occurrence kernel to consider the relative spatial arrangement. Chen *et al.* [44] proposed a pyramid of spatial relations model to describe spatial relationships of low-level features. Inspired by their works [43], [1], [44], as shown in Section III-D, we also incorporate rich spatial information to improve the discriminative ability.

In this paper, we consider to use the features extracted from two kinds of the latest network architectures, VGG-VD [10] and ResNet [11]. Both of them are pretrained on the ImageNet dataset and have shown strong discriminations. Different from [34], [35], and [40], we take full advantage of semantic information contained in pretrained CNNs, i.e., both of the intermediate convolutional features and the fully connected features

are encoded to form a global image representation. As a result, our method could have stronger discrimination than the state of the arts, which is validated by our experiments.

III. FEATURE CONSTRUCTION

The pipeline of the proposed method is shown in Fig. 2. It can be found that the EMR is based upon the intermediate convolutional features and fully connected features from a pretrained CNN. In the following, we will first explain the reason why using CNNs pretrained on natural scene images can still achieve good performance for remote sensing scene classification. Then, we introduce two typical CNN architectures, VGG-VD and ResNet, used in our method. A brief introduction of the VLAD encoding technique is given afterward. Finally, we elaborate how to encode richer spatial and semantic information for a global image representation.

A. Analysis of the Similarity Between the Natural Imagery and Remote Sensing Imagery

The remote sensing images are acquired from the sky whereas the natural images are usually at a horizontal view; hence, objects in remote sensing images have completely different shapes and poses from natural images. Nevertheless, as evidenced in the literature, using features of CNNs pretrained in natural images could still achieve good performance for remote sensing scene classification [34], [40], [35]. In addition to the good generalization of CNNs, we claim that an important reason is the local structure similarity between remote sensing images and natural images. This is because that if the local similarity between natural images and remote sensing images is higher, the first several convolutional filters learned from natural images can more precisely describe local structure information in remote sensing images, and the descriptions of the latter convolutional layers would be more meaningful, thus encoding more plentiful and precise semantic information into the final representation, which is critical for a better classification performance.

To support this point, we statistically study the local similarity between natural and remote sensing images by three typical datasets. The ImageNet dataset [45] consists of huge amounts of natural images from 1000 object categories. It contains about 1.2 million training images, 50 000 validation images and 100 000 test images. We randomly pick up 2100 images from the training set, resize them to 224×224 , and calculate their local structure distribution. Note that when an image is input to a CNN, it usually subtracts the mean image of the training set for preprocessing [9]. We adopt the same strategy to preprocess all images in the following experiments; hence, the local structure distribution is calculated in preprocessed images. We use the uniform LBP ($LBP_{8,1}^u$) and rotation-invariant uniform LBP ($LBP_{8,1}^{riu2}$) to describe fundamental properties of local image structure. Due to the space limitation, please refer to [39] for more details about $LBP_{8,1}^u$ and $LBP_{8,1}^{riu2}$. We calculate histograms of $LBP_{8,1}^u$ and $LBP_{8,1}^{riu2}$ from all pixels in the above-mentioned 2100 natural images. Similarly, the histograms of $LBP_{8,1}^u$ and $LBP_{8,1}^{riu2}$ are also calculated on the UC-Merced [1] and Brazilian Coffee Scenes [34] datasets, which are remote

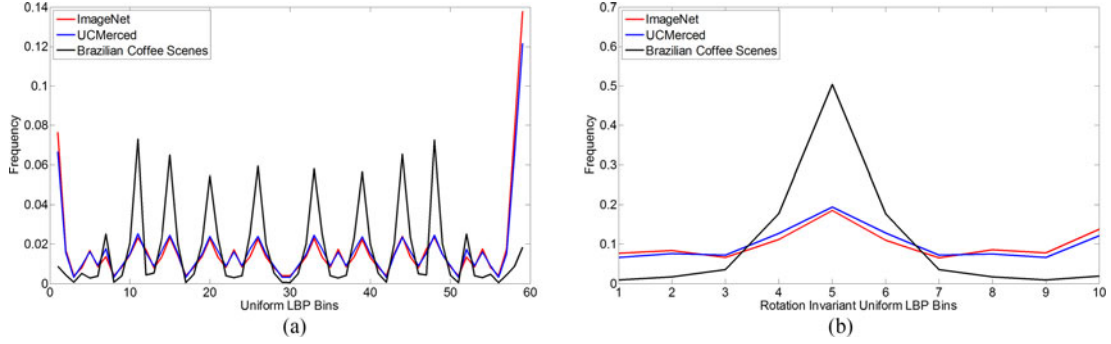


Fig. 3. Local structure distribution in natural images and remote sensing images: (a) histograms of $LBP^{u,2}_{8,1}$, 59 bins are used; and (b) histograms of $LBP^{riu,2}_{8,1}$, 10 bins are used.

TABLE I
INFLUENCE OF LOCAL STRUCTURE SIMILARITY

Dataset	The chi-square distance		Accuracy (%)			
	$LBP^{u,2}_{8,1}$	$LBP^{riu,2}_{8,1}$	Caffe [34]	OverFeat [34]	VGG-16	VGG-19
UC-Merced	0.0084	0.0069	93.42	90.91	90.67	89.62
Brazilian Coffee Scenes	0.4904	0.4860	84.82	81.2	86.78	85.98

sensing scene datasets. Fig. 3 presents the corresponding statistical results. The chi-square distance between histograms from UC-Merced and ImageNet as well as the chi-square distance between histograms from Brazilian Coffee Scenes and ImageNet are listed in Table I. It also contains classification accuracies on the UC-Merced and Brazilian Coffee Scenes datasets by using the last-layer fully connected features from CaffeNet, OverFeat, VGG-16, and VGG-19 pretrained on the ImageNet dataset. From Fig. 3 and Table I, it can be seen that compared with the Brazilian Coffee Scenes dataset, histograms of the UC-Merced dataset are more similar to those of ImageNet. As a result, such a similarity means that the convolutional filters learned from the ImageNet can be easier applied to UC-Merced than Brazilian Coffee Scenes and better classification performance can be achieved.

B. Typical CNNs

Due to the strong discriminative abilities shown in the ILSVRC classification task, two kinds of the latest CNN architectures, VGG-VD and ResNet, are used in our method, respectively. Both of them are implemented by MatConvNet [46] and pretrained on the ImageNet ILSVRC challenge dataset.

1) *VGG-VD*: Simonyan and Zisserman [10] designed very deep CNN architectures (VGG-VD) and won the second place in ILSVRC-2014. The most important design of their network architectures was the stack of convolutional layers with small receptive fields. Although the largest receptive field of all convolutional filters is 3×3 , stacking multiple convolutional layers could achieve a larger receptive field. This design could effectively increase the expressiveness by using more layers;

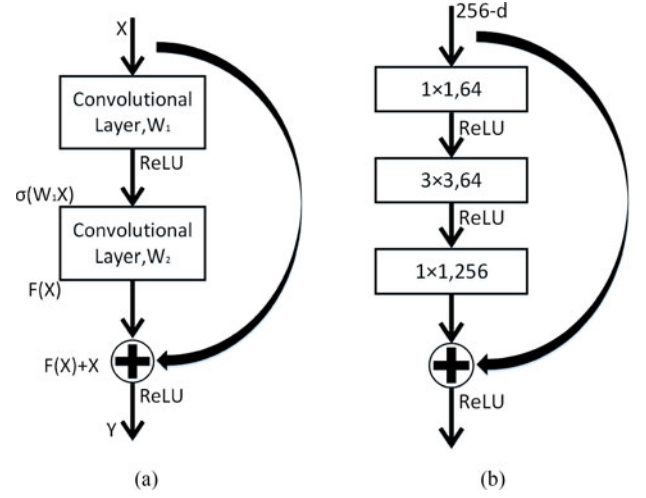


Fig. 4. Residual learning: (a) illustration of a residual learning block, and (b) practical block used in ResNet-152.

meanwhile, the number of parameters decreased rapidly, resulting in a deeper network with fewer parameters. Two successful VGG-VD architectures, known as VGG-16 (including 13 convolutional layers and 3 fully connected layers) and VGG-19 (including 16 convolutional layers and 3 fully connected layers), have showed strong discriminative abilities on the ImageNet ILSVRC challenge dataset; hence, we, respectively, use their pretrained network models to extract features in our method.

2) *ResNet*: He *et al.* [11] proposed the residual learning framework to ease CNN training; hence, the deeper network could be designed to increase the discrimination. As shown in Fig. 4(a), the feature map X is first input to a convolutional layer and a nonlinear layer, generating an intermediate feature map $\sigma(W_1 X)$. Here, σ is the nonlinear activation function ReLU [47], and W_i denotes the filter parameters of the i th convolutional layer and the biases are omitted for simplification. Then, $\sigma(W_1 X)$ is input to the following convolutional layer, obtaining a residual mapping $F(X) = W_2 \sigma(W_1 X)$. Finally, a shortcut connection is used to acquire an identity mapping $H(X) = F(X) + X$. Fig. 4(b) shows the practical residual block used in ResNet. In this paper, the pretrained ResNet-152 is used, which won the competition on the ILSVRC 2015 classification task.

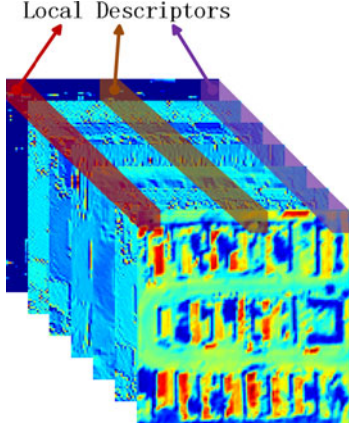


Fig. 5. Example of the convolutional feature maps.

Due to the space limitation, please refer to [10] and [11] for the detailed description of VGG-16, VGG-19, and ResNet-152. In our method, we extract features after the ReLU layers, because it achieves a slightly better performance than extracting features before the ReLU layers according to our experiments.

C. Vectors of Locally Aggregated Descriptors

Feature maps output by a convolutional layer can be viewed as a set of local descriptors. As shown in Fig. 5, the features across feature map channels at each location constitute a local descriptor, which describes semantic information of the corresponding spatial region (receptive field). Compared with traditional handcraft descriptors, local descriptors generated from convolutional feature maps incorporate richer semantic information; hence, encoding these descriptors into a global representation can help achieve a more discriminative feature vector. In this paper, we choose to use the well-studied VLAD to encode these descriptors.

Suppose $V = \{v_1, v_2, \dots, v_k\}$ is a codebook with k visual words, which is usually learned from a training set of local descriptors by k -means. For a local descriptor set $P = \{p_1, p_2, \dots, p_m\}$ extracted from an image, its locally aggregated descriptors (VLAD) feature [19] can be calculated as follows:

$$q_{i,j} = \sum_{p_l \text{ such that } NN(p_l)=v_i} p_{l,j} - v_{i,j} \quad (1)$$

where $p_{l,j}$ and $v_{i,j}$, respectively, denote the j th components of the local feature p_l and the visual word v_i , $NN(p_l) = v_i$ means that v_i is the nearest neighbor visual word of p_l . The final VLAD feature Q is a single vector concatenating all $q_{i,j}$ s and normalized to unit length. Therefore, for a d dimensional local feature with k visual words, the dimension of VLAD feature will be $k \times d$. When we use VLAD to encode intermediate convolutional features from pretrained CNNs, the dimension of the generated VLAD feature is independent of the size of intermediate convolutional features. As a result, for different sizes of input images, we can obtain a fixed dimensional VLAD feature. In this paper, the VLAD encoding technique is implemented

Algorithm 1: Encoded Mixed-Resolution Representation.

Input: Input Image I

Output: EMR Feature F

- 1: Input I to a pretrained CNN, the last convolutional feature maps in different resolution and the fully connected features are reserved.
 - 2: **for all** Reserved intermediate convolutional features **do**
 - 3: Reduce local features by PCA.
 - 4: **for all** Spatial divisions **do**
 - 5: Encode the corresponding local features by VLAD.
 - 6: **end for**
 - 7: Concatenate all VLAD vectors from different divisions and reduce its dimension by PCA.
 - 8: **end for**
 - 9: **for all** Reserved fully connected features **do**
 - 10: Reduce fully connected features by PCA.
 - 11: Average pooling these reduced features.
 - 12: L2-normalize the pooled feature.
 - 13: **end for**
 - 14: Concatenate all encoded features from intermediate convolutional layers and fully connected layers to form the final feature F .
-

by vlfeat [48]; meanwhile, a kd-tree is built to quickly find the nearest neighbor visual word of each local feature.

D. Encoded Mixed-Resolution Representation

Algorithm 1 gives the pseudocode of our proposed method and Fig. 2 shows the whole procedure by an example of using a pretrained VGG-16 network.

Due to the existence of pooling layers, CNNs can generate convolutional feature maps with different sizes. Obviously, these feature maps have different resolutions. From the first convolutional layer to the last convolutional layer, the generated feature maps have lower and lower resolutions with higher and higher semantic information. For example, feature maps of the first convolutional layer usually describe 3×3 or 5×5 regions and tend to extract low-level structure information, such as, corners and edges; feature maps of the latter convolutional layers can describe larger spatial regions, representing some parts of objects or even whole objects. Therefore, encoding local descriptors from convolutional feature maps in different resolutions can incorporate multiple levels of semantic and spatial information. In a deep CNN network, there are several convolutional layers exporting feature maps in the same resolution. In this case, we use the feature maps from the last convolutional layer in each resolution as they have the strongest expressiveness. Meanwhile, in order to incorporate richer spatial information, the feature maps can be spatially divided into multiple groups for generating their VLAD representations. As shown in Fig. 2, four kinds of spatial divisions are used and 11 spatial groups are generated accordingly. We extract VLAD features from these spatial groups and concatenate them into a single vector. As the feature dimension becomes too high, we use PCA to reduce its dimension. Finally, the computed features from different intermediate convolutional

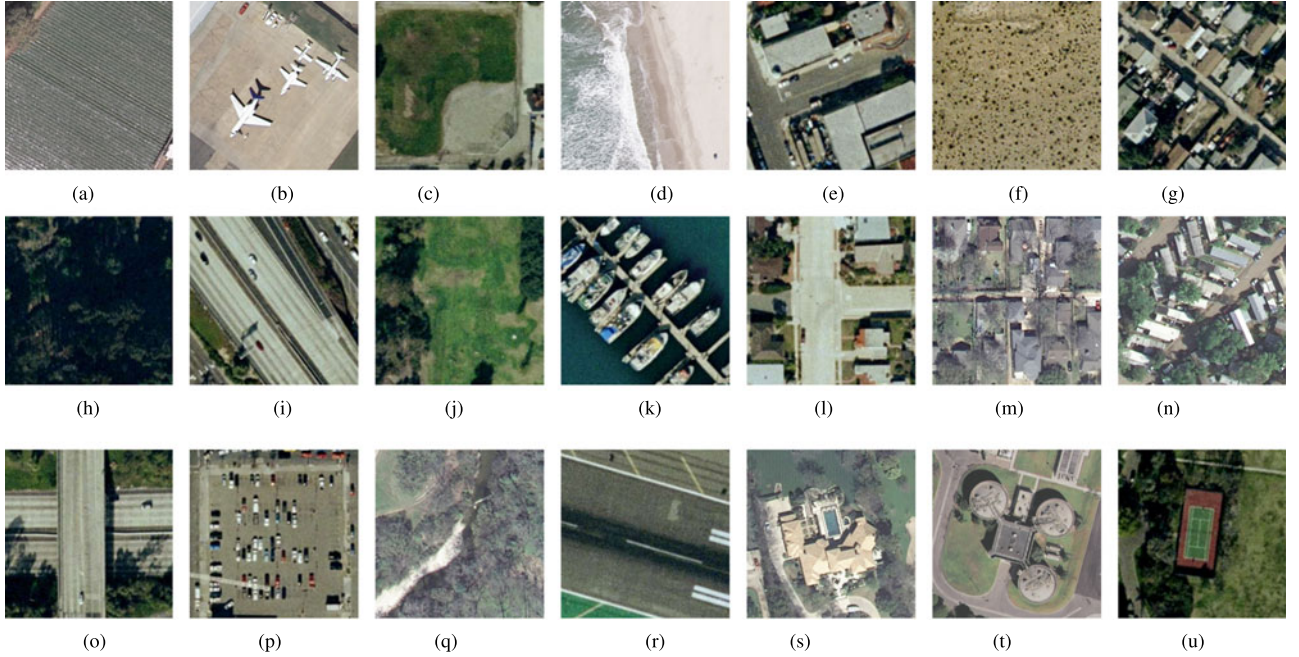


Fig. 6. Different scene categories contained in the UC-Merced dataset: (a) agricultural, (b) airplane, (c) baseball diamond, (d) beach, (e) buildings, (f) chaparral, (g) dense residential, (h) forest, (i) freeway, (j) golf course, (k) harbor, (l) intersection, (m) medium residential, (n) mobile home park, (o) overpass, (p) parking lot, (q) river, (r) runway, (s) sparse residential, (t) storage tanks, (u) tennis court.

layers are concatenated together as a single vector and we call it as encoded convolutional feature.

In addition to the intermediate convolutional features, the fully connected features of pretrained CNNs provide high-level generalized representations, which have shown strong discrimination in remote sensing classification [34], [35]. Therefore, we also use features from the fully connected layers in the proposed framework. To handle input images with different sizes, we consider to average pool features in each fully connected layer. Specifically, features from each fully connected layer are first reduced by PCA; then these reduced features in each fully connected layer are respectively average pooled to obtain a robust feature; finally, the pooled feature of each fully connected layer is L2 normalized. Fig. 2 has shown specific feature dimensions of the VGG-16 layers with an input image of $256 \times 256 \times 3$. It can be found that the average pooling strategy can eliminate the influence of the input image size and thus obtaining a fixed dimensional feature vector. The encoded fully connected feature concatenates normalized features of all fully connected layers.

The entire feature of our method is an alliance of the encoded convolutional feature and the encoded fully connected feature. It is worth to note that although PCA has been used many times in different stages, their projections are different. As our method encodes features from both the intermediate layers and the fully connected layers of a pretrained CNN, we call it as **Encoded Mixed-resolution Representation (EMR)**. By using VGG-16, VGG-19, and ResNet-152 as the pretrained CNNs, respectively, we obtain three corresponding features as VGG16_EMR, VGG19_EMR, and ResNet152_EMR, which are evaluated on two remote sensing classification datasets in the following experiments.



Fig. 7. Examples of the Brazilian Coffee Scenes dataset: (a) coffee, (b) non-coffee.

IV. EXPERIMENTS

A. Experimental Setup

To evaluate the performance of our method in the task of remote sensing scene classification, we conduct experiments on two widely used datasets. We report the accuracy rate calculated as the number of correct classifications divided by the total number of classifications.

We first evaluate on the UC-Merced dataset [1], which collects land-use images from the USGS National Map Urban Area Imagery with a pixel resolution of one foot. It consists of 21 distinctive scene categories, as shown in Fig. 6. Each class contains 100 images with the fixed size of 256×256 . Similar to the experimental setting in [1], fivefold cross validation is performed. Specifically, the dataset is first randomly partitioned into five equal-sized subgroups, each of which contains 20 images per category; then, the classifier is circularly trained on four subgroups and evaluated on the remaining subgroup; finally, the classification accuracy rate is the average over the above-mentioned five evaluation results.

In addition to UC-Merced dataset, we also evaluate on the Brazilian Coffee Scenes dataset [34], which includes two categories of scenes: coffee and noncoffee. Fig. 7 shows some examples. In this dataset, all images are multispectral ones. Each

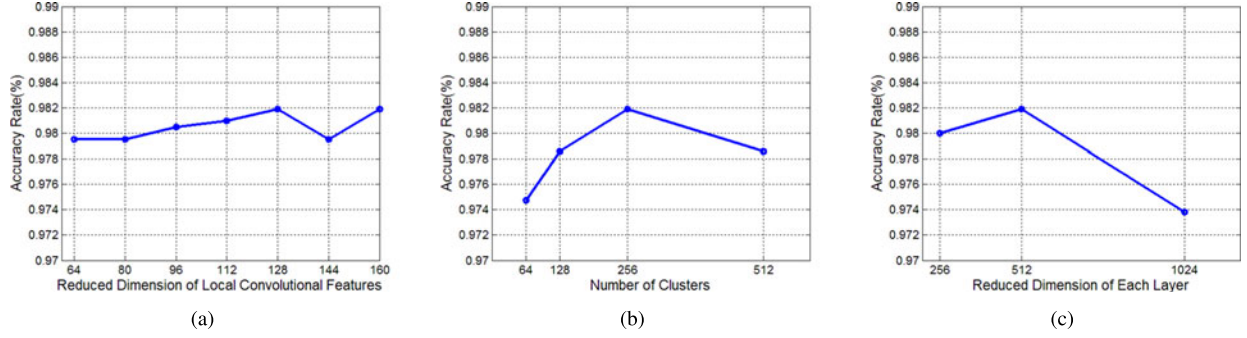


Fig. 8. Parameter evaluation for VGG16_EMR.

image incorporates the red, green, and near infrared bands with a fixed size of 64×64 . It contains five folds, four of which have 600 samples and the remaining one has 476 samples. The numbers of coffee and noncoffee samples in each fold are the same. Fivefold cross validation is used too.

The extracted image representations are combined with the linear SVM (implemented by liblinear [49]) to distinguish different scene categories. Although this processing is simple and standard, our method can still achieve the best classification performance in the following experiments, which reflects the strong discriminative ability of the proposed feature.

B. Parameter Evaluation

1) *Parameter Settings*: Three important parameters exist in the proposed feature: the reduced dimension of local convolutional features d by PCA, the number of clusters k , and the reduced dimension of each layer l by PCA (for simplicity, the reduced dimension of local fully connected features is set to l too). We use the UC-Merced dataset to estimate the most suitable parameter settings for our method. Since these parameters are independent of the used CNN architecture, we only evaluate the influence of different parameters by VGG-16.

When evaluating each parameter, we keep the other parameters at the best settings. Fig. 8 shows the accuracy rates with different parameter settings. It can be seen that the highest accuracy rate can be achieved at $d = 128$ and $d = 160$, thus we set $d = 128$ for its lower dimension. For convolutional layers whose feature dimension is lower than d , their original dimension is retained. From Fig. 8(b), we can find that increasing the number of clusters can notably improve VGG16_EMR's discriminative ability. The highest accuracy rate is reached when the number of clusters reaches 256. Then, the accuracy rate decreases if continue increasing the number, as a consequence, we set the number of clusters as $k = 256$. Accordingly, the dimension of the generated VLAD feature is equal to 32768 and concatenating VLAD vectors from different spatial divisions forms a 360448 dimensional feature vector for each convolutional layer. However, as shown in Fig. 8(c), the reduced feature of each layer only needs a low dimension of $l = 512$. Finally, we set $d = 128$, $k = 256$, and $l = 512$ in the following experimental evaluations.

2) *Influence of Different CNN Layers*: In this part, we analyze the discriminative abilities of features extracted from

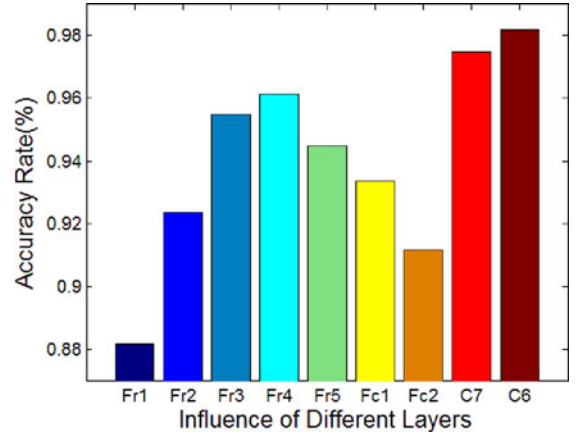


Fig. 9. Accuracy rates of features extracted from different layers of VGG-16 on the UC-Merced dataset.

different layers of VGG-16. To this end, we denote the reduced VLAD features in different resolutions as $Fr1, Fr2, \dots, Fr5$, separately. As shown in Fig. 2, $Fr1$ is from the largest resolution whereas $Fr5$ is from the smallest. The normalized features from the first two fully connected layers are, respectively, denoted as $Fc1$ and $Fc2$ (note that $Fc2$ is different from the fully connected feature of VGG-16 used in Section III-A because of the different input image sizes, i.e., one is 256×256 and the other is 224×224). Furthermore, $C7$ means the feature combining $Fr1, Fr2, \dots, Fr5, Fc1$ and $Fc2$ whereas $C6$ only combines $Fr2, Fr3, \dots, Fr5, Fc1$ and $Fc2$. Fig. 9 shows their classification results on the UC-Merced dataset and we have several observations as follows.

- 1) From $Fr1$ to $Fr4$, the accuracy rate has significantly increased, indicating the improved discriminative ability of convolutional features with deeper network.
- 2) The performance of $Fr5$ is poorer than $Fr3$ and $Fr4$. This may because the number of local convolutional features is a bit less (16×16), which cannot provide adequate semantic information.
- 3) The performance of $Fr3, Fr4$, and $Fr5$ is better than $Fc1$ and $Fc2$, which shows the significance of using semantic information in intermediate layers.
- 4) The accuracy rate of $Fc1$ is higher than that of $Fc2$, demonstrating a better generalization ability of the first fully connected layer than the second fully connected layer.

TABLE II
INFLUENCE OF DIFFERENT PRETRAINED CNNs

CNN	AlexNet	VGG-F	VGG-M	VGG-S	VGG-16	VGG-19	ResNet-152
Accuracy Rate(%)	96.86	97.05	96.95	96.95	98.19	97.62	98.76

5) If combining features of convolutional layers and fully connected layers ($C6$ and $C7$), the better classification performance can be obtained. However, the performance of $C7$ is a little poorer than $C6$ owing to the weaker discrimination of $Fr1$.

In summary, the intermediate layers of CNNs incorporate rich semantic information, thus by combining features extracted from multiple CNN layers, a feature with very strong discriminative ability can be obtained. Meanwhile, convolutional features in the largest resolution lack sufficient discriminative abilities. Therefore, we encode fully connected features and all intermediate convolutional features except for the largest resolution to construct VGG16_EMR, VGG19_EMR, and ResNet152_EMR.

3) *Influence of Pretrained CNNs*: Besides VGG-16, VGG-19, and ResNet-152, we also tried other CNNs in our proposed framework and evaluated their performance on the UC-Merced dataset. These networks include AlexNet [9], VGG-F, VGG-M, and VGG-S [32]. All of them are implemented by MatConvNet and pretrained on the ImageNet dataset too. The classification results are presented in Table II. It can be seen that using VGG-F, VGG-M, and VGG-S perform slightly better than using AlexNet, but worse than using VGG-16 and VGG-19. The features extracted by ResNet-152 achieves the best accuracy. It can be found that these results are exactly similar to the classification performance on the ImageNet dataset. This illustrates that if a CNN achieves better performance on the classification task, its intermediate features in each layer usually have stronger discrimination; hence, our proposed feature based on this CNN architecture can achieve a higher performance. This is the reason why we use VGG-16, VGG-19, and ResNet-152 to construct EMR features in this paper.

4) *Influence of Input Image Sizes*: In this part, we evaluate the influence of different input image sizes to EMR. The input sizes of VGG-16, VGG-19, and ResNet-152 are all 224×224 whereas the image size in the UC-Merced dataset is 256×256 ; hence, we evaluate the influence of input image sizes of both 224×224 and 256×256 pixels. The commonly used image size 320×320 is also included in our evaluation. For the Brazilian Coffee Scene dataset, its image size, 64×64 , is much smaller than the above-mentioned three sizes. Even though, we evaluate these different image sizes on the Brazilian Coffee Scenes dataset to study the influence of excessive zoom. Table III shows their classification accuracy rates. It can be found that VGG16_EMR and ResNet152_EMR achieve higher accuracy rates with size of 256×256 than those of 224×224 on the UC-Merced dataset. The possible reason is that when resizing images to smaller ones, some semantic information has been discarded. When resizing images in the UC-Merced dataset to 320×320 , VGG19_EMR

TABLE III
INFLUENCE OF DIFFERENT INPUT IMAGE SIZES

	The UC-Merced dataset			The Brazilian Coffee Scenes dataset		
	224	256	320	224	256	320
VGG16_EMR	97.76	98.19	98.14	92.28	92.03	91.80
VGG19_EMR	97.95	97.62	97.90	91.60	92.25	91.03
ResNet152_EMR	98.38	98.76	98.90	99.00	96.26	96.31

TABLE IV
CLASSIFICATION ACCURACY RATES OF DIFFERENT METHODS ON THE UC-MERCE DATASET

		Method	Accuracy(%)
Feature-based	PSR [44]		89.10
	UFL-SC [21]		90.26
	SIFT + SC [3]		81.67
	FK-S [5]		91.63
	VLAD [6]		92.50
	VLAT [6]		94.30
	LGFBOVW [22]		96.88
	GoogLeNet [13]		97.10
Network-based	GBRCN [23]		94.53
	Multiview deep learning [50]		93.48
	CaffeNet [34]		93.42
	OverFeat [34]		90.91
	OverFeat + Caffe [34]		99.43
	VGG-S + IFK(VGG-16) [40]		98.49
	VGG16_EMR		98.14
	VGG19_EMR		97.90
	ResNet152_EMR		98.90
	ResNet152(intermediate) + VGG16(fully connected)		99.43
	ResNet152(intermediate) + VGG19(fully connected)		99.48

and ResNet152_EMR perform slightly better than the original size. However, VGG16_EMR and ResNet152_EMR become worse when resizing the image to larger sizes on the Brazilian Coffee Scenes dataset. Based on the above-mentioned results, we evaluate our proposed method with the input image sizes of 320×320 on the UC-Merced dataset and 224×224 on the Brazilian Coffee Scenes dataset in the following experiments.

C. Comparison With the State of the Art

1) *UC-Merced*: As shown in Table IV, we compare our proposed features with the state of the arts on the UC-Merced dataset. The competitive methods can be classified into two categories: feature-based and network-based. The feature-based methods usually encode dense local descriptors into a global representation by different feature encoding methods, such as sparse coding [3], Fisher Vector [5], VLAD [6], etc., [21], [44]. However, due to the limited descriptive abilities of local descriptors, these methods do not work very well, which can be seen from Table IV. To alleviate the problem of limited descriptive abilities of using a single local descriptor, LGFBOVW [22] mixed several local descriptors in the BOVW framework and achieves the accuracy rate at 96.88%. Compared to the feature-based methods, the network-based methods perform much better. In [23] and [50], Zhang *et al.* and Luus *et al.* designed special CNN architectures for scene classification. However,

TABLE V
CLASSIFICATION ACCURACY RATES OF DIFFERENT METHODS ON THE
BRAZILIAN COFFEE SCENES DATASET

Method	Accuracy(%)
BIC [34]	87.0
BOVW + SIFT [34]	80.5
OverFeat [34]	84.82
CaffeNet [34]	81.2
OverFeat + CaffeNet [34]	83.04
GoogLeNet [13]	91.83
VGG16_EMR	92.28
VGG19_EMR	91.60
ResNet152_EMR	99.00

their network depths are insufficient; hence, their final classification performance is not good enough. Using pretrained GoogLeNet and further fine-tuning it on the UC-Merced dataset can help achieve a good performance with 97.10% accuracy rate, which is superior to all feature-based methods. Directly using the fully connected features from either CaffeNet or OverFeat to classify remote sensing scenes, it cannot achieve a good classification performance. However, once concatenating both of them, it can reach a surprising accuracy rate 99.43%. VGG-S+IFK(VGG-16) [40], which combines the fully connected feature from pretrained VGG-S and the Fisher vector encoding last convolutional features from pretrained VGG-16, also obtains a high classification accuracy rate up to 98.49%. Our proposed feature, VGG16_EMR and VGG19_EMR, achieve accuracy rates at 98.14% and 97.90%, respectively. Although their performance is worse than OverFeat+Caffe and VGG-S+IFK(VGG-16), they achieve the best performance on the UC-Merced dataset among all the methods using a single CNN model. ResNet152_EMR performs better than VGG16_EMR and VGG19_EMR, just slightly worse than OverFeat+Caffe. Because the fully connected feature of ResNet-152 is average pooled from its lower convolutional layer rather than learned from the dataset, combining it with its convolutional layer features is not able to encode additional information. When we replace the normalized features in ResNet152_EMR by the normalized features in VGG16_EMR or VGG19_EMR to form a mixed feature, higher accuracy rate (99.43% or 99.48%) can be achieved. The highest accuracy rate 99.48% is achieved by ResNet152(intermediate)+VGG19(fully connected).

2) *Brazilian Coffee Scenes*: The results are listed in Table V. The BOVW with dense SIFT performs the worst in this dataset. The simple color feature, border-interior pixel classification (BIC), improves the accuracy rate by 6.5%. Using fully connected features from pretrained OverFeat and CaffeNet achieve accuracy rates no more than 84.82%, which are not as good as their performance on the UC-Merced dataset. The reason is mainly caused by the difference between the characteristic of these two datasets. The spectral information in Brazilian Coffee Scenes is entirely different from UC-Merced, which is more similar to ImageNet. Such a different spectral information further results in a significant difference on local texture. As shown in Section III-A, the local similarity between Brazilian Coffee Scenes and ImageNet is much lower than that between UC-

Merced and ImageNet. Therefore, the fully connected features learned from natural images can be easily applied on the UC-Merced dataset than on the Brazilian Coffee Scenes dataset. Our proposed features, VGG16_EMR, VGG19_EMR, and ResNet152_EMR, can still achieve good performance with accuracy rates: 92.28%, 91.60%, and 99.00%. This indicates the importance of encoding intermediate convolutional features. Although GoogLeNet is fine-tuned on the Brazilian Coffee Scenes dataset, it still performs much worse than ResNet152_EMR, which validates the strong discrimination of our proposed method.

V. CONCLUSION

This paper proposes a novel image representation method named EMR to classify remote sensing scene images. To encode rich semantic information contained in CNNs, we consider to fully employ the low-level, middle-level, and high-level features from the intermediate convolutional features and fully connected features of CNNs. Specifically, the low-level and middle-level features from intermediate convolutional layers are, respectively, encoded into a global image representation by VLAD, followed by a PCA to reduce the feature dimension. Meanwhile, the fully connected features are average pooled and subsequently normalized to form new global features. The EMR combines all the above-mentioned global features. Owing to the good properties of VLAD and average pooling, our method can handle input image with different sizes. In addition, we have statistically analyzed that CNNs pretrained on the natural image dataset are more easily applied to the remote sensing dataset when the local structure similarity between two datasets is higher. With our proposed encoding method, using CNNs pretrained on the natural image dataset can also help achieve very good performance on remote sensing images even if their local structure similarity to natural images is relatively low.

REFERENCES

- [1] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. ACM Geograph. Inf. Syst.*, 2010, pp. 270–279.
- [2] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.
- [3] A. M. Cheriyyadat, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan. 2014.
- [4] M. L. Mekhalif, F. Melgani, Y. Bazi, and N. Alajlan, "Land-use classification with compressive sensing multifeature fusion," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 10, pp. 2155–2159, Oct. 2015.
- [5] B. Zhao, Y. Zhong, L. Zhang, and B. Huang, "The fisher kernel coding framework for high spatial resolution scene classification," *Remote Sens.*, vol. 8, no. 2, p. 157, 2016.
- [6] R. Negrel, D. Picard, and P.-H. Gosselin, "Evaluation of second-order visual features for land-use classification," in *Proc. IEEE 12th Int. Workshop Content-Based Multimedia Indexing.*, 2014, pp. 1–5.
- [7] E. Li, P. Du, A. Samat, Y. Meng, and M. Che, "Mid-level feature representation via sparse autoencoder for remotely sensed scene classification," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 10, no. 3, pp. 1068–1081, Mar. 2017.
- [8] J. Chen, L. Jiao, and Z. Wen, "High-level feature selection with dictionary learning for unsupervised sar imagery terrain classification," *J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 10, no. 1, pp. 145–160, Jan. 2017.

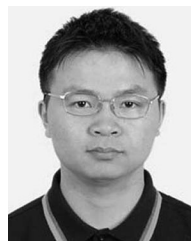
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [10] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015. arXiv:1512.03385.
- [12] Z. Yan *et al.*, "Hd-cnn: hierarchical deep convolutional neural networks for large scale visual recognition," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 2740–2748.
- [13] M. Castelluccio, G. Poggi, C. Sansone, and L. Verdoliva, "Land use classification in remote sensing images by convolutional neural networks," 2015. arXiv:1508.00092.
- [14] Y. Jia *et al.*, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. 22nd ACM Int. Conf. Multimedia*, 2014, pp. 675–678.
- [15] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 1–9.
- [16] I. Ševo and A. Avramović, "Convolutional neural network based automatic object detection on aerial images," *IEEE Geosci. Remote Sens. Soc. Lett.*, vol. 13, no. 5, pp. 740–744, May 2016.
- [17] G. Cheng, P. Zhou, and J. Han, "Learning rotation-invariant convolutional neural networks for object detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7405–7415, Dec. 2016.
- [18] K. Nogueira, O. A. Penatti, and J. A. dos Santos, "Towards better exploiting convolutional neural networks for remote sensing scene classification," *Pattern Recognit.*, vol. 61, pp. 539–556, 2017.
- [19] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *Proc. IEEE Comput. Vis. Pattern Recognit.*, 2010, pp. 3304–3311.
- [20] Z. Zhou, Y. Wang, Q. J. Wu, C.-N. Yang, and X. Sun, "Effective and efficient global context verification for image copy detection," *IEEE Trans. Inf. Forensics Security*, vol. 12, no. 1, pp. 48–63, Jan. 2017.
- [21] F. Hu, G.-S. Xia, Z. Wang, X. Huang, L. Zhang, and H. Sun, "Unsupervised feature learning via spectral clustering of multidimensional patches for remotely sensed scene classification," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 8, no. 5, pp. 2015–2030, May 2015.
- [22] Q. Zhu, Y. Zhong, B. Zhao, G.-S. Xia, and L. Zhang, "Bag-of-visual-words scene classifier with local and global features for high spatial resolution remote sensing imagery," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 6, pp. 747–751, Jun. 2016.
- [23] F. Zhang, B. Du, and L. Zhang, "Scene classification via a gradient boosting random convolutional network framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 3, pp. 1793–1802, Mar. 2016.
- [24] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *IEEE Proc.*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [25] S. Zha, F. Luisier, W. Andrews, N. Srivastava, and R. Salakhutdinov, "Exploiting image-trained CNN architectures for unconstrained video classification," 2015. arXiv:1503.04144.
- [26] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 91–99.
- [27] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, and S. Reed, "SSD: Single shot multibox detector," 2015. arXiv:1512.02325.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," in *Eur. Conf. Comput. Vis.*, Berlin, Germany: Springer, 2014, pp. 346–361.
- [29] M. Lin, Q. Chen, and S. Yan, "Network in network," 2013. arXiv:1312.4400.
- [30] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," 2013. arXiv:1312.6229.
- [31] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis.*, Berlin, Germany: Springer, 2014, pp. 818–833.
- [32] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," 2014. arXiv:1405.3531.
- [33] W. Zhao and S. Du, "Learning multiscale and deep representations for classifying remotely sensed imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 113, pp. 155–165, 2016.
- [34] O. A. Penatti, K. Nogueira, and J. A. dos Santos, "Do deep features generalize from everyday objects to remote sensing and aerial scenes domains?" in *Proc. IEEE Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 44–51.
- [35] P. Napoletano, "Visual descriptors for content-based retrieval of remote sensing images," 2016. arXiv:1602.00970.
- [36] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *Eur. Conf. Comput. Vis.*, Berlin, Germany: Springer, 2014, pp. 392–407.
- [37] Y. Cheng, R. Cai, X. Zhao, and K. Huang, "Convolutional fisher kernels for RGB-D object recognition," in *Proc. IEEE Int. Conf. 3D Vis.*, 2015, pp. 135–143.
- [38] D. Yoo, S. Park, J.-Y. Lee, and I. So Kweon, "Multi-scale pyramid pooling for deep convolutional representation," in *Proc. IEEE Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 71–80.
- [39] T. Ojala, M. Pietikainen, and T. Maenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Analysis Mach. Intell.*, vol. 24, no. 7, pp. 971–987, Jul. 2002.
- [40] F. Hu, G.-S. Xia, J. Hu, and L. Zhang, "Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery," *Remote Sens.*, vol. 7, no. 11, pp. 14 680–14 707, 2015.
- [41] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik, "Hypercolumns for object segmentation and fine-grained localization," in *Proc. IEEE Comput. Vis. Pattern Recognit.*, 2015, pp. 447–456.
- [42] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Comput. Vis. Pattern Recognit.*, 2015, pp. 3431–3440.
- [43] X. Wang and E. Grimson, "Spatial latent Dirichlet allocation," in *Proc. Adv. Neural Inf. Process. Syst.*, 2008, pp. 1577–1584.
- [44] S. Chen and Y. Tian, "Pyramid of spatial relations for scene-level land use classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, pp. 1947–1957, Apr. 2015.
- [45] O. Russakovsky *et al.*, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [46] A. Vedaldi and K. Lenc, "Matconvnet: Convolutional neural networks for MATLAB," in *Proc. 23rd ACM Int. Conf. Multimedia*, 2015, pp. 689–692.
- [47] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proc. 27th Int. Conf. Mach. Learn.*, 2010, pp. 807–814.
- [48] A. Vedaldi and B. Fulkerson, "VLFeat: An open and portable library of computer vision algorithms," in *Proc. ACM Multimedia*, 2010, pp. 1469–1472.
- [49] R. Fan, K. Chang, C. Hsieh, X. Wang, and C. Lin, "LIBLINEAR: A library for large linear classification," *J. Mach. Learn. Res.*, vol. 9, pp. 1871–1874, 2008.
- [50] F. Luus, B. Salmon, F. Van Den Bergh, and B. Maharaj, "Multiview deep learning for land-use classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 12, pp. 2448–2452, Dec. 2015.



Guoli Wang received the B.S. degree in electronic information science and technology from China University of Mining and Technology, Xuzhou, China, in 2007. He is currently working toward the Ph.D. degree in pattern recognition and intelligent systems in the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China.

He is currently also in the School of Computer and Control Engineering, University of Chinese Academy of Sciences, Beijing, China. His research

interests include remote sensing image processing, pattern recognition, and machine learning.



Bin Fan (M'10–SM'16) received the B.S. degree in automation from Beijing University of Chemical Technology, Beijing, China, in 2006, and the Ph.D. degree in pattern recognition and intelligent systems from the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China, in 2011.

He is currently an Associate Professor in the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences. His research interests include remote sensing image

processing, computer vision, and pattern recognition.

Dr. Fan is an Associate Editor of *Neurocomputing* and was an Area Chair of IEEE Winter Conference on Applications of Computer Vision 2016.



Shiming Xiang received the B.S. degree in mathematics and the M.S. degree in computational mechanics from Chongqing Normal University, Chongqing, China, in 1993 and 1996, respectively, and the Ph.D. degree in computer application technology from the Institute of Computing-Technology, Chinese Academy of Sciences, Beijing, China, in 2004.

From 1996 to 2001, he was a Lecturer at the Huazhong University of Science and Technology, Wuhan, China. From 2004 to 2006, he was a Postdoctorate Candidate in the Department of Automation, Tsinghua University, Beijing, China. He is currently a Professor in the Institute of Automation, Chinese Academy of Sciences, Beijing, China. His interests include pattern recognition and machine learning.



Chunhong Pan (M'14) received the B.S. degree in automatic control from Tsinghua University, Beijing, China, in 1987, the M.S. degree in optical information processing from the Shanghai Institute of Optics and Fine Mechanics, Chinese Academy of Sciences, Beijing, China, in 1990, and the Ph.D. degree in pattern recognition and intelligent system from the Institute of Automation, Chinese Academy of Sciences, in 2000.

He is currently a Professor in the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences. His research interests include computer vision, image processing, computer graphics, and remote sensing.