

Object Detection in Optical Remote Sensing Images Based on Weakly Supervised Learning and High-Level Feature Learning

Junwei Han, Dingwen Zhang, Gong Cheng, Lei Guo, and Jinchang Ren

Abstract—The abundant spatial and contextual information provided by the advanced remote sensing technology has facilitated subsequent automatic interpretation of the optical remote sensing images (RSIs). In this paper, a novel and effective geospatial object detection framework is proposed by combining the weakly supervised learning (WSL) and high-level feature learning. First, deep Boltzmann machine is adopted to infer the spatial and structural information encoded in the low-level and middle-level features to effectively describe objects in optical RSIs. Then, a novel WSL approach is presented to object detection where the training sets require only binary labels indicating whether an image contains the target object or not. Based on the learnt high-level features, it jointly integrates saliency, intraclass compactness, and interclass separability in a Bayesian framework to initialize a set of training examples from weakly labeled images and start iterative learning of the object detector. A novel evaluation criterion is also developed to detect model drift and cease the iterative learning. Comprehensive experiments on three optical RSI data sets have demonstrated the efficacy of the proposed approach in benchmarking with several state-of-the-art supervised-learning-based object detection approaches.

Index Terms—Bayesian framework, deep Boltzmann machine (DBM), object detection, weakly supervised learning (WSL).

I. INTRODUCTION

THE rapid development of remote sensing technologies has rendered many satellite and aerial sensors to provide optical imagery with high spatial resolution, facilitating a wide range of applications such as disaster control, land planning, urban monitoring, and traffic planning [1]–[3]. In these applications, automatic detection of natural or man-made objects is a fundamental task and has received increasing research interests.

Early attempts [2]–[4] detected objects in optical remote sensing images (RSIs) in an unsupervised manner, which often

started from generating region of interest by grouping pixels into clusters and then detected objects of interest based on the shape and spectral information. Afterward, many supervised learning methods have been adopted to learn the object model effectively with the help of prior information obtained from training examples [1], [5], [6]. By heavily relying on the human-labeled training examples, which are statistically representatives of the classification problem to solve, the supervised learning methods can achieve more promising performance than the unsupervised approaches.

The recent advance of remote sensing technology has led to the explosive growth of satellite and aerial images in both quantity and quality. It brings about two increasingly serious problems for the object detection task in optical RSIs. First, supervised-learning-based object detection approaches often require a large number of training data with manual annotation of labeling a bounding box around each object to be detected. However, manual annotation of objects in large image sets is generally expensive and sometimes even unreliable. For example, for the natural objects such as landslide, the proper manual annotation generally requires considerable expertise. In addition, manual annotation is also difficult for the man-made objects such as airplane and car, where the coverage of target object appears to be very small, particularly when complex textures are contained in the image background. As a result, it is difficult to achieve accurate annotation on such small regions. Moreover, the manual annotations may tend to be less accurate and unreliable when the targets are occluded or camouflaged. As a result, it is a great interest in training object detectors with weak supervision for large-scale optical satellite and aerial image data sets.

The second problem is that the rich information contained in the optical RSIs with high spatial resolution has more details of objects, whereas feature descriptors used by existing object detectors are still insufficiently powerful to characterize the structural information of the objects. The limited understanding of the spatial and structural patterns of objects in optical RSIs leads to a tremendous semantic gap for the object detection task. It can be observed that man-made facilities, such as airplanes, vehicles, and airports, always have intrinsic structural property with specific semantic concepts, which has obvious difference from the background areas in optical RSIs. Consequently, building of the high-level structural features is a promising way for object detection task.

Manuscript received May 23, 2014; revised September 8, 2014; accepted October 31, 2014. This work was supported in part by the National Science Foundation of China under Grants 91120005, 61473231, and 61401357.

J. Han, D. Zhang, G. Cheng, and L. Guo are with the School of Automation, Northwestern Polytechnical University, Xi'an 710072, China (e-mail: junwei.han2010@gmail.com).

J. Ren is with the Department of Electronic and Electrical Engineering, University of Strathclyde, Glasgow G1 1XW, U.K.

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/TGRS.2014.2374218

In this paper, we tackle the manual annotation problem for object detection in optical RSIs by proposing a weakly supervised learning (WSL) framework. As one of the most cost-effective learning approaches, WSL only requires a weak label for the training images to specify whether the image contains the object of interest or not. To this end, unlike conventional supervised learning approaches, which rely on manually labeled bounding boxes for training object detector, accurate locations and sizes of the target objects are not needed in the WSL framework. Object detection using WSL tends to solve localization of the objects of interest in each positive training image (automatic annotation) and object detector training using automatic annotations (detector learning) simultaneously. In practice, WSL is implemented as follows. Given the weak label only indicating whether a certain category of object is contained in an image or not, an initial annotation is first obtained automatically, based on which, a detector is trained. The trained detector is then used as the annotator to refine the annotation, whereas the detector is iteratively trained using refined annotations until the model drift is detected. In this paper, we propose a Bayesian framework by jointly exploring saliency, intraclass compactness, and interclass separability to initialize a training example set. Afterward, we propose a novel detector evaluation method, which is able to cease the iterative learning process when the detector starts to drift to bad results, and thus, we can obtain final object detector with satisfactory performance.

To tackle the problem of insufficiently powerful feature descriptors, we explore the spatial and structural information within image patches via high-level feature learning. Unlike existing works to extract structural features solely based on human design [7], [8], the proposed approach derives high-level features by applying unsupervised representation learning approach, where spatial and structural patterns from the low-level and middle-level features can be automatically captured. Here, we adopt deep Boltzmann machine (DBM) to learn high-level feature because it has been demonstrated to have the potential of learning useful distributed feature representations and become a promising way in solving object and speech recognition problems [9]–[11].

In summary, the main contributions of this paper are threefold.

- 1) We propose a novel WSL framework based on Bayesian principles for detecting objects from optical RSIs, which extensively reduces human labors for annotating training data while achieving performance comparable with that of the fully supervised learning approaches.
- 2) We propose unsupervised feature learning via DBM to build high-level feature representation for various geospatial objects. The learned high-level features capture the structural and spatial patterns of objects in an effective and robust fashion, which leads to further improvement of object detection performance.
- 3) Extensive evaluations on three optical RSI data sets with different spatial resolutions and objects of interest are carried out to validate the effectiveness of the proposed methodology.

The rest of this paper is organized as follows. Section II gives a brief review of the related work. Section III introduces the proposed framework. Section IV proposes the unsupervised feature learning. Section V describes the WSL framework for object detection in optical RSIs. Experimental results are presented in Section VI. Finally, conclusions are drawn in Section VII.

II. RELATED WORK

Object or target detection in optical RSIs has been extensively studied in the past decades. For example, Li *et al.* [2] developed an algorithm for straight road edge detection from optical RSIs based on the ridgelet transform with the revised parallel-beam Radon transform. Liu *et al.* [4] detected inshore ships in optical satellite images by using shape and context information that is extracted in the segmented image. Liu *et al.* [3] presented robust automatic vehicle detection in QuickBird satellite images by applying morphological filters for road line removing and histogram representation for separating vehicle targets from background. All these methods are performed in an unsupervised manner. They are effective for detecting the designed object category in simple scenario.

With the advancement of machine learning techniques, many approaches started to cast object detection as a classification problem. In these approaches, a set of features that can characterize the objects is extracted first. Then, classification is performed using the extracted features and predefined classifiers. For example, Han *et al.* [1] proposed to detect multiple-class geospatial objects based on visual saliency modeling and discriminative learning of sparse coding. Cheng *et al.* [5] used histogram of oriented gradients (HOG) feature and latent support vector machine (SVM) to train deformable part-based mixture models for each object category. Based on the prior information obtained from a large number of human-labeled training examples, the supervised-learning-based approaches normally can achieve better performance. However, collection of large-scale training examples is often difficult and very time consuming.

A few efforts [12]–[15] have been performed to alleviate the work of human annotation. One interesting idea is to adopt the semisupervised learning model [16]. Such methods apply a self-learning or an active learning scheme where machine learning algorithms can automatically pick the most informative unlabeled examples based on a limited set of available labeled examples. Then, these picked unlabeled examples are combined with the initial labeled examples for the training of object detector or classifier. Specifically, Liao *et al.* [14] proposed a semisupervised local discriminant analysis method for feature extraction in hyperspectral RSI. Dópido *et al.* [13] adapted active learning methods to semisupervised learning for hyperspectral image classification. Jun and Ghosh [17] presented a semisupervised spatially adaptive mixture model to identify land covers from hyperspectral images.

Although semisupervised learning methods can considerably reduce the labor of human annotation, they still inevitably require a number of precise and concrete labeled training examples where each object is manually labeled by a bounding box

in positive training images. WSL is desirable to further reduce the human labor significantly, where the training set needs only binary labels indicating whether an image contains the object of interest. Although a few WSL approaches have been applied to natural scene image analysis [18]–[22], those existing methods cannot be directly used to the field of RSI analysis as they have insufficient capability to handle the challenges in RSIs, which contain large-scale complex background and a number of target objects with arbitrary orientation. As an initial effort, in our previous work in Zhang *et al.* [23], WSL was adopted and heuristically combined with saliency-based self-adaptive segmentation, negative mining algorithm, and negative evaluation mechanism for target detection in RSIs. This work lacks a principled framework and ignores some important information, which, thus, can be largely improved. In this paper, we propose a novel principled WSL framework for detecting targets from RSIs. Compared with [23], our improvements in this paper are threefold: 1) we propose powerful high-level feature learning using DBM; 2) we propose a probabilistic approach via the Bayesian rule to jointly integrate saliency, intraclass compactness, and interclass separability to initialize the training examples; and 3) we propose a novel scheme for model drift detection using the information from both negative training images and positive training images. The experimental results in Section VI can demonstrate these improvements.

III. OVERVIEW OF THE PROPOSED METHOD

Given a training optical RSI set with weak label only indicating whether a certain category of object is contained in an image or not, the objective of the proposed work is to detect target objects of the same class within the testing images. Because these images generally have very large scale and contain multiple objects of interest, a straightforward way of processing is to decompose the images into small patches by sliding windows and then predict whether each patch contains the object of interest. As suggested in [1] and [12], we adopt the multiscale sliding window scheme to handle the variation of the object size.

The proposed object detection framework consists of two major components: unsupervised feature learning and WSL-based object detection. The flowchart of the feature learning component (see Section IV) is shown in Fig. 1. In order to obtain more structural and semantic representation of the image patches, we extract the low-level and middle-level features to capture the spatial information and then use DBM to learn the hidden patterns of the middle-level features, which can abstract more structural and semantic information and lead to the desired high-level feature.

Based on the obtained high-level features, the component of WSL-based object detection shown in Fig. 2 (see Section V) contains two stages: training and testing. The objective of the training stage is to learn an object detector. In the testing stage, the learned object detector is applied to detect objects in a given testing image. The training stage includes two major steps: training example initialization (see Section V-A) and iterative object detector training (see Section V-B). For the first step, a Bayesian approach is proposed to integrate three kinds of

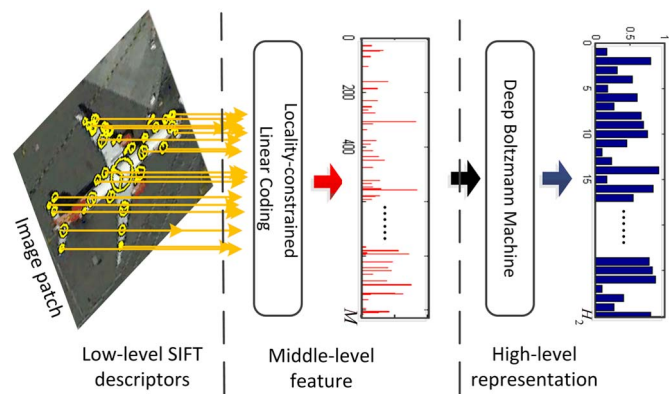


Fig. 1. High-level feature representation of the image patch.

important information of saliency, intraclass compactness, and interclass separability, which estimates the probability of an image patch being the object of interest. After initializing the training examples, we are inspired by the bootstrapping method [24] to train the object detector in an iterative process. In each iteration, the detector is utilized as an annotator to refine the positive training set, which is then used to retrain the detector. Thus, both the training examples and the object detector could be gradually updated to be more precise and strong. Afterward, a novel detector evaluation method is proposed to detect the model drift and stop the iterative process automatically for obtaining the final object detector.

IV. HIGH-LEVEL FEATURE REPRESENTATION

The performance of the existing feature descriptions in RSI analysis is still far from satisfactory. The main issue lies in the insufficiency in extracting features using only the pixel-based spectral information, which ignores the contextual spatial information and thus fails to capture the more important structural pattern of the object. With the advancement of the remote sensing technology, optical satellite and aerial imagery with high spatial resolution makes capturing spatial and structural information possible. Nowadays, accurate interpretation of optical RSI relies on effective spatial feature representation to capture the most structural and informative property of the regions in each image. A number of such approaches have started to explore the spatial information by applying some low-level descriptors (such as scale-invariant feature transform (SIFT), HOG, and gray-level co-occurrence matrices in [5], [25], and [26]) or middle-level features (such as bag of word (BOW) in [26]) to represent image patches. Although to some extent these human-designed features can improve the classification and detection accuracies in optical RSIs, they still suffer from several problems. Specifically, the low-level descriptors only catch limited local spatial geometric characteristics, which cannot be directly used to describe the structural contents of image patches. The middle-level features are usually extracted based on the statistic property of the low-level descriptors to capture the structural information of the spatial region. However, it cannot provide enough strong description and generalization ability for object detection in complex backgrounds.

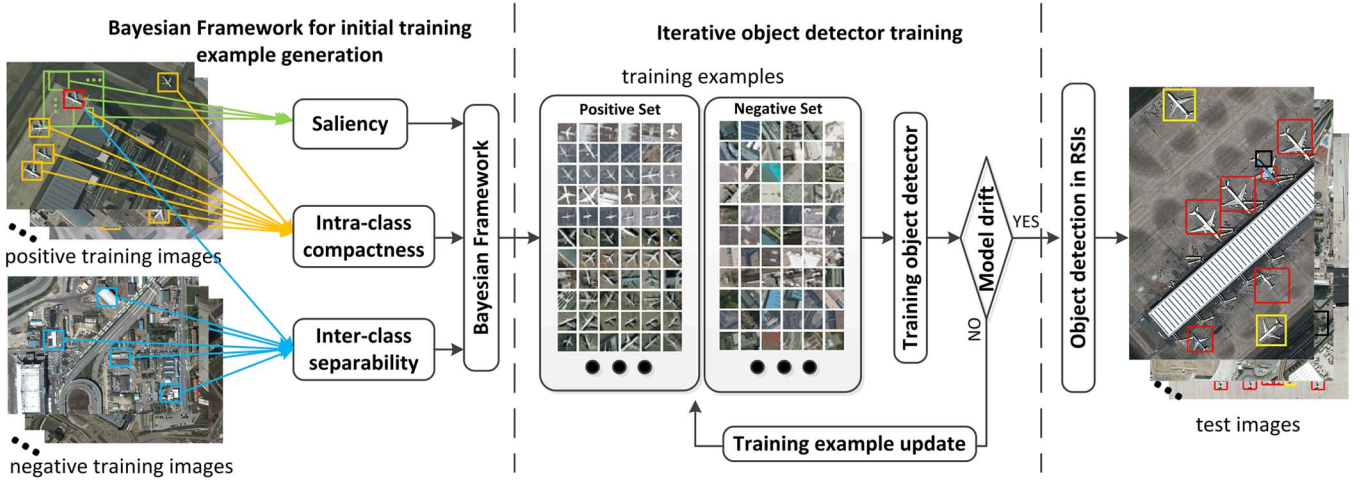


Fig. 2. Flowchart of WSL-based object detection.

To tackle these problems, we build high-level feature representation via DBM to capture the spatial and structural patterns encoded in the low-level and middle-level features. DBM is one type of neural networks with deep architecture that learns feature representation in an unsupervised manner and demonstrated to be promising for building high-level feature descriptors [9]–[11]. We therefore use it to map the middle-level features to the high-level representation that is highly accurate in characterizing different scenes or objects in optical RSIs. Specifically, the extraction of high-level feature representation (see Fig. 1) is carried out in three main stages, i.e., low-level descriptor extraction, middle-level feature generation, and high-level feature learning.

A. Low-Level Descriptor Extraction

We use low-level features to characterize the local region of each key point in image patches. Due to its ability to handle variations in terms of intensity, rotation, scale, and affine projection, the SIFT descriptor [27] is adopted in the proposed algorithm as the low-level descriptor to detect and describe the key points. According to existing work [12], [28], the SIFT descriptor has been demonstrated to outperform a set of existing descriptors, which are also widely used in analyzing RSIs.

B. Middle-Level Feature Generation

To alleviate the unrecoverable loss of discriminative information, we apply the locality-constrained linear coding (LLC) model [29] to encode the local descriptors into image patch representation. Specifically, all the extracted low-level descriptors are clustered to generate a codebook by using the K-means method. Let $\mathbf{D} = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_N]$ denote a set of N extracted low-level descriptors in one image patch. Given a codebook $\mathbf{CB} = [\mathbf{cb}_1, \mathbf{cb}_2, \dots, \mathbf{cb}_M]$ with M entries, LLC converts each descriptor into a M -dimensional code to generate the image patch representation by the following three steps. 1) For each input low-level descriptor \mathbf{d}_n , $n \in [1, N]$, its five nearest neighbors in \mathbf{CB} are used as the local bases \mathbf{LB}_n to

form a local coordinate system [29]. 2) The local code $\tilde{\mathbf{c}}_n$ is obtained by solving an objective function

$$\min \sum_{n=1}^N \|\mathbf{d}_n - \tilde{\mathbf{c}}_n \cdot \mathbf{LB}_n\|^2 \text{ s.t. } \sum_{n=1}^N \tilde{\mathbf{c}}_n = 1. \quad (1)$$

Then, the full code \mathbf{c}_n is generated, which is an $M \times 1$ vector with five nonzero elements whose values correspond to $\tilde{\mathbf{c}}_n$. 3) The final middle-level image patch representation is yielded by max pooling all the generated codes within the patch.

C. High-Level Feature Learning

A DBM [10] is a neural network with deep structure constructed by stacking multiple restricted Boltzmann machines (RBMs). In our framework, a three-layered DBM is adopted to learn high-level representations by capturing the structural and spatial patterns from middle-level features in an unsupervised manner. It contains a visible layer $\mathbf{v} \in \{0, 1\}^M$ and two hidden layers $\mathbf{h}^1 \in \{0, 1\}^{H_1}$ and $\mathbf{h}^2 \in \{0, 1\}^{H_2}$, where H_1 and H_2 indicate the numbers of units of the first and second hidden layers, respectively. The energy of the state $\{\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2\}$ is defined as

$$E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \Theta) = -\mathbf{v}^T \mathbf{W}^1 \mathbf{h}^1 - \mathbf{h}^{1T} \mathbf{W}^2 \mathbf{h}^2 \quad (2)$$

where $\Theta = \{\mathbf{W}^1, \mathbf{W}^2\}$ are the model parameters, denoting visible-to-hidden and hidden-to-hidden symmetric interaction terms. The probability that the model assigns to a visible vector \mathbf{v} is given by the Boltzmann distribution

$$\Pr(\mathbf{v}; \Theta) = \frac{1}{Z(\Theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2} \exp(-E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \Theta)) \quad (3)$$

where $Z(\Theta) = \sum_{\mathbf{v}} \sum_{\mathbf{h}^1, \mathbf{h}^2} \exp(-E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \Theta))$ is the partition function.

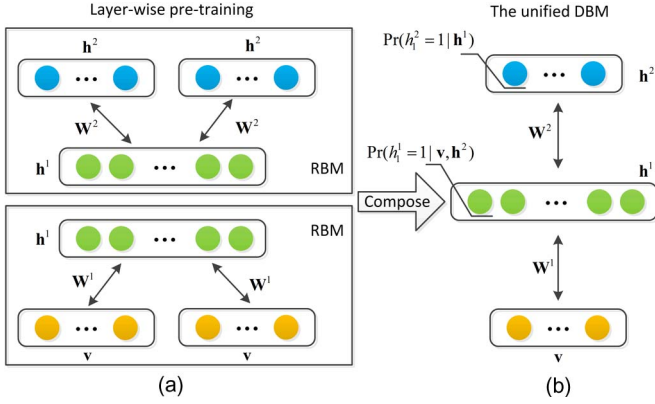


Fig. 3. Learning processes for DBM.

The conditional distributions over the visible units and the two sets of hidden units are given by

$$\Pr(h_i^1 = 1 | \mathbf{v}, \mathbf{h}^2) = \text{sigm} \left(\sum_{m=1}^M W_{mi}^1 v_m + \sum_{j=1}^{H_2} W_{ij}^2 h_m^2 \right) \quad (4)$$

$$\Pr(h_j^2 = 1 | \mathbf{h}^1) = \text{sigm} \left(\sum_{i=1}^{H_1} W_{mj}^2 h_i^1 \right) \quad (5)$$

$$\Pr(v_m = 1 | \mathbf{h}^1) = \text{sigm} \left(\sum_{i=1}^{H_1} W_{mi}^1 h_i^1 \right) \quad (6)$$

where $\text{sigm}(\cdot)$ is a sigmoid function.

Given a set of training data, learning of DBM is a process to determine the related model parameters $\Theta = \{\mathbf{W}^1, \mathbf{W}^2\}$ in (2). Although exact maximum-likelihood estimation of these parameters is intractable, efficient approximate learning of DBMs can be carried out by using mean-field inference together with the Markov chain Monte Carlo algorithms [10]. Furthermore, the entire model can be efficiently pretrained in a layer-by-layer unsupervised manner by minimizing the energy function in each individual RBM model [see Fig. 3(a)]. Composing the RBM models afterward forms a unified DBM model [see Fig. 3(b)], which can be used to extract high-level feature representation.

In the proposed algorithm, all the middle-level features extracted from the image patches in training images are used as the input data to train DBM, where the second hidden layer is used to build the final high-level feature representation for each image patch.

V. WSL-BASED OBJECT DETECTION

A. Training Example Initialization

By applying sliding windows as preprocessing, the training images are divided into many patches. Thus, the patch-level training data $X^+ = \{x_p^+ | p \in [1, P]\}$ and $X^- = \{x_q^- | q \in [1, Q]\}$ can be generated from the positive training images and the negative training images, respectively. Our first task is to select potential target object patches from X^+ to generate the initial positive training set X_0^+ . Typically, three different information

cues, namely, saliency, intraclass compactness, and interclass separability [20], [21], are used to initialize the positive training examples. Based on the assumption that the object to be detected is one kind of foreground objects in the image, saliency information ensures that the selected positive example is a foreground region. It acquires generic knowledge about the sizes and locations of the objects. The intraclass compactness enforces the selected positive examples to be visually similar to each other, whereas the interclass separability ensures that all selected positive examples are different from negative examples. In this paper, a novel Bayesian framework is proposed to combine these three types of information simultaneously to initialize the positive example training set as follows.

Let y_p^+ denote whether an image patch x_p^+ belongs to one specified object. According to Bayes' rule

$$\Pr(y_p^+ = 1 | x_p^+) = \frac{\Pr(x_p^+ | y_p^+ = 1) \Pr(y_p^+ = 1)}{\Pr(x_p^+)} \quad (7)$$

$$\begin{aligned} \Pr(y_p^+ = 1 | x_p^+) &= 1 - \Pr(y_p^+ = 0 | x_p^+) \\ &= 1 - \frac{\Pr(x_p^+ | y_p^+ = 0) \Pr(y_p^+ = 0)}{\Pr(x_p^+)}. \end{aligned} \quad (8)$$

After adding the preceding two equations and omitting the constant term, we have

$$\begin{aligned} \Pr(y_p^+ = 1 | x_p^+) &\propto \underbrace{\frac{1}{\Pr(x_p^+)}}_{\text{Saliency}} \underbrace{\Pr(x_p^+ | y_p^+ = 1)}_{\text{Intra-class compactness}} \\ &\times \underbrace{\Pr(y_p^+ = 1)}_{\text{Prior Probability}} - \underbrace{\Pr(x_p^+ | y_p^+ = 0)}_{\text{Inter-class separability}} \underbrace{\Pr(y_p^+ = 0)}_{\text{Prior Probability}}. \end{aligned} \quad (9)$$

In the information theory, $-\log \Pr(x_p^+)$, which is the log form of $1/\Pr(x_p^+)$, is known as the self-information of the random variable x_p^+ [30], [31]. Self-information increases when the probability of a patch decreases. In other words, patches discriminative from surroundings are more informative and thus more likely to be objects. Therefore, the term of $1/\Pr(x_p^+)$ in (9) is associated with the saliency information. The term $\Pr(x_p^+ | y_p^+ = 1)$ indicates the likelihood that favors image patches sharing the similar characteristic with the class of target object. Hence, it can be considered as the metric of intraclass compactness. Similarly, $\Pr(x_p^+ | y_p^+ = 0)$ reflects the distinctness of image patches in positive and negative images; thus, it corresponds to the metric of interclass separability. Finally, the remaining two prior probabilities $\Pr(y_p^+ = 1)$ and $\Pr(y_p^+ = 0)$ are treated as the weights of the intraclass and interclass metrics, respectively.

1) *Saliency*: As we assume that objects to be detected are normally one kind of foreground objects, our objective then becomes quantifying how likely each image patch is a foreground object. Foreground objects are generally informative and salient from the surrounding background, as shown in Fig. 4. In computer vision, saliency detection technique can be used to estimate the saliency for each image patch. In recent years, it is also employed for the analysis in the domain of

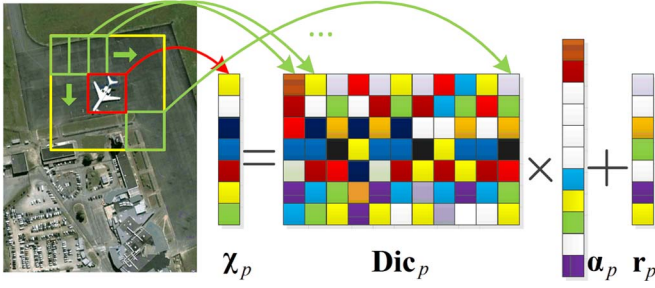


Fig. 4. Illustration of saliency calculation.

remote sensing [32], [33]. Inspired by [32], we adopt sparse coding theory to calculate saliency based on the raw pixels to reveal the structural difference between an image patch and its surrounding. For each image patch x_p^+ (the patch indicated by red frame in Fig. 4), it is sparsely coded with its adjacent half-overlapped surrounding patches (patches indicated by green frames in Fig. 4) by

$$\chi_p \approx \text{Dic}_p \alpha_p \quad (10)$$

where χ_p are the raw pixels within x_p^+ , whereas Dic_p and α_p indicate the dictionary constructed by all surrounding patches and the sparse codes, respectively.

The rationale behind (10) is to represent χ_p approximately by its surrounding patches. According to [32], the coding sparseness $\|\alpha_p\|_0$ and the coding residual $\mathbf{r}_p = \chi_p - \text{Dic}_p \alpha_p$ indicates the saliency of the image patch x_p^+ with respect to its surrounding. Therefore, we estimate the saliency by

$$1/\Pr(x_p^+) = \|\alpha_p\|_0 \cdot \|\mathbf{r}_p\|_1. \quad (11)$$

2) *Intraclass Compactness*: Termed as $\Pr(x_p^+|y_p^+ = 1)$, the intraclass compactness metric aims to constrain the similarity between positive examples. As positive examples of a specific object category should be visually similar, we use the Gaussian mixture model (GMM) to estimate the probability distribution of all positive examples. Then, $\Pr(x_p^+|y_p^+ = 1)$ measures how likely each image patch is a positive example. Image patches with large $\Pr(x_p^+|y_p^+ = 1)$ may be selected as positive examples. We use the high-level feature \mathbf{f}_p^+ to represent each image patch x_p^+ as this feature can handle the variations in scale and orientation and capture the spatial and structural patterns of each image patch. As patterns learned by DBM are approximately independent, the joint probability is simplified to the product of probability of each hidden unit's response, i.e.,

$$\Pr(x_p^+|y_p^+ = 1) = \prod_{j=1}^{H_2} \Pr([\mathbf{f}_p^+]_j | y_p^+ = 1) \quad (12)$$

where $[\mathbf{f}_p^+]_j$ indicates the j th-dimensional value of \mathbf{f}_p^+ , and H_2 indicates the dimensionality of \mathbf{f}_p^+ . The distribution of each hidden unit's response is estimated using GMM with adaptive component $K_{j=1, \dots, H_2}^+$ by

$$\Pr([\mathbf{f}_p^+]_j | y_p^+ = 1) = \sum_{k=1}^{K_j^+} \pi_{jk}^+ N([\mathbf{f}_p^+]_j | \mu_{jk}^+, \sigma_{jk}^{2+}) \quad (13)$$

where π_{jk}^+ , μ_{jk}^+ , and σ_{jk}^{2+} are parameters of the GMM in the k th component for the j th-dimensional feature. All parameters are inferred based on object candidates in \tilde{X}^+ by using the expectation-maximization algorithm and Bayesian inference [34]. Here, \tilde{X}^+ denotes the set of object candidates and will be described in Section V-A5.

3) *Interclass Separability*: The interclass separability metric is to enforce that the selected positive examples are dissimilar to negative examples. In WSL, the most confident information comes from the negative training images because they definitely do not contain the target. It is also reasonable to believe that the positive examples containing target objects should be different from the negative image patches in the negative images. Consequently, we can collect a large number of negative image patches to estimate the probability distribution of negative examples via a GMM. Then, we formulate the interclass metric as the likelihood term $\Pr(x_p^+|y_p^+ = 0)$, which reflects the probability of a certain image patch appearing in negative training images. The high probability of the appearance in negative images would lead to low interclass difference and separability. Similar to $\Pr(x_p^+|y_p^+ = 1)$, $\Pr(x_p^+|y_p^+ = 0)$ can be decided based on the high-level feature by

$$\Pr(x_p^+|y_p^+ = 0) = \prod_{j=1}^{H_2} \Pr([\mathbf{f}_p^+]_j | y_p^+ = 0) \quad (14)$$

$$\Pr([\mathbf{f}_p^+]_j | y_p^+ = 0) = \sum_{k=1}^{K_j^-} \pi_{jk}^- N([\mathbf{f}_p^+]_j | \mu_{jk}^-, \sigma_{jk}^{2-}) \quad (15)$$

where parameters π_{jk}^- , μ_{jk}^- , σ_{jk}^{2-} , and K_j^- are inferred by GMM based on all the negative image patches.

4) *Prior Probability*: $\Pr(y_p^+ = 1)$ and $\Pr(y_p^+ = 0)$ are two prior terms in the proposed Bayesian framework. According to [34], Bayesian methods would result in poor performance when inappropriate choices of prior are applied without any prior belief. Therefore, inspired by [35], we define the prior terms to reflect the prior belief. Our prior belief is that $\Pr(y_p^+ = 0)$ should be high when the content of certain image patch x_p^+ has small distance to the negative image patches in \tilde{X}^- , and $\Pr(y_p^+ = 1)$ should become high when the content of x_p^+ is close to the object candidates in \tilde{X}^+ . Hence, we simply adopt the nearest neighbor distance [20] to estimate these prior probabilities as

$$\Pr(y_p^+ = 0) = \exp\{-\|x_p^+ - \text{Nn}(x_p^+)\|_1\} \quad (16)$$

$$\Pr(y_p^+ = 1) = \exp\{-\|x_p^+ - \text{Np}(x_p^+)\|_1\} \quad (17)$$

where $\|\cdot\|_1$ is the L_1 norm. Same as in [20], $\text{Nn}(x_p^+)$ and $\text{Np}(x_p^+)$ refer to the nearest neighbors of x_p^+ in \tilde{X}^- and \tilde{X}^+ (in terms of the high-level feature), respectively. Finally, these two prior terms are used as the weights of the interclass and intraclass metrics in order to reflect the prior probability that an image patch belongs to the positive and negative training examples, respectively.

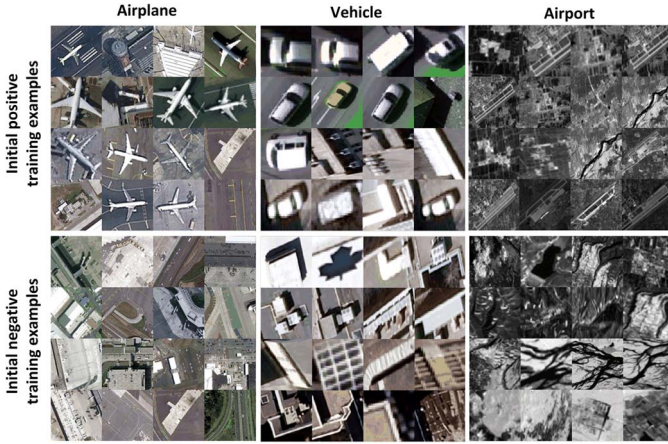


Fig. 5. Some examples in initial positive and negative training sets.

5) *Implementation Details*: In terms of (9), the postprobability $\Pr(y_p^+ = 1|x_p^+)$ is estimated by integrating the saliency, intraclass, and interclass metrics. Note that, before calculating the intraclass compactness metric, \tilde{X}^+ needs to be available. The work [21] proposed an exhaustive searching strategy to generate one object candidate for each image. However, it lacks accuracy and efficiency for the large-scale RSIs, particularly when it contains multiple target objects located at quite scattered positions. To tackle this challenging problem, we in practice implement our work in two stages.

In the first stage, we calculate $\Pr(y_p^+ = 1|x_p^+)$ approximately by only using the saliency and interclass separability metrics to generate \tilde{X}^+ . As initially $\Pr(x_p^+|y_p^+ = 1)$ and $\Pr(y_p^+ = 0)$ are unknown, we omit them by following [30] and [31], which is equivalent to assuming a uniform likelihood distribution for the unspecified object category. The overall formulation reduces to

$$\Pr(y_p^+ = 1|x_p^+) \propto \frac{1}{\Pr(x_p^+)} [1 - \Pr(x_p^+|y_p^+ = 0)\Pr(y_p^+ = 0)]. \quad (18)$$

Hence, \tilde{X}^+ can be further determined by choosing a probability threshold τ , i.e.,

$$\tilde{X}^+ = \{x_p^+ | \Pr(y_p^+ = 1|x_p^+) \geq \tau\}. \quad (19)$$

Once \tilde{X}^+ is obtained, we fully implement the proposed Bayesian framework in the second stage, where all the three types of information are explored and integrated for calculating $\Pr(y_p^+ = 1|x_p^+)$ by (9). Similar to the first stage, a threshold τ is chosen to generate the initial positive training set X_0^+ by

$$X_0^+ = \{x_p^+ | \Pr(y_p^+ = 1|x_p^+) \geq \tau\}. \quad (20)$$

By considering the fact that imbalanced positive and negative training data may reduce the performance of the object detector, we follow the previous work of [24] to generate the initial negative training set X_0^- by random undersampling of X^- to the same size as X_0^+ . Some examples in the initial training set are shown in Fig. 5.

B. Iterative Detector Training

After obtaining the initial training examples, the object detector is iteratively trained in the proposed framework (see Fig. 2). In each iteration, the training set generated by the previous iteration is used to train the current object detector, which, in turn, updates the training examples for the next iteration. The iteration process stops when a model drift is detected by a novel detector evaluation method. Then, the object detector obtained before model drift is regarded as the final object detector.

1) *Training Example and Object Detector Update*: As shown in Fig. 5, although most of the examples in the initial positive training set generated by the proposed work are the objects of interest, it still contains several noise examples. Consequently, promising object detector cannot be obtained by directly using the initial training data. Inspired by [19] and [21], we train the object detector in an iterative process, which can update the training set and the object detector iteratively. Linear SVM is adopted in the proposed algorithm because it has very low training costs and has been demonstrated to be both efficient and effective in RSI analysis [1], [36]. Based on the SVM formulation, we use the following score function for object annotation and detection:

$$\text{Score}(x_p^+) = \mathbf{w}_1^T \mathbf{f}_p^+ + b_1 \quad (21)$$

where the variables \mathbf{w}_1 and b_1 are defined as the SVM decision plane and its bias, respectively, which are learnt from the initial training data. With the score function, binary class label y_p^+ is assigned to the image patch x_p^+ based on the sign of the function, i.e.,

$$y_p^+ = \begin{cases} 1, & \text{Score}(x_p^+) \geq 0 \\ 0, & \text{Score}(x_p^+) < 0. \end{cases} \quad (22)$$

In order to obtain more precise object patches as the updated positive training examples, an adaptive threshold is used to determine image patches that have higher confidence to be the object of interest, i.e.,

$$X_1^+ = \left\{ x_p^+ | \text{Score}(x_p^+) \geq \frac{\sum_{p=1}^P y_p^+ \cdot \text{Score}(x_p^+)}{\sum_{p=1}^P y_p^+} \right\} \quad (23)$$

where X_1^+ is the updated positive training set after the first iteration. Afterward, the same number of negative examples randomly selected from X^- is used to generate the new negative training set X_1^- . Alternating the update of object detector and training examples progressively improves their accuracy until the end of the iteration. Combination of these two stages in an iterative way is very similar to the bootstrapping [24] or the active learning [12] strategy, which allows the proposed WSL-based object detection in optical RSIs to achieve good performance that is even superior to the traditional supervised learning methods in some cases.

2) *Detector Evaluation*: Similar to the model drift phenomenon in adaptive object tracking, the performance of the trained object detector is improved in the first several iterations, continually, and then begins to degrade. Consequently, generating reasonable evaluation mechanism to detect the model drift

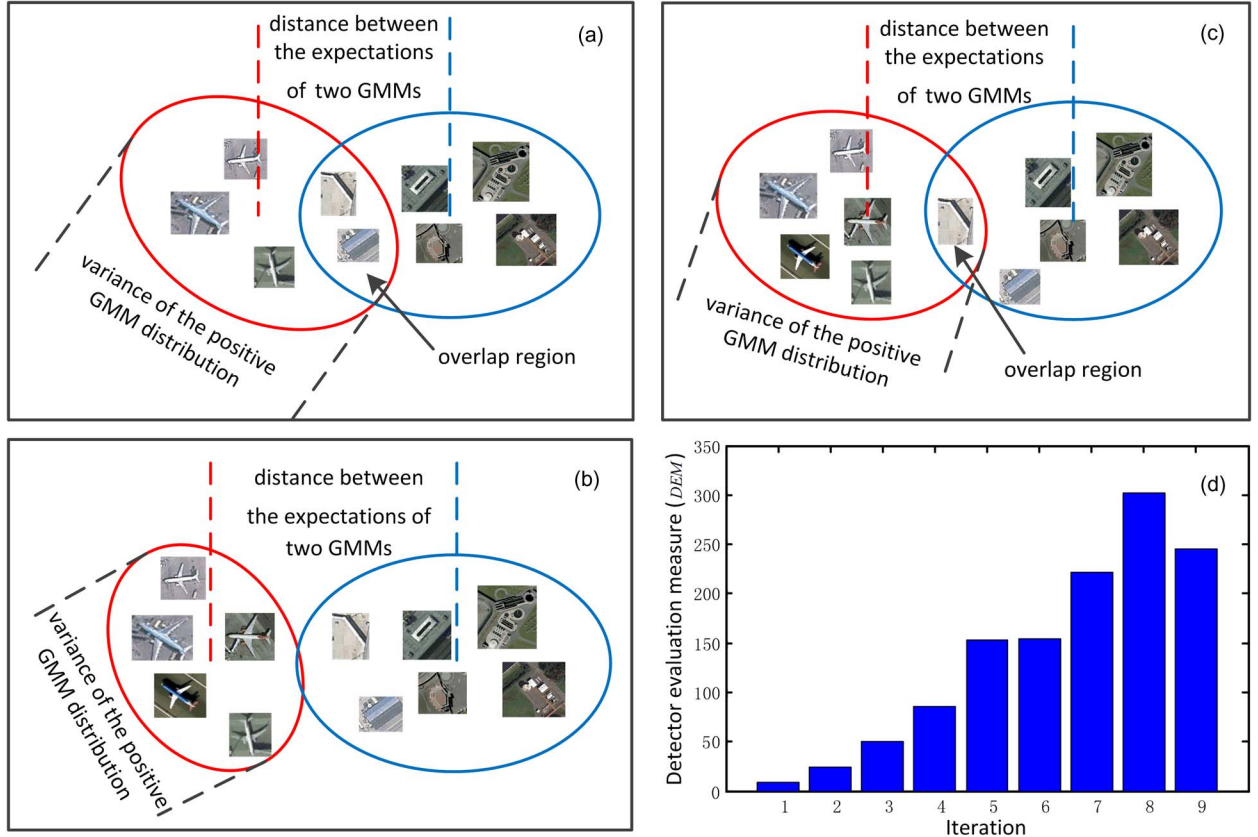


Fig. 6. Simple illustration of model drift on GMM distribution. In (a)–(c), the distribution of positive GMM is in red color, whereas the distribution of negative GMM is in blue color. In (d), how the value changes in each iteration is shown. It is based on the iterative training of airport detector. See text for detailed explanation.

is important. As the exact location of the objects of interest in each positive training image is unknown, thus, it is impossible to measure directly whether a stronger object detector has been obtained after each iteration. It brings great challenge for evaluating the object model and detecting model drift in WSL.

First, we use a negative-example-based evaluation mechanism to estimate the performance of the trained object detector in each iteration. In general, a good object detector is expected to obtain detection results with high true positives and low false positives. In the WSL, we can only obtain precise negative image patches, which certainly contain no object of interest. As a result, the negative evaluation mechanism is adopted here to approximately evaluate the false positive rate for the object detector. Specifically, for each iteration, the trained object detector is applied to classify image patches with the negative training images and then calculate the false positive rate FR by

$$FR = |X_{false}^-| / |X^-| \quad (24)$$

$$X_{false}^- = \{x_q^- | \text{Score}(x_q^-) \geq 0\} \quad (25)$$

where $|\cdot|$ denotes the number of elements.

Another evaluation mechanism is based on the estimation of the object detector's performance in positive training images. Here, we define GMM^+ and GMM^- as the distributions inferred by GMM based on the positive and negative examples, respectively. As the positive training examples are updated after

each iteration, the distribution of the j th-dimensional high-level feature is modified along the iteration as

$$GMM_j^+ = \sum_{k=1}^{\tilde{K}_j^+} \tilde{\pi}_{jk}^+ N(\tilde{\mu}_{jk}^+, \tilde{\sigma}_{jk}^{2+}) \quad (26)$$

where $\tilde{\pi}_{jk}^+$, $\tilde{\mu}_{jk}^+$, $\tilde{\sigma}_{jk}^{2+}$, and \tilde{K}_j^+ are inferred based on the updated positive training examples after each iteration. In contrast, GMM^- is fixed as

$$GMM_j^- = \sum_{k=1}^{K_j^-} \pi_{jk}^- N(\mu_{jk}^-, \sigma_{jk}^{2-}) \quad (27)$$

where π_{jk}^- , μ_{jk}^- , σ_{jk}^{2-} , and K_j^- are inferred based on the constant negative image patches in X^- . In the first iteration, the object detector trained on the initial training examples is not very accurate. Thus, the trained object detector may work unsatisfyingly, and the updated positive examples generate the GMM^+ distribution having amount of overlap with the GMM^- distribution, as shown in Fig. 6(a). After several iterations, if the object detector is getting stronger, the overlap between the two distributions should become less, as shown in Fig. 6(b). Finally, when the detector starts to drift toward some noise patches without containing objects of interest, the overlap tends to become large again, as shown in Fig. 6(c). Consequently, we evaluate the object detector and monitor the model drift by estimating

the overlap between the two GMM distributions in each iteration. As the GMM^- distribution is fixed, the distance between the expectations of the two distributions and the variance of the GMM^+ distribution are used to approximately predict the overlap. Intuitively, the GMM^+ distribution with expectation away from that of GMM^- and small variance has small overlap with the distribution of GMM^- and vice versa (see Fig. 6). According to [34], the expectation and variance of GMM^+ for the j th-dimensional high-level feature are decided by

$$\text{Ex}(\text{GMM}_j^+) = \sum_{k=1}^{\tilde{K}_j^+} \tilde{\pi}_{jk}^+ \tilde{\mu}_{jk}^+ \quad (28)$$

$$\text{Var}(\text{GMM}_j^+) = \sum_{k=1}^{\tilde{K}_j^+} \tilde{\pi}_{jk}^+ (\tilde{\sigma}_{jk}^{2+} + \tilde{\mu}_{jk}^{2+}) - \text{Ex}(\text{GMM}_j^+)^2 \quad (29)$$

where $\tilde{\pi}_{jk}^+$, $\tilde{\mu}_{jk}^+$, $\tilde{\sigma}_{jk}^{2+}$, and \tilde{K}_j^+ are inferred by the updated positive training examples after each iteration. Similarly, the expectation of GMM^- for the j th-dimensional high-level feature is obtained by

$$\text{Ex}(\text{GMM}_j^-) = \sum_{k=1}^{K_j^-} \pi_{jk}^- \mu_{jk}^- \quad (30)$$

By combining the aforementioned two evaluation mechanisms, the final detector evaluation measure (DEM) in WSL is determined as

$$\text{DEM} = \frac{\sqrt{\sum_{j=1}^{H_2} (\text{Ex}(\text{GMM}_j^+) - \text{Ex}(\text{GMM}_j^-))^2}}{\text{FR} \times \sum_{j=1}^{H_2} \text{Var}(\text{GMM}_j^+)}. \quad (31)$$

Based on DEM, we can evaluate the object detector trained in each iteration. Higher DEM indicates better performance of the current object detector and vice versa. Being consistent with the preceding analysis, the DEM value of the object detector trained in the first iteration should be relatively small. Then, it increases as the detector is gradually refined in the following iterations. When the DEM value starts to decrease, the model drift is detected, and the iteration process is terminated. The final object detector is determined as the one obtained before the model drift [see Fig. 6(d)].

VI. EXPERIMENTS

A. Data Sets and Experimental Setup

Three optical RSI data sets established in [23] with different spatial resolutions and various objects of interest were used in our experiments. The details of these data sets are shown in Table I. The first data set consists of 120 very high resolution images from the publicly available Google Earth service. This data set is adopted to train and test the airplane detector. Seventy randomly selected images were weakly labeled and used as the training set (50 images containing airplanes as positive training images and 20 images not containing any airplanes as negative training images), and the remaining 50 images were used as the testing images. The second data set, which

TABLE I
INFORMATION ABOUT THE THREE EVALUATION DATABASES

Data Set	Dimension (pixels)	Spatial Resolution	Target Area (pixels)
Google Earth	about 1000×800	About 0.5m	700~25488
ISPRS	about 900×700	8-15cm	1150~11976
Landsat	400×400	30m	1760~15570

is called the ISPRS data set, is a very high resolution aerial image data set, which contains 100 images of vehicle objects provided by the German Association of Photogrammetry and Remote Sensing (DGPF) [37]. We randomly selected 60 weakly labeled images as the training data (45 positive training images and 15 negative training images) to train vehicle detector. The remaining 40 images were used as the testing data. The third data set consists of 180 shortwave-infrared imageries from Landsat-7 satellite. One hundred thirty-three randomly selected images were weakly labeled (98 positive training images and 35 negative training images) and used as the training data to train the airport detector. The remaining 47 images were used as the testing data. For all the three data sets, we also manually labeled bounding box for each target object in both training data and testing data to form the ground truth for the following evaluations.

In the experiments, as suggested in [1], we used square sliding windows with side lengths of {60, 100, 135} for airplane detection, {60, 80} for vehicle detection, and {60, 100, 130} for airport detection, respectively, where the sliding step size was also set to be 1/3 of the window side length. When building the high-level feature for the image patches, we set the number of entries $M = 1024$ empirically.

In the test phase, the proposed object detector trained using our WSL framework was performed to classify each image patch in the test images generated by the multiscale sliding window scheme. For sliding windows in different sizes, there may be significant overlap on detected targets. To solve this problem, we adopt a nonmaximum suppression step, as suggested in [1], [5], and [12], to retain the sliding window with the highest score.

B. Key Parameter Analysis

In the implementation of training example initialization (see Section V-A), several parameters may affect the performance and thus have to be set properly. These include the number of units H_1, H_2 in each hidden layer of DBM and the probability threshold τ in (19) and (20). To show how their values affect the performance of the proposed approach, we performed experiments on all the three data sets and evaluated the F1-measure by

$$\text{F1 - measure} = \frac{2 \cdot \text{PRE} \cdot \text{REC}}{\text{PRE} + \text{REC}} \quad (32)$$

$$\text{PRE} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{REC} = \frac{\text{TP}}{\text{NP}} \quad (33)$$

where TP, FP, and NP denote the number of true positives, the number of false positives, and the number of total positives under the threshold τ , respectively. As suggested in [1] and [5],

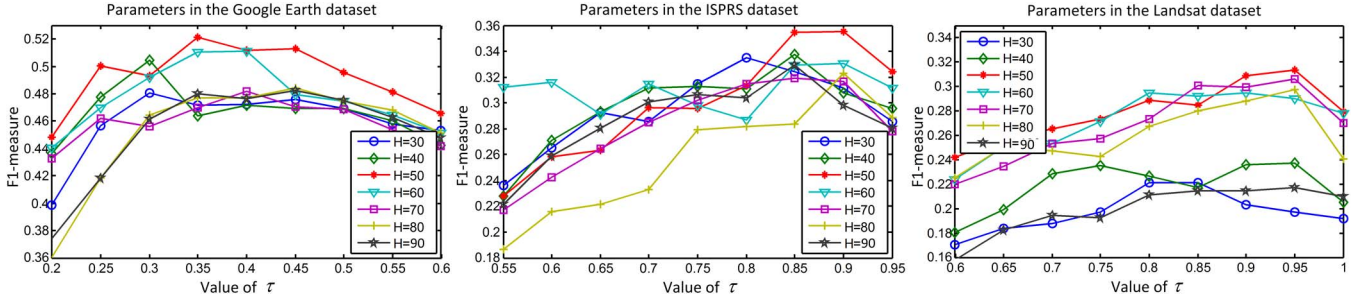


Fig. 7. Influence of key parameters to training example initialization.

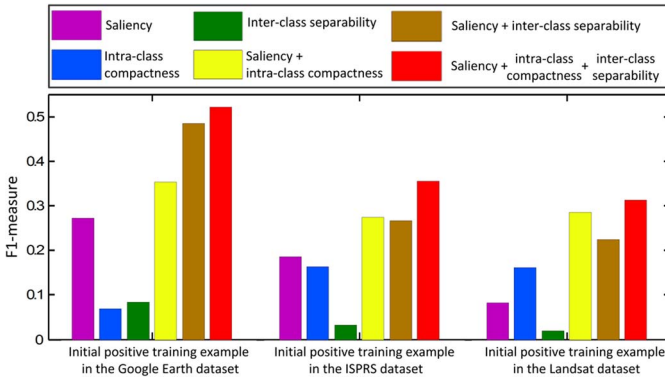


Fig. 8. Evaluation of the proposed Bayesian framework.

an annotation or a detection is marked as a true positive when its corresponding image patch can cover more than 50% of a ground truth. PRE and REC denote the precision and recall rate, respectively. As suggested in [38], equal number of units is used in each hidden layer ($H_1 = H_2 = H$) in our implementation, and the experimental results are shown in Fig. 7. We empirically set $H = 50$ for all the data sets and $\tau = 0.45, 0.90$, and 0.95 for the Google Earth data set, the ISPRS data set, and the Landsat data set, respectively, based on which the best detection performance can be achieved. We used this set of parameters in subsequent experiments.

C. Evaluation of the Bayesian Framework

Here, we evaluated the performance of the proposed Bayesian framework by comparing it with the baseline methods. Since the proposed Bayesian framework integrates the saliency, intraclass compactness, and interclass separability information for the positive training example initialization (indicated by the bins in red in Fig. 8), we evaluated its performance on the training sets. Here, we treat the methods that initialize positive examples by using the saliency information only, the interclass information only, the intraclass information only, fusing the saliency and interclass information, and fusing the saliency and intraclass information as the baseline methods. Note that the last two baseline methods were also implemented by using the proposed Bayesian framework. Based on the criterion of F1-measure, the experimental results are shown in Fig. 8.

From Fig. 8, we can observe the following: 1) the impact of the three single information on the initialization results changes with the variation of the data set and object of interest. For example, saliency makes the biggest contribution for the initialization results on the Google Earth data set, whereas the intraclass information contributes mostly on the Landsat data set; 2) in comparison with those using one of the three kinds of information, the performance of the fusion methods is better; and 3) the fusion of all the three information always achieves the best results regardless of the variation of data sets and objects of interest.

D. Evaluation of the High-Level Feature

In order to demonstrate the effectiveness of the proposed high-level feature, we compared it with three state-of-the-art features, which include the BOW [26], the pyramid HOG (pHOG) [5], [39], and the LLC [29]. Specifically, the BOW feature characterizes each training data by using a histogram of visual words; the pHOG feature represents the shape property of the image patches by using HOG, whereas the LLC feature is described in Section IV-B. For quantitative evaluation, we plot the precision–recall (PR) curve of the object detection results and the calculated average precision (AP) value, as shown in Fig. 10, for comparisons. Specifically, the PR curve is plotted based on the values of PRE and REC under different thresholds, whereas the AP is calculated by the area under the PR curve [1], [12]. The four different features were compared using the proposed WSL framework and the same sets of training and testing data. As shown in Fig. 9, the proposed high-level feature always outperforms the other three state-of-the-art features.

E. Evaluation of the Object Detector

We evaluated the performance of the proposed weakly supervised object detector by comparing it with one existing WSL-based method and several supervised-learning-based methods. First, we compared the proposed approach with the WSL-based method in [23]. For the fair comparison, in the experiment, we utilized the same experimental settings, including the same feature representation built by DBM, the same sliding window scheme, and the same testing image set. Fig. 10 gives the PR curves of the experiment results. The corresponding AP values are shown in Table II.

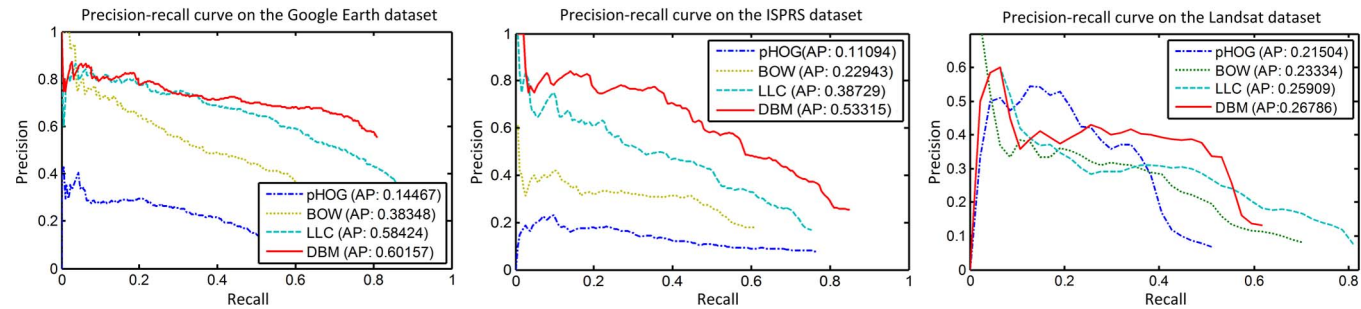


Fig. 9. PR curves for different types of feature on three data sets. Here, DBM indicates the high-level feature learned by the proposed work.

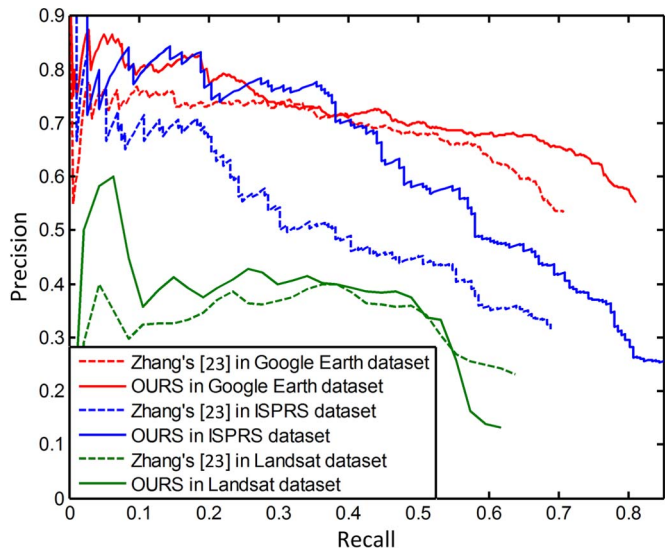


Fig. 10. PR curves for the comparisons with the WSL method.

TABLE II
DETAILED TARGET DETECTION RESULTS
IN TERMS OF THE METRIC OF AP

Objects of interest	WSL based object detector		Supervised learning based object detector		
	OURS	Zhang's	Baseline	Xu's	Han's
Airplane	0.6016	0.5128	0.6194	0.5275	0.5421
Vehicle	0.5332	0.3833	0.5774	0.4581	0.4626
Airport	0.2679	0.2439	0.2884	0.2710	0.3257
Overall	0.4676	0.3801	0.4951	0.4189	0.4435

We also compared the proposed WSL approach with several existing supervised-learning-based object detection methods, including a baseline method from Xu *et al.* [26] and the method of Han *et al.* [1]. The baseline method was implemented by training object detector (linear SVM) based on the proposed high-level feature in a manner of fully supervised learning, where the human annotations (manually labeled bounding box for each target in training images) are provided in the training images. The object detector trained by the method of Xu *et al.* was based on the spectral and texture local feature descriptor and SVM with radial basis function kernel. The method of Han *et al.* trained object detector via discriminative sparse coding, which has small within-class scatter and large between-class scatter. All comparison methods were evaluated using the same sets of training and testing data. Fig. 11 illustrates the PR

curves of the experiment results. The corresponding AP values are shown in Table II.

From Figs. 10 and 11 and Table II, we can observe that the proposed WSL approach can achieve much better performance than the state-of-the-art WSL-based method and comparable performance with the state-of-the-art fully supervised learning based methods. Specifically, the object detection accuracy of the proposed WSL approach achieves about 97.13%, 92.34%, and 92.89% of what the baseline approach does in the Google Earth data set, the ISPRS data set, and the Landsat data set, respectively. It also improves the performance of the previous WSL-based approach [23] significantly, i.e., 0.088 (8.88%), 0.1499 (14.99%), and 0.024 (2.4%) in terms of AP in the Google Earth data set, the ISPRS data set, and the Landsat data set, respectively. More encouragingly, the proposed WSL approach performs even better than the other two state-of-the-art fully supervised methods in some cases. Specifically, for airplane detection in the Google earth data set, it outperforms the methods of Xu *et al.* and Han *et al.* by 0.0741 (7.41%) and 0.0595 (5.95%), respectively. In the ISPRS data set, the proposed WSL approach outperforms the methods of Xu *et al.* and Han *et al.* by 0.0751 (7.51%) and 0.0706 (7.06%), respectively.

From the overall results among the three data sets, it can be seen that, due to the powerful high-level feature representation built by DBM, the supervised baseline method yields the best results on these data sets. Benefited by the Bayesian framework to generate accurate initial training examples and the iterative training scheme to gradually refine the object detector, the proposed WSL algorithm achieves detection performance that outperforms the previous WSL-based target detection method [23] and approaches to the fully supervised baseline method. Furthermore, based on the combination of the high-level feature representation and the proposed WSL framework, the overall performance of weakly supervised detector apparently outperforms the other two existing state-of-the-art supervised methods.

Finally, some experimental results from the proposed approach for airplane, vehicle, and airport detection are shown in Fig. 12, respectively. In these figures, the red rectangles indicate the true positive results; whereas the black and yellow rectangles denote the false positive and miss alarm results, respectively. As can be seen, the object detector trained via the proposed WSL approach can effectively detect objects of interest from all the data sets with different spatial resolutions and cluttered backgrounds.

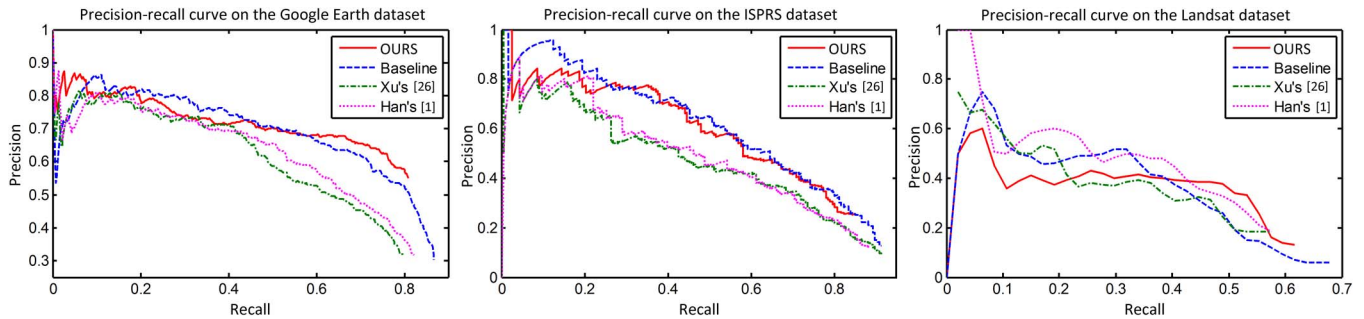


Fig. 11. PR curves for the comparisons with the supervised learning methods.

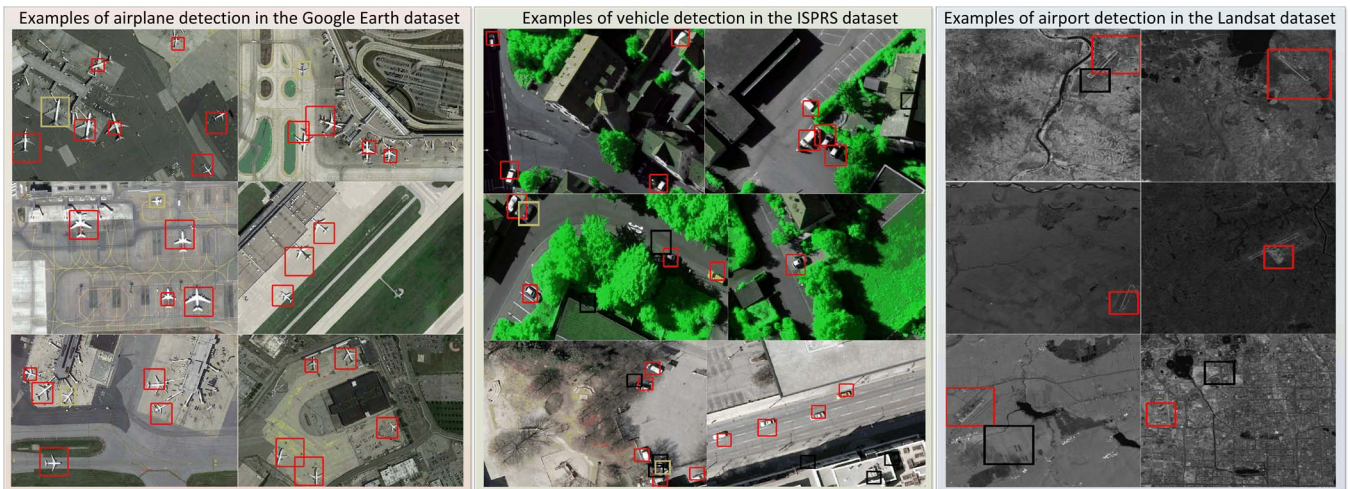


Fig. 12. Some samples from the three benchmark data sets.

VII. CONCLUSION

In this paper, we have proposed a novel framework to tackle the problem of object detection in optical RSIs. The novelties that distinguish the proposed work from previous works lie in two major aspects. First, instead of using traditional supervised or semisupervised learning methodology, this paper developed a WSL framework that can substantially reduce the human labor of annotating training data while achieving the outstanding performance. Second, we developed a deep network to learn high-level features in an unsupervised manner, which offers a more powerful descriptor to capture the structural information of objects in RSIs. It thus can improve the object detection performance further. Experiments on three different RSI data sets have demonstrated the effectiveness of the proposed work.

Our future work will focus on three directions. First, we will extend it to the joint learning of multiple categories of object detectors. Second, we will combine the rich spectral information provided by RSIs with spatial information for robust object detection. Third, transfer learning [40] will be used to further improve the WSL framework.

ACKNOWLEDGMENT

The authors would like to thank the provision of the Vaihingen data set provided by the German Society for Photogrammetry, Remote Sensing and Geoinformation and the provision of the Downtown Toronto data set by Optech Inc., First Base Solutions Inc., GeoICT Lab, and ISPRS WG III/4.

REFERENCES

- [1] J. Han *et al.*, "Efficient, simultaneous detection of multi-class geospatial targets based on visual saliency modeling and discriminative learning of sparse coding," *ISPRS J. Photogramm. Remote Sens.*, vol. 89, pp. 37–48, Mar. 2014.
- [2] X. Li, S. Zhang, X. Pan, P. Dale, and R. Cropp, "Straight road edge detection from high-resolution remote sensing images based on the ridgelet transform with the revised parallel-beam Radon transform," *Int. J. Remote Sens.*, vol. 31, no. 19, pp. 5041–5059, Oct. 2010.
- [3] W. Liu, F. Yamazaki, and T. T. Vu, "Automated vehicle extraction and speed determination from QuickBird satellite images," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 4, no. 1, pp. 75–82, Mar. 2011.
- [4] G. Liu *et al.*, "A new method on inshore ship detection in high-resolution satellite images using shape and context information," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 3, pp. 617–621, Mar. 2014.
- [5] G. Cheng *et al.*, "Object detection in remote sensing imagery using a discriminatively trained mixture model," *ISPRS J. Photogramm. Remote Sens.*, vol. 85, pp. 32–43, Nov. 2013.
- [6] J. Leitloff, S. Hinz, and U. Stilla, "Vehicle detection in very high resolution satellite images of city areas," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 7, pp. 2795–2806, Jul. 2010.
- [7] X. Huang and L. Zhang, "An SVM ensemble approach combining spectral, structural, and semantic features for the classification of high-resolution remotely sensed imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 1, pp. 257–272, Jan. 2013.
- [8] P. Zhang, Z. Lv, and W. Shi, "Object-based spatial feature for classification of very high resolution remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 6, pp. 1572–1576, Nov. 2013.
- [9] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1798–1828, Aug. 2013.
- [10] R. Salakhutdinov and G. E. Hinton, "Deep Boltzmann machines," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2009, pp. 448–455.
- [11] L. Shao, D. Wu, and X. Li, "Learning deep and wide: A spectral method for learning deep networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 12, pp. 2303–2308, Dec. 2014.

- [12] X. Bai, H. Zhang, and J. Zhou, "VHR object detection based on structural feature extraction and query expansion," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 10, pp. 1–13, Oct. 2014.
- [13] I. Dórido *et al.*, "Semisupervised self-learning for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 7, pp. 4032–4044, Jul. 2013.
- [14] W. Liao, A. Pizurica, P. Scheunders, W. Philips, and Y. Pi, "Semisupervised local discriminant analysis for feature extraction in hyperspectral images," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 1, pp. 184–198, Jan. 2013.
- [15] E. Pasolli, F. Melgani, N. Alajlan, and N. Conci, "Optical image classification: A ground-truth design framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 6, pp. 3580–3597, Jun. 2013.
- [16] F. Zheng *et al.*, "A semi-supervised approach for dimensionality reduction with distributional similarity," *Neurocomputing*, vol. 103, pp. 210–221, Mar. 2013.
- [17] G. Jun and J. Ghosh, "Semisupervised learning of hyperspectral data with unknown land-cover classes," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 1, pp. 273–282, Jan. 2013.
- [18] B. Alexe, T. Deselaers, and V. Ferrari, "Measuring the objectness of image windows," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2189–2202, Nov. 2012.
- [19] T. Deselaers, B. Alexe, and V. Ferrari, "Weakly supervised localization and learning with generic knowledge" *Int. J. Comput. Vis.*, vol. 100, no. 3, pp. 275–293, Dec. 2012.
- [20] P. Siva, C. Russell, and T. Xiang, "In defense of negative mining for annotating weakly labeled data," in *Proc. Eur. Conf. Comput. Vis.*, pp. 594–608.
- [21] P. Siva and T. Xiang, "Weakly supervised object detector learning with model drift detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2011, pp. 343–350.
- [22] F. Zhu, and L. Shao, "Weakly-supervised cross-domain dictionary learning for visual recognition" *Int. J. Comput. Vis.*, vol. 109, no. 1–2, pp. 42–59, Aug. 2014.
- [23] D. Zhang *et al.*, "Weakly supervised learning for target detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 4, pp. 701–705, Apr. 2015.
- [24] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2010.
- [25] B. Sirmacek and C. Unsalan, "Urban-area and building detection using SIFT keypoints and graph theory," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 4, pp. 1156–1167, Apr. 2009.
- [26] S. Xu, T. Fang, D. Li, and S. Wang, "Object classification of aerial images with bag-of-visual words," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 2, pp. 366–370, Apr. 2010.
- [27] D. G. Lowe, "Distinctive image features from scale-invariant keypoints" *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, Nov. 2004.
- [28] Y. Yang and S. Newsam, "Geographic image retrieval using local invariant features," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 2, pp. 818–832, Feb. 2013.
- [29] J. Wang *et al.*, "Locality-constrained linear coding for image classification," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit.*, 2010, pp. 3360–3367.
- [30] J. Han *et al.*, "An object-oriented visual saliency detection framework based on sparse coding representations," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 23, no. 12, pp. 2009–2021, Dec. 2013.
- [31] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell, "SUN: A Bayesian framework for saliency using natural statistics" *Int. J. Comput. Vis.*, vol. 8, no. 7, pp. 32, Dec. 2008.
- [32] I. Rigas, G. Economou, and S. Fotopoulos, "Low-level visual saliency with application on aerial imagery," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 6, pp. 1389–1393, Nov. 2013.
- [33] Z. Li and L. Itti, "Saliency and gist features for target detection in satellite images," *IEEE Trans. Image Process.*, vol. 20, no. 7, pp. 2017–2029, Jul. 2011.
- [34] C. M. Bishop, "Approximate inference" in *Proc. Pattern Recognit. Mach. Learn.*, New York, NY, USA, 2006, pp. 461–522.
- [35] Y. Xie, H. Lu, and M. Yang, "Bayesian saliency via low and mid level cues," *IEEE Trans. Image Process.*, vol. 22, no. 5, pp. 1689–1698, May 2013.
- [36] A. M. Cheriyyadath, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan. 2014.
- [37] M. Cramer, "The DGPF-test on digital airborne camera evaluation overview and test design" *Photogramm.-Fernerkundung-Geoinf.*, vol. 2010, no. 2, pp. 73–82, May 2010.
- [38] F. Del Frate, F. Pacifici, G. Schiavon, and C. Solimini, "Use of neural networks for automatic classification from high-resolution images," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 4, pp. 800–809, Apr. 2007.
- [39] A. Bosch, A. Zisserman, and X. Munoz, "Representing shape with a spatial pyramid kernel," in *Proc. 6th ACM Int. Conf. Image Video Retrieval*, 2007, pp. 401–408.
- [40] L. Shao, F. Zhu, and X. Li, "Transfer learning for visual categorization: A survey," *IEEE Trans. Neural Netw. Learn. Syst.*, to be published.



Junwei Han received the Ph.D. degree from Northwestern Polytechnical University, Xi'an, China, in 2003.

He is currently a Professor with Northwestern Polytechnical University. His research interests include computer vision and multimedia processing.



Dingwen Zhang received the M.S. degree from Northwestern Polytechnical University, Xi'an, China, in 2011. He is currently working toward the Ph.D. degree at Northwestern Polytechnical University.

His research interests include computer vision and weakly supervised learning.



Gong Cheng received the Ph.D. degree from Northwestern Polytechnical University, Xi'an, China, in 2013.

He is currently a Postdoctoral Fellow with Northwestern Polytechnical University. His main research interests are computer vision and remote sensing image analysis.



Lei Guo received the Ph.D. degree from Xidian University, Xi'an, China, in 1994.

He is currently a Professor with Northwestern Polytechnical University, Xi'an. His research interests include computer vision, pattern recognition, and medical image processing.



Jinchang Ren received the B.E., M.Eng., and D.Eng. degrees from Northwestern Polytechnical University, Xi'an, China, and the Ph.D. degree in electronic imaging and media communication from Bradford University, Bradford, U.K.

He is currently with the University of Strathclyde, Glasgow, U.K. His research interests focus on visual computing and multimedia signal processing, particularly on semantic content extraction for video analysis and understanding and, more recently, hyperspectral imaging.