

The Bag-of-Visual-Words Scene Classifier Combining Local and Global Features for High Spatial Resolution Imagery

Qiqi Zhu, Yanfei Zhong*, Bei Zhao, Guisong Xia, Liangpei Zhang

State Key Laboratory of Information Engineering in Surveying, Mapping, and Remote Sensing, Wuhan University
Collaborative Innovation Center of Geospatial Technology, Wuhan University
Wuhan, China

*Corresponding author E-mail: zhongyanfei@whu.edu.cn

Abstract—Scene classification has been proved to be an effective method for high spatial resolution (HSR) image semantic interpretation. Considering the complex structure and abundant information, three issues should be discussed for HSR imagery: 1) Which kind of features should be combined to comprehensively describe the HSR imagery? 2) How to efficiently fuse the different types of features? 3) Which scene classification method is best for capturing the distinctive characteristics of HSR image scenes? In this paper, an easy but effective local-global-feature bag-of-visual-words classifier (LGFBOVW) is proposed to fuse the complementary features at the histogram level. The LGFBOVW representation is then classified by support vector machine (SVM) with a histogram intersection kernel (HIK) for HSR image scene classification. LGFBOVW can incorporate distinctive features with different characteristics, whether these features are local or global, continuous or discrete. The proposed approach introduces the novel use of shape-based invariant texture index (SITI) which was originally used to analyze the natural images. SITI is captured as the global texture feature descriptor for the challenging scene representation. The mean and standard deviation values (MeanStd) is utilized as the local spectral feature descriptor, and the dense scale-invariant feature transform (SIFT) is utilized as the local structural feature. The experimental results demonstrate that the proposed method is superior to the state-of-the-art methods with UCMERGED dataset.

Keywords—scene classification, local and global features, SITI, bag-of-visual-words, fusion, high spatial resolution imagery

I. INTRODUCTION

With the ongoing development of high spatial resolution (HSR) remote sensing technology, huge quantities of HSR remote sensing images can provide detailed spatial information. Nevertheless, this type of data demonstrates the phenomenon of a complex spatial arrangement, which poses a big challenge for image classification. Object-based and contextual-based methods are available for precise object recognition [1], [2]. However, these methods have no access to the semantics in the image. This leads to the so-called “semantic gap”, namely the divergence between the information of data and the high-level knowledge [3]. How to bridge the semantic gap and make use of the strengths of HSR images is of great significance. Scene

classification, which is aimed at recognize an image among a collection of semantic classes [4], is proposed for image processing. Scene classification has now been well explored for natural image interpretation [5]–[7]. With the variability and complexity of HSR imagery, scene classification is successfully introduced to HSR image classification for its capacity of bridging the semantic gap [8]–[10]. For instance, Yang *et al.* [29] represent land-use scenes using the spatial pyramid co-occurrence. Cheriyyadat [32] exploits local spatial and structural patterns for scene classification by exploring an unsupervised feature learning approach.

Among the scene classification methods, object-based scene classification utilizes a relevant model to define the spatial relationship between the objects. For instance, Aksoy *et al.* [11] proposed an object-based scene classification method under a Bayesian framework. For object-based scene classification, prior information of objects is needed. Therefore, the bag-of-visual-words (BOVW) model has been gaining more popularity. Based on the BOVW model, the probabilistic topic model (PTM), including the probabilistic latent semantic analysis (PLSA) [12] and the latent Dirichlet allocation (LDA) [13], treat the images as a set of topics and have received wide attention [14]–[16].

BOVW is a classical and efficient intermediate feature representation method. BOVW treats the image as a document and directly represents images with a bag of visual words [17]. With no prior information required, the BOVW based scene classification method circumvents the object recognition and spatial modeling. Through direct modeling of the image scenes based on low-level features, BOVW can capture a more compact and robust representation of the image scenes compared with the PTM.

BOVW as a popular method has been very successful applied in natural image scene classification [18], [19]. In general, BOVW quantizes the features of visual words into different bins, and acquires a 1-D histogram to describe the images. Hence, the capture of the feature descriptors is very important for BOVW. In general, a single feature is used in the BOVW based scene classification, which is inadequate. For instance, in [20], Sridharan *et al.* proposed bag of lines (BoL)

to represent various linear structures in the scenes based on the low-level line feature from the scenes. Chen *et al.* [21] proposed a pyramid-of-spatial-relatons (PSR) model utilizing spatial relatons for the remotely sensed land use scene classification.

Methods combining multiple features have also been proposed, and have been developed for natural scenes. However, there are differences in the angle of view, the resolution, and the atmospheric effect between natural and HSR image scenes. The combination and fusion of different features for HSR imagery should be carefully designed, according to the distinct characteristics of HSR scenes. Zhao *et al.* [22] utilized multiple local features, i.e. SIFT descriptor, color moments, and local binary patterns (LBPs), for the BOVW based land-use scene classification. These local features can only capture the local characteristics of the HSR images. Hence, the idea of capturing the global features is motivated to comprehensively understand the semantics of the scenes.

In this paper, an easy but efficient local-global-feature bag-of-visual-words classifier (LGFBOVW) is presented for HSR image scene classification. The proposed method introduces the novel use of the shape-based invariant texture index (SITI) of texture feature to the BOVW based scene classification for challenging HSR imagery. SITI as a global scene feature focuses more on the shape characteristics of a scene, as opposed to the popular scale-invariant feature transform (SIFT) feature. SIFT focuses more of the edge attributes. LGFBOVW efficiently combines the mean and standard deviation values (MeanStd) of spectral feature, SITI of texture feature, and SIFT of structural feature. By capturing the characteristic of HSR imagery from both the local and global, discrete and continuous perspective, LGFBOVW is able to capture the distinct spatial arrangement of HSR imagery and presents a robust feature description for HSR imagery. The three features are quantized separately to generate three 1-D histograms during the k -means clustering. The LGFBOVW representation is then obtained by concatenating the three histograms, which circumvents the inadequate fusion capacity of k -means clustering. Finally, the LGFBOVW representation is classified by support vector machine (SVM) with a histogram intersection kernel (HIK).

The remainder of this paper is organized as follows. Section II describes the proposed LGFBOVW for HSR imagery scene classification. A description of the dataset and an analysis of the experiments are presented in Section III. Finally, Section IV provided the conclusions.

II. BAG-OF-VISUAL-WORDS SCENE CLASSIFIER COMBINING LOCAL AND GLOBAL FEATURES FOR HSR IMAGERY

A. Bag-of-Visual-Words Model

The bag-of-words (BOW) model has shown remarkable success in natural language processing and information retrieval. In BOW, the grammar and word order are ignored, and the document is denoted as the occurrences of the words. BOW has been widely applied to image interpretation. An image is regarded as a document, and is denoted as a collection of visual words. We call this “bag-of-visual-words” (BOVW).

The extracted features of the visual words are mapped to a codebook into different bins. The image is then transferred into a 1-D histogram of the visual word occurrences.

Hence, given an HSR image, LGFBOVW is conducted in three steps: 1) the extraction of the local features and global feature; 2) the fusion of the local features and global feature; and 3) scene classification with the classifier.

B. Local and Global Features Extraction and Fusion

Here, we specify the extraction and fusion of the multiple features. Considering the abundant spectral characteristics and the complex spatial arrangement of HSR imagery, three local and global features are designed for the HSR image scene classification:

1) *Local spectral feature*: The spectral feature is the reflection of the attributes that constitute the ground components and structures. In our experiments, the images were first uniformly sampled with a patch spacing of 4 pixels and a patch size of 8×8 pixels for the **UCMERCED** dataset. These parameter settings performed well for scene classification in our experiments. The first- statistics (the mean values) and second-order statistics (the standard deviation values) of the patches are computed in each spectral channel as the spectral feature. Let n be the number of pixels in the sampled patch, and $v_{i,j}$ denotes the j -th band value of the i -th pixel in a patch. The mean and standard deviation of the patch are calculated according to (1) and (2), respectively.

$$mean = \frac{\sum_{i=1}^n v_i}{n} \quad (1)$$

$$std = \sqrt{\frac{\sum_{i=1}^n (v_i - mean)^2}{n}} \quad (2)$$

2) *Local structural feature* [23]: The SIFT feature can circumvent the noise and changes in the illumination, as well as compensate for the deficiency of the spectral feature. During the uniform grid sampling in our experiments, the patch size and patch spacing were set to 16 and 8 for the **UCMERCED** dataset. The gray dense SIFT descriptor with 128 dimensions was extracted as the structural feature. This was inspired by previous work, in which dense features performed better for scene classification [12].

3) *Global texture feature*: The texture feature indicates the spatial distribution characteristic of the image, which can give consideration to both the macroscopic properties and fine structure. SITI [24] is employed as the texture feature, which is invariant to local contrast change and local geometrical transform. The parameters in our experiments were set according to the recommendation of the author.

The three features are then separately extracted and should be transformed into a 1-D histogram for an image. As for the local continuous features, such as the spectral or SIFT descriptor, the local features are quantized into a 1-D histogram with V_1 or V_2 bins. Here, V_1 or V_2 represents the

number of visual words. The histogram of the local discrete feature is directly utilized to describe the image. In addition, the global continuous features are stretched into a 1-D histogram with a certain scale. Suppose there are M images with a texture histogram of V_3 bins, then three histograms are fused to generate an LGFBOVW representation with $(V_1+V_2+V_3) \times M$ dimensions for all the images.

C. Local-Global-Feature Bag-of-Visual-Words Scene Classifier for HSR Imagery Scene Classification

After obtaining the LGFBOVW representation, the final classification step utilizes the effective SVM classifier with an HIK to predict the scene label. The HIK measures the degree of similarity between two histograms, to deal with the scale changes. Let $\tilde{\mathbf{V}} = (\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_M)$ be the LGFBOVW representation vectors of M images, the HIK is calculated according to (3). The HIK has been successfully applied in image classification using color histogram features [25]. With the generated LGFBOVW representation as an extended histogram, SVM with an HIK is able to enlarge the discrimination of the LGFBOVW representation vector. The core idea of SVM [26]–[28] is to effectively train a linear learning classifier, which can solve the pattern classification problem in a nonlinear way, as well as give consideration to the generalization and optimization performance. Barla *et al.* showed that histogram intersection has the required mathematical properties to be used as a kernel function for SVM [25]. The scene classification based on LGFBOVW is shown in Fig. 1.

$$K(\tilde{\mathbf{v}}_i, \tilde{\mathbf{v}}_j) = \sum_k \min(\tilde{v}_{i,k}, \tilde{v}_{j,k}) \quad (3)$$

III. EXPERIMENTS AND ANALYSIS

A. Experimental Setup

LGFBOVW was evaluated with the **UCMERCED** dataset. The challenging **UCMERCED** dataset was extracted from large images in USGS National Map Urban Area Imagery collection [26]. It consists of 21 land-use scenes (Fig. 3). Each class separately contains 100 images, which were cropped to 256×256 pixels and a 1 ft resolution.

In the experiments, considering the accuracy and the efficiency, the number of visual words for MeanStd and SIFT were optimally set to 1000, with 350 for SITI. The penalty parameter for SVM was set to 10. To test the stability of the proposed LGFBOVW, the different methods were executed 100 times by a random selection of training samples, to obtain convincing results for the **UCMERCED** dataset. Following the experimental setup in [29], 80 samples were randomly selected per class from the **UCMERCED** dataset for training.

B. Experimental Result Analysis

The classification performances of different procedures with the **UCMERCED** dataset are reported in Table I. **LDA(SIFT)**, **LDA(MeanStd)**, and **LDA(SITI)** denote the scene classification methods in [30], using SIFT, MeanStd, and SITI features, respectively. **BOVW(SIFT)**, **BOVW(MeanStd)**, and **BOVW(SITI)** denote the BOVW methods using SIFT, MeanStd, and SITI features, respectively. **LGFBOVW-Linear** represents **LGFBOVW** utilizing SVM with a linear kernel. The classification performances of the different methods with the **UCMERCED** dataset are reported in Table II. **SPM** [31] employed dense gray SIFT, and the spatial pyramid layer was optimally selected as one.

From Table I, it can be seen that **BOVW(SIFT)**, **BOVW(MeanStd)**, and **BOVW(SITI)** outperform **LDA(SIFT)**, **LDA(MeanStd)**, and **LDA(SITI)**, respectively. This indicates that BOVW is a competitive method compared to LDA when applied to HSR image scene classification. **LGFBOVW** is superior to **BOVW(SIFT)**, **BOVW(MeanStd)**, **BOVW(SITI)**, and **LGFBOVW-Linear**. This proves that the proposed **LGFBOVW** can overcome the insufficient feature capture of the single-feature strategies. In addition, the simple HIK outperforms the commonly used linear kernel in the SVM classifier. In Table II, the proposed **LGFBOVW** is superior to the performance of **SPM** [31], the **Yang and Newsam** method [29], and the **Cheriyadat** method [32], and exceeds the state-of-the-art with **UCMERCED** dataset.

TABLE I. THE CLASSIFICATION PERFORMANCES FOR THE DIFFERENT PROCEDURES WITH THE UCMERCED DATASET

Methods	Accuracy (%)
LDA(SIFT)	81.83 ± 1.88
BOVW(SIFT)	87.31 ± 1.40
LDA(MeanStd)	81.27 ± 2.01
BOVW(MeanStd)	85.30 ± 1.67
LDA(SITI)	79.73 ± 2.64
BOVW(SITI)	81.54 ± 1.46
LGFBOVW-Linear	88.33 ± 1.56
LGFBOVW	96.88 ± 1.32

TABLE II. THE CLASSIFICATION PERFORMANCES FOR THE DIFFERENT METHODS WITH THE UCMERCED DATASET

Methods	Accuracy (%)
SPM [31]	82.30 ± 1.48
Yang and Newsam [29]	81.19
Cheriyadat [32]	81.67 ± 1.23
LGFBOVW	96.88 ± 1.32

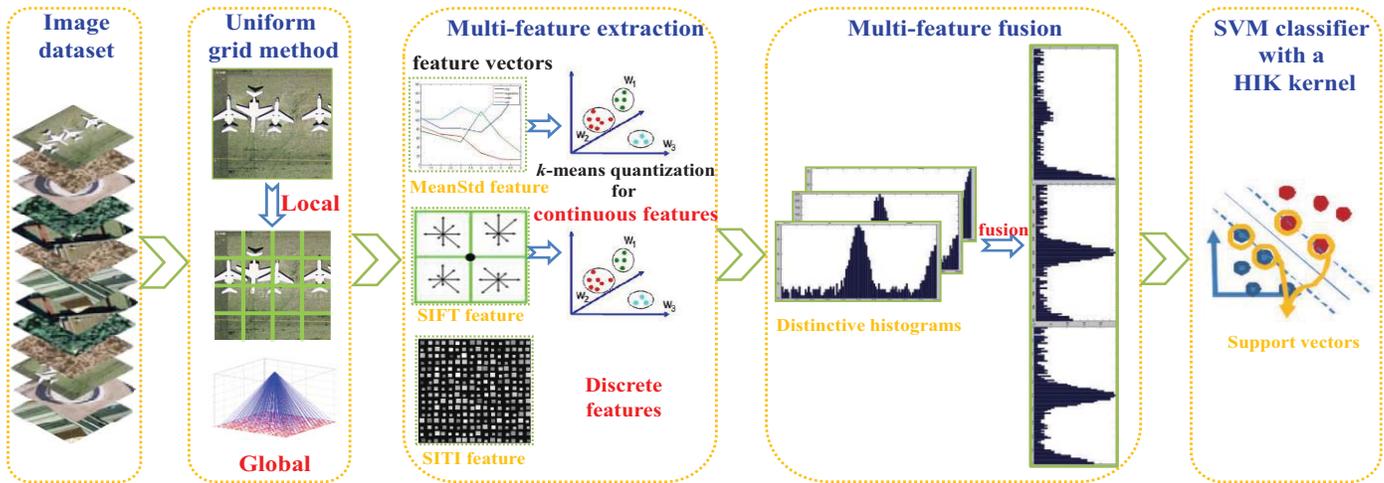


Fig. 1. The proposed HSR scene classification based on LGFBOVW.

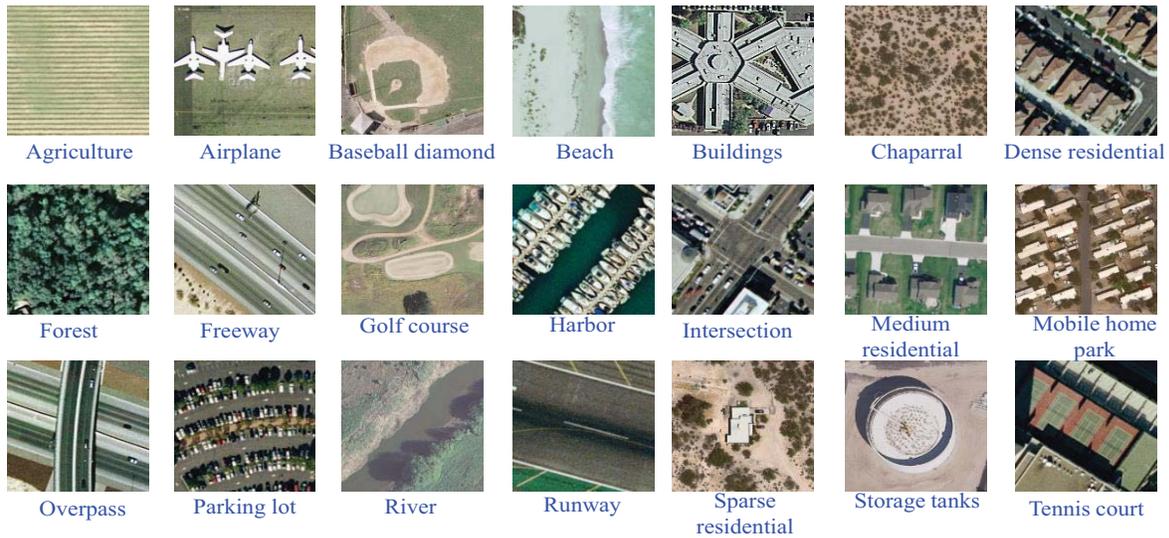


Fig. 2. Example images of the UCMERCED dataset.

IV. CONCLUSION

Considering the three issues of the HSR imagery mentioned in the Abstract, we have designed an easy but effective method—local-global-feature bag-of-visual-words classifier (LGFBOVW)—for challenging HSR image scene classification. The novel use of the shape-based invariant texture index (SITI) of texture feature is introduced to the BOVW, which can capture the global characteristics of the HSR imagery.

LGFBOVW captures three complementary features from both the local and global, discrete and continuous perspective, based on investigating the distinct characteristics of HSR imagery from text information and natural scenes. The fusion of the three features at the

histogram level and the incorporation of SVM with a HIK are effective in increasing the discrimination of different scenes. The experimental results show that SITI can achieve good performance for the HSR image scene classification. And LGFBOVW is superior to the state-of-the-art methods with UCMERCED dataset. In our future work, we plan to consider other models which can relax the normalization constraints of the probabilistic topic model.

ACKNOWLEDGEMENT

This work was supported by National Natural Science Foundation of China under Grant No. 41371344, the Fundamental Research Funds for the Central Universities under Grant No.2042014kf00231, Program for Changjiang Scholars and Innovative Research Team in University under

Grant No. IRT1278, 863 High Technology Program of the People's Republic of China under Grant No. 2013AA12A301, and State Key Laboratory of Earth Surface Processes and Resource Ecology under Grant No. 2015-KF-02.

REFERENCES

- [1] T. Blaschke, "Object based image analysis for remote sensing," *ISPRS J. Photogramm. Remote Sens.*, vol. 65, no. 1, pp. 2–16, Jan. 2010.
- [2] J. C. Tilton, Y. Tarabalka, P. M. Montesano, and E. Gofman, "Best merge region-growing segmentation with integrated nonadjacent region object aggregation," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 11, pp. 4454–4467, Nov. 2012.
- [3] D. Bratananu, I. Nedelcu, and M. Datcu, "Bridging the semantic gap for satellite Image annotation and automatic mapping applications," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens. (JSTARS)*, vol. 4, no. 1, pp. 193–204, Mar. 2011.
- [4] A. Bosch, X. Munoz, and R. Martí, "Which is the best way to organize/classify images by content?," *Image Vision Comput.*, vol. 25, pp. 778–791, Jul. 2007.
- [5] D. Tao, L. Jin, Z. Yang, X. Li and L. Xuelong, "Rank Preserving Sparse Learning for Kinect Based Scene Classification," *IEEE Trans. Cybern.*, vol. 43, no. 5, pp. 1406–1417, Oct. 2013.
- [6] B. Fernando, E. Fromont, and T. Tuytelaars, "Mining Mid-level Features for Image Classification," *Int. J. Comput. Vis.*, pp. 1–18, Feb. 2014.
- [7] J. Luo, M. Boutell, R. T. Gray, and C. Brown, "Image transform bootstrapping and its applications to semantic scene classification," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 35, no. 3, pp. 563–570, Jun. 2005.
- [8] W. Shao, W. Yang, G.-S. Xia, and G. Liu, "A Hierarchical Scheme of Multiple Feature Fusion for High-Resolution Satellite Scene Categorization," *Comput. Vision Syst.*, pp. 324–333, 2013.
- [9] D. Dai and W. Yang, "Satellite image classification via two-layer sparse coding with biased image representation," *IEEE Geosci. Remote Sens. Lett.*, vol. 8, no. 1, pp. 173–176, Jan. 2011.
- [10] G. Sheng, W. Yang, T. Xu, and H. Sun, "High-resolution satellite scene classification using a sparse coding based multiple feature combination," *Int. J. Remote Sens.*, vol. 33, no. 8, pp. 2395–2412, 2012.
- [11] S. Aksoy, K. Koperski, C. Tusk, G. Marchisio, and J. C. Tilton, "Learning Bayesian classifiers for scene classification with a visual grammar," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 3, pp. 581–589, Mar. 2005.
- [12] T. Hofmann, "Unsupervised learning by probabilistic latent semantic analysis," *Mach. Learn.*, vol. 42, pp. 177–196, Jan. 2001.
- [13] D. Blei, A. Ng, and M. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Jan. 2003.
- [14] M. Liénou, H. Maître, and M. Datcu, "Semantic annotation of satellite images using latent Dirichlet allocation," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 1, pp. 28–32, Jan. 2010.
- [15] Y. Zhong, Q. Zhu, and L. Zhang, "Scene classification based on the multifeature fusion probabilistic topic model for high spatial resolution remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. PP, no. 99, pp. 1–16, 2015.
- [16] C. Văduva, I. Gavăt, and M. Datcu, "Latent Dirichlet allocation for spatial analysis of satellite Images," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 5, pp. 2770–2786, May. 2013.
- [17] L. Weizman and J. Goldberger, "Urban-area segmentation using visual words," *Remote Sens. Lett.*, vol. 6, no. 3, pp. 388–392, 2009.
- [18] L. Zhou, Z. Zhou, and D. Hu, "Scene classification using a multi-resolution bag-of-features model," *Pattern Recognit.*, vol. 46, no. 1, pp. 424–433, Jan. 2013.
- [19] Y. Huang, K. Huang, C. Wang, and T. Tan, "Exploring relations of visual codes for image classification," in *Proc. CVPR*, 2011, pp. 1649–1656.
- [20] H. Sridharan and A. Cheriyyadat, "Bag of Lines (BoL) for Improved Aerial Scene Representation," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 3, pp. 676–680, 2015.
- [21] 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, 2015.
- [22] L.-J. Zhao, P. Tang, and L.-Z. Huo, "Land-use scene classification using a concentric circle-structured multiscale bag-of-visual-words model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 12, pp. 4620–4631, 2014.
- [23] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91–110, 2004. L. Fei-Fei and P. Perona, "A Bayesian hierarchical model for learning natural scene categories," in *Proc. CVPR*, 2005.
- [24] G.-S. Xia, J. Delon, and Y. Gousseau, "Shape-based invariant texture indexing," *Int. J. Comput. Vision*, vol. 88, no. 3, pp. 382–403, Jul. 2010.
- [25] A. Barla, F. Odone, and A. Verri, "Histogram intersection kernel for image classification," in *Proc. IEEE Int. Conf. Image Processing*, vol. 3, pp. III–513–16, 2003.
- [26] V. Kecman, *Learning and Soft Computing: Support Vector Machines, Neural Networks and Fuzzy Logic Models*. Cambridge, MA: MIT Press, 2001.
- [27] L. Wang, *Support Vector Machines: Theory and Application*. Berlin, German: Springer, 2005.
- [28] J. Suykens and J. Vanderwalle, "Least squares support vector machines classifiers," *Neural Process. Lett.*, vol. 9, no. 3, pp. 293–300, Jun. 1999.
- [29] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. ACM SIGSPATIAL GIS*, 2010, pp. 270–279.
- [30] B. Zhao, Y. Zhong, and L. Zhang, "Scene classification via latent Dirichlet allocation using a hybrid generative/discriminative strategy for high spatial resolution remote sensing imagery," *Remote Sens. Lett.*, vol. 4, no. 12, pp. 1204–1213, Dec. 2013.
- [31] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2006, vol. 2, pp. 2169–2178.
- [32] A. M. Cheriyyadat, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan. 2014.