

# FAST BINARY CODING FOR SATELLITE IMAGE SCENE CLASSIFICATION

Fan Hu<sup>1,2</sup>, Zifeng Wang<sup>1</sup>, Gui-Song Xia<sup>1</sup>, Bin Luo<sup>1</sup>, Liangpei Zhang<sup>1</sup>

<sup>1</sup> State Key Laboratory LIESMARS, Wuhan University, Wuhan 430072, China

<sup>2</sup> Electronic Information School, Wuhan University, Wuhan 430072, China

## ABSTRACT

Feature extraction is at the core of satellite scene classification task. In this paper, we propose a fast binary coding (FBC) method to effectively generate the global discriminative feature representation of image scenes. Equipped with unsupervised feature learning technique, we first learn a set of optimal “filters” from large quantities of randomly sampled image patches, and then we obtain feature maps by convolving image scene with the learned filter bank. After binarizing the feature maps, a simple skillful conversion of binary-valued feature map to integer-valued feature map is performed. The final statistical histograms, which are considered as the global feature representations of scenes, are computed on the integer-valued feature map similar to the conventional BOW model. Experiments on two datasets demonstrate that the proposed FBC achieve satisfying classification performance as well as has much faster computational speed compared with traditional scene classification methods.

**Index Terms**— Scene classification, filters, feature representation, binary coding

## 1. INTRODUCTION

Scene classification of large high-resolution remotely sensed images [1–3] is a fundamental yet challenging problem in intelligent remote sensing field, playing a significant role in urban planning, land resource management, computer cartography, and many more. Remote sensing scenes in this case refer to some separated subareas which contain specific semantic meaning, such as the residential area, industrial area, commercial area and green land in a typical urban area satellite image. However, the high complexity of spatial and structural patterns in the massive HR images makes the intelligent scene understanding and classification a challenging problem. In order to accurately obtain the scene classes, generating discriminative holistic feature representation for each scene is a key step.

The typical bag-of-visual-words (BOW) [4] scheme, which represents each image scene with a histogram where each bin counts the occurrence frequency of features on a codeword, is probably the most popular and effective scene

classification framework. There are three basic steps in the BOW pipeline for scene classification: extracting local feature descriptors, generating codewords, encoding local features. Of the three steps, feature extraction is the core part and greatly influences the final classification performance. For the purpose of high classification performance for different image scene datasets, it is crucial to choose or design powerful feature descriptors. However, designing good features needs too much human effort and expert-domain knowledge. Nowadays learning features automatically from only unlabeled images in unsupervised ways has been a new tide. Moreover, in BOW framework the step of generating codewords, where the codewords are typically generated by clustering (e.g., K-means) over local features, is usually time-consuming. Some binarized feature representations method [5, 6] have become increasingly popular, which are very simple and efficient to compute.

In this paper, we present a fast binary coding scheme for feature representation of image scenes. We first randomly sample amount of local image patches from images in dataset, and apply proper unsupervised learning techniques to learn a dictionary, which is regarded as a set of filters. Then convolve each image scene with the learnt filters and binarize the filter responses according to a predefined threshold. Finally, we convert the binary responses back into a single decimal number, and then compute the histogram of the decimal values for each image scene. The final histogram is considered as the holistic feature representation of the image, which can be fed into the classifier for training and testing. In contrast to the typical BOW pipeline, we neither use any hand-crafted features nor feature encoding techniques, and therefore greatly improve the computational efficiency. When the set of filters have been generated, the holistic histogram of each image scene can be yielded extremely fast on common CPUs. Extensive experiments show that we can obtain comparable classification performance at a low computational cost.

This work was partially founded by NSFC project No.91338113.

## 2. FAST BINARY CODING FOR SCENE CLASSIFICATION

### 2.1. Fast Binary Coding (FBC)

Suppose that we are given a image scene  $I$  of size  $m \times n$  and  $K$  linear filters  $\{\mathbf{W}^{(k)}\}_{k=1}^K$ , the  $k$ th filter response  $f^{(k)}$  (also referred to feature map) through convolving the image with the set of filter bank can be described as:

$$f^{(k)} = I * \mathbf{W}^{(k)}, k = 1, \dots, K \quad (1)$$

where the operator “ $*$ ” denotes image 2D convolution. In order to make the size of  $f^{(k)}$  identical to  $I$ , we pad zeros around the boundary of  $I$  before the convolution. Each value in the filter response can be interpreted as a descriptor of the local region centered the corresponding pixel in the  $I$ . Finally the image scene  $I$  outputs  $K$  real-valued filter responses, and in other words we can get a  $K$ -dimensional real-valued feature for each pixel. We binarize all the  $K$  responses and obtain the binarized maps  $\{B^{(k)}\}_{k=1}^K$ :

$$B^{(k)} = H\left(f^{(k)}\right), k = 1, \dots, K \quad (2)$$

where  $H(\cdot)$  is a Heaviside step function, which outputs one if  $f^{(k)} > 0$  and zero otherwise.

For each pixel, the  $K$ -dimensional real-valued feature is now transformed into a  $K$ -bit binary string. We can consider the binary strings as a binary-valued number, and convert it back into a single integer value by the following fomula:

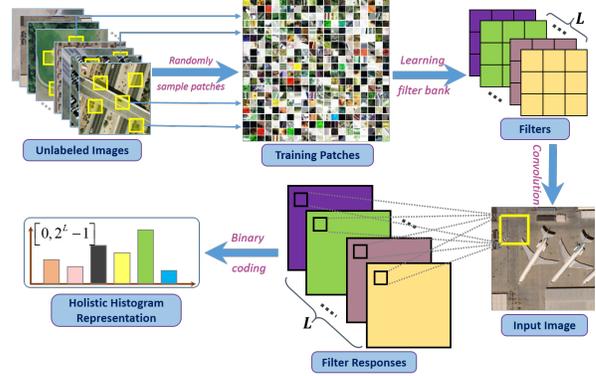
$$I^{new}(x, y) = \sum_{k=1}^K 2^{k-1} \cdot B^{(k)}(x, y) \quad (3)$$

where  $I^{new}$  is new generated “image” after the conversion of binary maps, and the  $x, y$  denote the position coordinate index. We can note that the value of pixels in  $I^{new}$  is within the range of  $[0, 2^K - 1]$ . Analog to the conventional BOW model, each integer value is regarded as a codeword, and thereby the size of codebook results in  $2^K$ . A statistical histogram  $Y \in \mathbb{R}^{2^K}$  is computed on this codebook, which is the resulting global feature representation for image  $I$ .

In our FBC method, the length of binary strings (i.e., the number of bit) depend on the number of predefined filters  $K$ . In addition, the dimensions of the final features for image scene is also closely related to  $K$ , and exponentially grow with  $K$ . Our empirical experience suggests that we should set a relatively small  $K$  to make the length of the global features acceptable, and avoid overfitting of the classifier. The whole stage for generating feature representation via FBC is illustrated in Fig. 1.

### 2.2. Scene Classification

We now can efficiently compute the global feature for each image scene in a dataset via the proposed FBC method. These



**Fig. 1.** Illustration of the fast binary coding for global feature representation of image scenes

histogram features can be directly fed into an off-the-shelf classifier for the classification task. In this paper, the SVM classifier with histogram intersection kernel (HIK) is utilized to train and predict labels for new image scenes. The HIK is defined as:

$$I(Y_i, Y_j) = \sum_{p=1}^{2^K} \min(Y_i(p), Y_j(p)) \quad (4)$$

where  $H_i$  and  $H_j$  is histogram feature representation for image  $i$  and  $j$ .

### 2.3. Learn filers via unsupervised learning methods

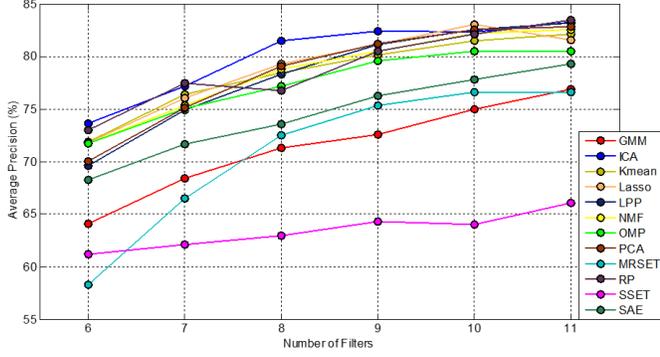
In the FBC pipeline, any kind of linear filters can used to convolve with image, and hence we want to discover what kind of filters are most suitable to FBC feature representation. For purpose of generating desirable filters, we are inspired by the unsupervised feature learning (UFL) methods and attempt to learn suitable filters by some widely-used unsupervised learning algorithms. Similar to UFL pipeline, we first perform some steps to learn a dictionary:

- Randomly extract a large number of image patches from the image dataset
- Normalize each patch to zero mean and unit variance as well as the ZCA whitening as a pre-processing stage
- Train the dictionary  $D$  via a proper unsupervised learning method

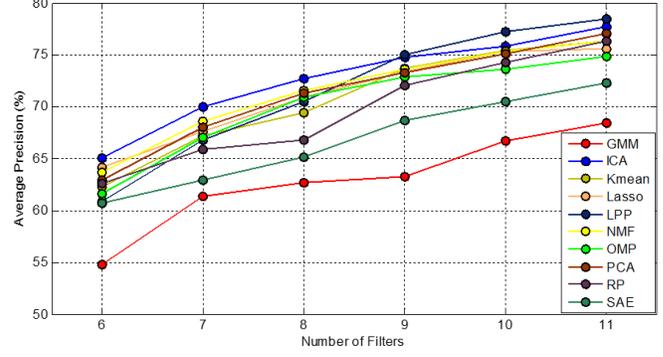
here we presents three typical method (but not limited to) to learn dictionary.

#### 2.3.1. $K$ -means

Given the normalized patch vectors  $x^{(i)}, i = 1, 2, \dots, M$ ,  $K$ -means can learn a dictionary  $D$  containing  $K$  cluster centers.



(a) Classification accuracies on UCM dataset



(b) Classification accuracies on WHU-RS dataset

**Fig. 2.** Classification results on two datasets with different learned linear filters.

The objective function is formulated by minimize the distance between training samples and their cluster centroids, which is defined as:

$$\min_{D,c} \sum_i \|Dc^{(i)} - x^{(i)}\|_2^2 \quad (5)$$

$$s.t. \|D^{(k)}\|_2 = 1, \forall k, \text{ and, } \|c^{(i)}\|_0 \leq 1, \forall i \quad (6)$$

where  $c^{(i)}$  is the assignments of the sample  $x^{(i)}$  to the clusters. Each learned centroid  $D^{(k)}$  can be regarded as a filter  $\mathbf{W}^{(k)}$  by simply resize it to the original size of image patch.

### 2.3.2. Sparse Coding

Given the normalized patch vectors  $x^{(i)}, i = 1, 2, \dots, M$ , the objective of learning the dictionary  $D$  can be defined as:

$$\min_{D,\alpha} \sum_i \|D\alpha^{(i)} - x^{(i)}\|_2^2 + \lambda \|\alpha^{(i)}\|_1 \quad (7)$$

$$s.t. \|D^{(k)}\|_2 \leq 1, \forall k \quad (8)$$

where  $\alpha$  denotes the sparse vectors. We can easily optimize this objective by online learning techniques. Each basis vector  $D^{(k)}$  in the dictionary  $D$  is a learned linear filter  $\mathbf{W}^{(k)}$ .

### 2.3.3. Principal Component analysis (PCA)

PCA is probably the most common unsupervised linear dimensionality reduction technique. The main goal of PCA is to iteratively find orthogonal directions maximizing variance of samples, or it can be cast as low-rank matrix factorization problem:

$$\min_U \sum_i \|X - UU^T X\|_2^2 \quad (9)$$

$$s.t. UU^T = \mathbf{I} \quad (10)$$

where  $\mathbf{I}$  is the identity matrix and  $X$  is a matrix of concatenating all training patch vectors. We use the first  $K$  principal eigenvectors of matrix  $XX^T$  as the set of linear filters.

## 3. EXPERIMENTS AND ANALYSIS

We evaluate the proposed FCB method on two public land-use scene datasets, which are: (1) *UCM dataset* [7], consisting of 21 scene categories with 100 samples (size of  $256 \times 256$  pixels) per class. (2) *WHU-RS dataset* [8], consisting of 19 satellite scene categories with 50 samples (size of  $600 \times 600$  pixels) per class.

We perform 10 typical unsupervised learning methods, which are Gaussian mixture model (GMM), K-means, PCA, independent component analysis (ICA), locality preserving projection (LPP), Lasso, orthogonal matching pursuit (OMP), sparse auto-encoder (SAE), non-negative matrix factorization (NMF) and random projection (RP) [9] to learn a set of filters prepared for FCB, aiming to find the most appropriate method which can explore the natural statistical properties of image patches. These methods can learn different linear filters that has various filtering properties. In our experiments, all the image patches used for learning filters are beforehand mean-removed and ZCA whitened. At the classification stage, We choose 80 samples of each class for SVM training and the rest for testing.

The classification results with different learned filters are reported in Fig. 2. The results show that performs better than other unsupervised techniques consistently on two datasets when filter size is less than or equal to 9. When filter size is greater than 9, PCA, NMF, LPP and ICA show comparative performance on UCM dataset; LPP has a leading performance on WHU-RS dataset. GMM leads to the worst performance, and apparently GMM is not a proper one for learning filters. On the whole, the performance gradually improves as the number of filter increases. In addition, it is interesting that the RP whose filters are generated from zero-mean, unit-variance normal distribution achieves acceptable accuracies. We also compare the capability of the hand-engineered filters with the learned filters in our FBC framework. The maximum response filters (MR) and the Schmid filters (S), which are specially designed for texture recognition, are investigated;

the results in Fig. 2(a) show that a majority of learned filters perform far better (over 10 percent accuracy) than these two specially-designed filters.

The effect of the filter size on classification performance is shown in Fig. 3. Filters learned by K-means and Lasso are respectively evaluated here. The results show that classification performance consistently improves as number of filter grows except some rare cases. In general, we can see that performance severely decreases with too small or too large filter size. The optimal size of filter by K-means and Lasso is different: filters of size  $5 \times 5$  yield the best accuracy when learned by K-means; filters of size  $7 \times 7$ ,  $9 \times 9$ ,  $11 \times 11$  lead to comparative performance when learned by Lasso. We can infer that filters learned via different learning algorithm determine differential best filter sizes.

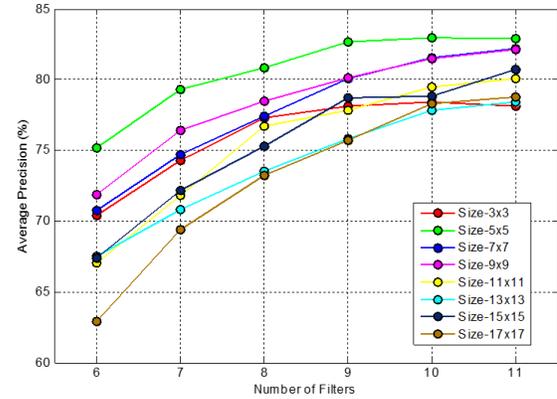
Note that when the filters are obtained, the whole stage of feature representation by FBC for images scenes is fairly fast because FCB is a totally feed-forward process and do not contain the extraction of low-level features as well as complex feature encoding and pooling steps compared with the typical scene classification pipeline. In our experiments, it only takes less than 1 minute extracting features for all 2100 images of UCM dataset, and is nearly  $60 \times$  faster than BOW model.

#### 4. CONCLUSION

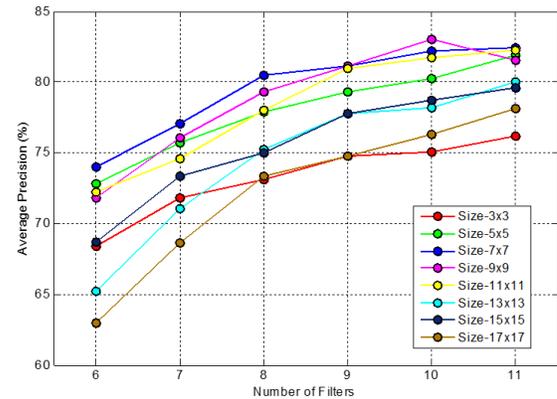
This paper presents an effective method for extracting global feature representation of satellite image scenes. In our FBC pipeline, we introduce unsupervised learning techniques to automatically learn optimal filters from a large amount of unlabeled image patches, and through binarizing the feature responses, we can skillfully compute the feature representations for image scenes in a computationally efficient way. The number of filters and filter size are two important parameters in FBC. Filters learned via different ways also lead to much different classification performance. A lot of experiments show that scene classification based on FBC can achieve satisfactory accuracy and save much computation cost compared with the classic scene classification model.

#### 5. REFERENCES

- [1] A.M. Cheriyyadath, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan 2014.
- [2] F. Hu, G.-S. Xia, Z. Wang, L. Zhang, and H. Sun, "Unsupervised feature coding on local patch manifold for satellite image scene classification," in *IEEE Geoscience and Remote Sensing Symposium, IGARSS 2014*, 2014, pp. 1273–1276.
- [3] W. Yang, X. Yin, and G.-S. Xia, "Learning high-level features for satellite image classification with limited labeled samples," *IEEE T. Geoscience and Remote Sensing*, vol. 53, no. 8, pp. 4472–4482, 2015.



(a) Classification results with K-means



(b) Classification results with Lasso

**Fig. 3.** Results with different filter size when the filters are learned via K-means and Lasso.

- [4] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2006, vol. 2, pp. 2169–2178.
- [5] Ojala T., Pietikäinen M., and Mäenpää T., "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 971–987.
- [6] J. Kannala and E. Rahtu, "Bsfif: Binarized statistical image features," in *Pattern Recognition (ICPR), 2012 21st International Conference on*, Nov 2012, pp. 1363–1366.
- [7] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2010, pp. 270–279.
- [8] G.-S. Xia, W. Yang, J. Delon, Y. Gousseau, H. Sun, and H. Maitre, "Structural high-resolution satellite image indexing," in *ISPRS, TC VII Symposium Part A: 100 Years ISPRS - Advancing Remote Sensing Science*, July 2010.
- [9] L. Liu and P. Fieguth, "Texture classification from random features," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 3, pp. 574–586, 2012.