

BEYOND LOCAL IMAGE FEATURES: SCENE CLASSIFICATION USING SUPERVISED SEMANTIC REPRESENTATION

Chunjie Zhang¹, Jing Liu², Chao Liang², Jinhui Tang³, Hanqing Lu²

¹Graduate University of Chinese Academy of Sciences, 100049, Beijing, China

²National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences
P.O. Box 2728, Beijing, China

³School of Computer Science, Nanjing University of Science and Technology, Nanjing, China
ivazhangchunjie@gmail.com, {jliu, cliang, luhq}@nlpr.ia.ac.cn, jinhuitang@mail.njust.edu.cn

ABSTRACT

The use of local features for image representation has been proven very effective for a variety of visual tasks such as object localization and scene classification. However, local image features carry little semantic information which is potentially not enough for high level visual tasks. To solve this problem, in this paper, we propose to use a supervised semantic image representation for scene classification, where an image is represented as a response histogram. This response histogram is a combination of the prediction of pre-trained generic object classifiers and classifiers generated by supervised learning. Besides, the use of sparsity constraints makes the proposed representation more efficient and effective to compute. Performances on the UIUC-Sports dataset, the MIT Indoor scene dataset and the Scene-15 dataset demonstrate the effectiveness of the proposed method.

Index Terms— Scene classification, sparse, supervised learning, semantic representation

1. INTRODUCTION

As an important problem for computer vision, scene classification has received considerable attention in the last few years. The state-of-the-art methods often used local feature based image representation for scene classification [1-3] and model an image as an order-less collection of local features. The local features are quantized based on a clustering method, such as *k*-means clustering. The clustering centers are used as visual words. Each local feature is then assigned to the nearest visual words. This bag-of-visual-words representation is inspired by the bag-of-words model in text retrieval [4]. However, the histogram representation of images ignores the spatial information, to overcome this problem, Lazebnik *et al.* [2] proposed a spatial pyramid matching (SPM) method which is widely used by researchers since its introduction.

Although inspired by text retrieval, the visual word has no explicit semantic meanings. Besides, the local features carry little semantic information which is potentially not enough for semantically classifying scene images. To overcome these problems, researchers resorted to techniques from text processing literature, such as Latent Dirichlet Allocation (LDA) [5], Probabilistic Latent Semantic Analysis (pLSA) [6]. This latent representation of image helps to bridge the semantic gap between the visual features and the semantic concepts. However, this latent representation is still not explicit for human understanding which limits its discriminative power.

The use of semantic representation of images becomes popular in recent years [7-9]. Instead of using visual word histogram, each image is represented by the probability distribution of pre-defined object classes. These pre-defined object classes can be obtained by either using the training images [7, 8] or by using images generated from other sources [9]. Rasiwasia and Vasconcelos [7] tried to classify scene with low-dimensional semantic spaces and weak supervision from casual image annotations. They introduced a low dimensional semantic “theme” image representation and represent each image as vectors of posterior theme probabilities which outperformed the unsupervised latent-space methods. Carneiro *et al.* [8] proposed a probabilistic formulation for semantic image annotation and retrieval. Each class is defined as the group of database images labeled with a common semantic label. This method is conceptually simple and do not require prior semantic segmentation of training images. Motivated by [7, 8], Li *et al.* [9] proposed the Object Bank which represented an image as a scale-invariant response map of a large number of pre-trained generic object detectors which is blind to the testing task. We can observe that researchers used the training images only [7, 8] or used generic images without considering the discriminative power of training images. The performance can be further improved if we combine the discriminative power of training images with generic images in a unified manner.

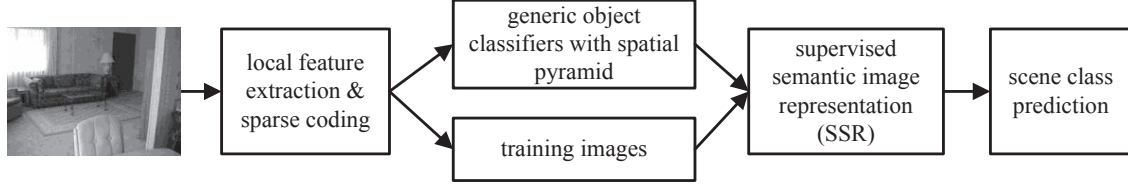


Figure 1. Flowchart of the proposed supervised semantic image representation for scene classification method.

In this paper, we propose a novel scene classification method by using supervised semantic representation of images. To represent an image, we use a histogram which is the combination of the prediction values of pre-trained generic object classifiers and classifiers generated by the training images. Sparsity constraints are used to choose the most discriminative classes and speed up computation. The effectiveness of the proposed method is demonstrated by the experimental results on the UIUC-Sports dataset, the MIT Indoor scene dataset and the Scene-15 dataset. Figure 1 gives the flowchart of the proposed method.

The rest of the paper is organized as follows. We give the details of the proposed scene classification using supervised semantic image representation method in Section 2. The experimental results on the UIUC-Sports dataset, the MIT Indoor scene dataset and the Scene-15 dataset are given in Section 3. Finally, we conclude in Section 4.

2. SUPERVISED SEMANTIC REPRESENTATION FOR SCENE CLASSIFICATION

We give the details of the proposed supervised semantic representation for scene classification method in this section. In order to represent an image, we combine the discriminative power of generic object classifiers [9] with classifiers generated by training data in an efficient way. Sparsity constraints are then used to choose the most discriminative features for efficient scene classification.

2.1. Supervised semantic representation of images

The use of semantic representation of images has been proven very effective by many researchers [7-9]. Some researchers [7, 8] used the training images for semantic representation while others [9] leveraged generic object categories. To combine the benefits of these methods, we propose to use the supervised semantic representation of images for scene classification.

The generic object classifier bank is proposed by Li *et al.* [9]. Two state-of-the-art methods are used to generate the object classifiers, the latent SVM object detectors [10] for blobby objects such as tables, cars, etc and texture classifier [11] for objects such as sky, road, and sand. The word “object” is used in very general form. We use the 200 objects at 3 spatial pyramid levels and 12 detection scales

provided by [9]. Each image is then represented as a response histogram which has the dimension of $200 \times 3 \times 12 = 7200$. Formally, we use $h_{g,i}$, $i = 1, 2, \dots, N$ as the generic histogram representation of the i -th image, where N is the number of images. This object bank based image representation ignores the training images and represents every image using the same classifiers. The discriminative power can be further improved by using the discriminative information of training images along with this generic image representation.

To combine the discriminative information of training images, we train linear SVM classifier with histogram intersection kernel using the training data. We use the linear SVM classifier for efficient computation. Each image category can then be predicted by the learned classifiers. Let $p_{i,j} \in \mathbb{R}^{M \times 1}$, $i = 1, 2, \dots, N$, $j = 1, 2, \dots, M$ be the predicted values of the i -th image with the j -th class, where M is the number of training image categories. We try to combine the supervised information with generic image representation by using $p_{i,j}$ to reweight the generic image representation $h_{g,i}$. Let $h_{i,j} = p_{i,j} \times h_{g,i}$, the proposed supervised semantic image representation is a long vector $h_i = [h_{i,1}, \dots, h_{i,j}, \dots, h_{i,M}]$ which has $7200 \times M$ dimension.

2.2. Scene classification using supervised semantic image representation

After representing images using the supervised semantic representation, we can predict the categories of images by training classifiers. Let $H = [h_1, h_2, \dots, h_N]$ be the supervised semantic representation of the N training images. $Y = [y_1, y_2, \dots, y_N] \in \{0, 1\}^N$ is the label of the corresponding images. We try to learn a linear classifier f so that:

$$y_i = f(h_i) = \alpha^T h_i \quad (1)$$

where $\alpha \in \mathbb{R}^{7200M \times 1}$ is the function parameter to be learned. However, for a particular scene classification task, not all of the features are equally useful; some of them may

even hinder the discriminative power. It is more efficient to choose the most discriminative features for prediction. Usually, this can be achieved by solving problem as:

$$\alpha = \arg \min_{\alpha} \sum_{i=1}^N L(\alpha^T h_i, y_i) + \lambda R(\alpha) \quad (2)$$

where $L(.,.)$ is the loss function and $R(.)$ is the regularization term with λ as the regularization parameter. We use:

$$L(\alpha^T h_i, y_i) = \log(1 / (\exp(0.5 y_i \times \alpha^T h_i))) \quad (3)$$

as the loss function for computational efficiency. As to the regularization term in problem (2), we choose to use the sparsity constraints. The use of sparsity has become popular and shown to be very effective in recent years [9, 12]. We follow [9] and use the joint sparsity via $\ell^1 / \ell^2 + \ell^1$ which has the form as:

$$R(\alpha) = \|\alpha\|_{1,2} + \|\alpha\|_1 \quad (4)$$

Where $\|\alpha\|_{1,2} = \sum_{m=1}^M \|\alpha^m\|_2$. This jointly sparsity term

$R(\alpha)$ controls the sparsity of α and also ensures that features predicted by the same classifier to be jointly zero, hence can choose the most discriminative features for scene classification. The optimization problem of (2) can then be solved using the coordinate descent algorithm [9].

3. EXPERIMENTS

We evaluate the proposed supervised semantic image representation (SSR) method for scene classification on 3 public datasets: the UIUC-Sports dataset [13], the MIT Indoor scene dataset [14] and the Scene-15 dataset [2]. We densely extract SIFT descriptors [15] on overlapping 16×16 pixels with an overlap of 8 pixels. The codebook size is set to 1,024 for the three datasets. Sparse coding along with max pooling [16] is used to generate codebook and encode local features. We use the one-versus-all rule for multi-class classification. A classifier is learned to separate each class of images from the rest images. The test image is assigned the label of classifiers with the highest responses. The average of per-class classification rates is used for quantitative performance comparison.

3.1. UIUC-Sports Dataset

The UIUC-Sports dataset has eight categories of 1,792 images with the eight categories as: *badminton, bocce, croquet, polo, rock climbing, rowing, sailing and snow boarding*. The number of images per categories ranges from 137 to 250. We follow the same experimental setup as in [13] and randomly choose 70 images per class for training and the rest images for testing. For fair comparison, we repeat this process for five times.

Table 1. Performance comparison on the UIUC-Sports dataset. OB: object bank; TIM: the integrative model; SCSPM: sparse coding spatial pyramid matching re-implemented by [17].

Methods	Performance
OB [9]	76.30
TIM [13]	73.40
SCSPM [17]	82.74 ± 1.46
SSR	83.53 ± 1.27

Table 1 gives the scene classification comparison of the proposed method and methods in [9, 13, 17] for the UIUC-Sports dataset. Li *et al.* [9] used Object Bank for scene classification while Li and Fei-Fei [13] tried to use an integrative model for scene relationship modeling. The SCSPM algorithm used the sparse coding method to reduce local feature quantization loss during the traditional nearest neighbor assignment process. Note that [9, 13] did not randomly choose training and testing images while [17] and the proposed SSR used randomization for reliable performance comparison. We can see the proposed SSR outperforms OB [9] by 7 percent and SCSPM [17] by 0.8 percent. This is because we use the supervised information for better semantic representation of images. Besides, the use of sparse coding along with max pooling can reduce the information loss during the feature quantization process which also helps to improve the performance of scene classification.

3.2. MIT Indoor Scene Dataset

This dataset has over 67 indoor scenes of 15620 images. The MIT Indoor scene dataset is more challenging than the UIUC-Sports dataset and the Scene-15 dataset both for more image categories and inter and intra class variation. Some indoor scenes are characterized by global spatial information while others by the objects within the scenes. We follow the experimental setting as in [14] and choose 80 images from each class for training and 20 images for testing. This process is repeated for five times.

We give the classification rates in Table 2 of the proposed method and methods in [9, 14] for the MIT Indoor scene dataset. Quattoni and Torralba [14] proposed a prototype based model which jointly combines both the local and global discriminative information. We can see that the proposed SSR again outperforms the OB method. Since the MIT Indoor scene dataset has large inter and intra class variations, the combination of training data for better representation is very useful for final classification.

3.3. Scene-15 Dataset

Table 2. Performance comparison on the MIT-Indoor Scene dataset. OB: object bank; PbM: prototype based model.

Methods	Performance
OB [9]	37.6
PbM [13]	26.0
SSR	38.2 \pm 0.7

This dataset has a number of 15 natural scene classes. We follow [2] and use 100 images in each class for classifier training and use the rest of images for testing. We repeat this process for ten times, as did in [2, 16].

The performance comparison is given in Table 3 for the Scene-15 dataset. We can see from Table 3 that the proposed SSR achieved the state-of-the-art method. That is because we combine the discriminative power of training data with the semantic image representation in a unified way and use sparsity for efficient feature selection hence improves the classification results.

4. CONCLUSION

This paper proposed a scene classification method by supervised semantic representation of images. Each image is represented as a response histogram which is generated by a combination of the prediction of pre-trained generic object classifiers and classifiers generated by supervised learning. Sparsity is also used to choose the most discriminative features for efficient scene classification. Experimental results on the UIUC-Sports dataset, the MIT Indoor scene dataset and the Scene-15 dataset show the effectiveness of the proposed method.

Our future work will concentrate on how to use the training images more efficiently and choose the most discriminative semantic space for image representation.

5. ACKNOWLEDGEMENT

This work was supported by the National Natural Science Foundation of China 973 Program (Project No. 2010CB327905) and the National Natural Science Foundation of China (Grant No. 60903146, 60835002)

6. REFERENCES

[1] J. S.ivic and A. Zisserman. "Video Google: A text retrieval approach to object matching in videos. In *Proc. ICCV*, volume 2, pages 1470-1477, 2003.

[2] S. Lazebnik, C. Schmid, J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. CVPR*, 2006.

[3] J. Gemert, C. Veenman, A. Smedulders and J. Geusebroek. Visual word ambiguity. In *IEEE Transactions on PAMI*, 2010.

Table 3. Performance comparison on the Scene-15 dataset. SPM: spatial pyramid matching; KC: kernel codebook; OB: object bank; SCSPM: sparse coding spatial pyramid matching.

Methods	Performance
SPM [2]	81.40 \pm 0.50
KC [3]	76.70 \pm 0.40
OB [9]	80.9
SPM [16]	76.73 \pm 0.65
SCSPM [16]	80.28 \pm 0.93
SSR	81.91 \pm 0.55

[4] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*, McGraw-Hill, 1983.

[5] D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. In *The Journal of Machine Learning Research*, 2003.

[6] T. Hofmann. Probabilistic latent semantic indexing. In *Proc. ACM SIGIR*, pages 50-57, 2006.

[7] N. Rasiwasia and N. Vasconcelos. Scene classification with low-dimensional semantic spaces and weak supervision. In *Proc. CVPR*, 2008.

[8] G. Carneiro, A. Chan, P. Morena, and N. Vasconcelos. Supervised learning of semantic classes for image annotation and retrieval. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Vol. 29, No. 3, 2007.

[9] L. Li, H. Su, E. Xing, and Li Fei-Fei. Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *Proc. ECCV*, 2010.

[10] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. In *JAIR*, 29, 2007.

[11] D. Hoiem, A. Efros, and M. Hebert. Automatic photo pop-up. In *Proc. SIGGRAPH*, 24(3):577-584, 2005.

[12] J. Wright, A. Yang, A. Ganesh, S. Satry, and Y. Ma. Robust face recognition via sparse representation. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, February, 2009.

[13] L. J. Li and L. Fei-Fei. What, where and who? Classifying events by scene and object recognition. In *Proc. ICCV*, Rio de Janeiro, Brazil, October 14-20, 2007.

[14] A. Quattoni and A. Torralba. Recognizing indoor scenes. In *Proc. CVPR*, 2009.

[15] D. G. Lowe. Distinctive image features from scale-invariant keypoints. In *Proceeding of ECCV Workshop on Statistical Learning in Computer Vision*, 60(2):91-110, 2004.

[16] J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *Proc. CVPR*, 2009.

[17] Local features are not lonely-Laplacian sparse coding for image classification. In *Proc. CVPR*, 2010.