

Image Classification Using Probability Higher-order Local Auto-Correlations

Tetsu Matsukawa¹ and Takio Kurita²

¹ University of Tsukuba,
1-1-1 Tennodai, Tsukuba, Japan

² National Institute of Advanced Industrial Science and Technology,
1-1-1 Umezono, Tsukuba, Japan
{t.matsukawa, takio-kurita}@aist.go.jp

Abstract. In this paper, we propose a novel method for generic object recognition by using higher-order local auto-correlations on probability images. The proposed method is an extension of bag-of-features approach to posterior probability images. Standard bag-of-features is approximately thought as sum of posterior probabilities on probability images, and spatial co-occurrences of posterior probability are not utilized. Thus, its descriptive ability is limited. However, using local auto-correlations of probability images, the proposed method extracts richer information than the standard bag-of-features. Experimental results show the proposed method is enable to have higher classification performances than the standard bag-of-features.

1 Introduction

Genetic object recognition technologies are important for automatic image search. Despite many methods have been researched until now, the performance is still inferior to human recognition system.

The most popular approach for generic object recognition is bag-of-features [3], because of its simplicity and effectiveness. Bag-of-features is originally inspired from text recognition method “bag-of-words”, and uses orderless collection of quantized local features. The main steps of bag-of-features are : 1) Detection and description of image patches. 2) Assigning patch descriptors to a set of pre-determined codebooks with a vector quantization algorithm. 3) Constructing a bag of features, which counts the number of patches assigned to each codebook. 4) Applying a classifier by treating the bag of features as the features vector, and thus determine which category to assign to the image.

It is known that the bag of features method is robust for background clutter, pose changes, intra-class variations and produces good classification accuracy. However, several problems are existed for applying to image representation. To solve these problems, many methods are proposed. Some of these methods are spatial pyramid binning to utilize location informations [7], higher level codebook creation based on local co-occurrence of codebooks [1][13][18], improvement of codebook creation[9][10][11] and region of interest based matching [14].

In this paper, we present a novel improvement of bag-of-features. The main novelty of the proposed method is to utilize probability images for feature extraction. The standard bag-of-features is approximately thought as a method so called “sum of posterior probabilities” on probability images. So the method does not utilize local co-occurrence on probability images. We applied higher-order local auto-correlations on probability images, thus richer information of probability images can be extracted. We call this image representation method as “Probability Higher-order Local Auto-correlations (PHLAC)”. PHLAC has desirable property for recognition, namely shift-invariance, additivity and synonymy [19] invariance. We show this image representation method PHLAC has the significantly better classification performance than the standard bag-of-features.

The proposed method gives the different direction of improvement to the currently proposed methods of bag-of-features (e.g. Correlation of codebooks, improvement of clustering and spatial pyramid binning), so this method can be combined with those methods in the future.

2 Related Work

The image feature extraction using local co-occurrence is recognized as an important concept [6] for recognition. Recently, several methods have been proposed using correlation. These are categorized to feature level co-occurrence and codebook level co-occurrence. The examples of feature level co-occurrence are local self similarity [12] and GLAC [5]. We can use these features in the codebook creation process, then the codebook level co-occurrence and the feature level co-occurrence is thought as another concept. The examples of codebook level co-occurrence are correlations [13] and Visual Phrases [18]. When using codebook level co-occurrence, we need large number of dimensions, e.g. in proportion to codebook size \times codebook size when we consider only co-occurrence of two codebooks. Thus, features selection method or dimension reduction method is necessary and current researches are focused on how to mining frequent and distinctive codebook sets [17][18][19]. The expressions of co-occurrence using a generative model have also been proposed [1] [16]. But, these methods require a complex latent model and expensive parameter estimations. On the other hand, our method can be easy implemented and is relatively low dimension but effective for classifications, because it is based on auto-correlations on posterior probability images. The methods which give posterior probability to a codebook have also been proposed [15][14], but these methods are not using auto-correlation of codebooks.

3 Probability High-order Local Auto-Correlations

3.1 Probability Images

Let I be an image region and $\mathbf{r} = (x, y)^t$ be a position vector in I . The image patch whose center is \mathbf{r}_k are quantized to M codebook $\{V_1, \dots, V_M\}$ by local feature extraction and vector quantization algorithm $VQ(\mathbf{r}_k) \in \{1, \dots, M\}$. These

steps are same as the standard bag-of-features [7]. Posterior probability $P(c|V_m)$ of category $c \in \{1, \dots, C\}$ is assigned to each codebook V_m using image patches on training images. Several forms of estimating posterior probability can be taken. (a) Codebook plausibility. The posterior probability is estimated by Bayes' theorem as follows.

$$P(c|V_m) = \frac{P(V_m|c)P(c)}{P(V_m)}, \quad (1)$$

where, $P(c) = 1/C$, $P(V_m) = (\# \text{ of } V_m) / (\# \text{ of all patches})$, $P(V_m|c) = (\# \text{ of class } c \wedge V_m) / (\# \text{ of class } c \text{ patches})$. Here, $P(c)$ is common constant, so set to 1.

(b) Codebook uncertainty. In our method, the probability is not restricted to the theoretical definition of probability. The pseudo probability which indicates the degree of supporting to each category from a codebook is considered. Codebook uncertainty is the percentage of class c in given codebook. This is defined as follows.

$$P(c|V_m) = \frac{P(V_m|c)P(c)}{\sum_{c=1}^C P(V_m|c)}. \quad (2)$$

(c) SVM weight. The weight of each codebook when learning by one-against-all linear SVM [4] is used to define pseudo probability. Assume we use K local image patches from one image, then the histogram of bag-of-features $\mathbf{H} = (H(1), \dots, H(M))$ becomes as follows.

$$H(m) = \sum_{k=1}^K \begin{cases} 1 & \text{if } (VQ(\mathbf{x}_k) = m) \\ 0 & \text{otherwise} \end{cases}. \quad (3)$$

Using the histogram of bag-of-features, the classification function of one-against-all linear SVM becomes as follows.

$$\arg \max_{c \in C} \{f_c(\mathbf{H}) = \sum_{m=1}^M \alpha_{c,m} H(m) + b_c\}, \quad (4)$$

where, $\alpha_{c,m}$ is the weight for each histogram bins and b_c is the learned threshold. We transform the weight of each histogram to non-negative by $\alpha_{c,m} \leftarrow \alpha_{c,m} - \min\{\alpha_c\}$ and normalize it by $\alpha_{c,m} \leftarrow \frac{\alpha_{c,m}}{\sum_{m=1}^M \alpha_{c,m}}$. Then we can obtain the pseudo probability by SVM weight as follows.

$$P(c|V_m) = \frac{\alpha_{c,m} - \min\{\alpha_c\}}{\sum_{m=1}^M (\alpha_{c,m} - \min\{\alpha_c\})}. \quad (5)$$

We used SVM weight as pseudo probability because the proposed method becomes a complete extension of the standard bag-of-features when using this pseudo probability (Sec.3.3).

In this paper, we assume to use grid sampling of local features [7] per p pixel interval, because of simplicity. We denote the set of sample points as I_p and we call the map of (pseudo) posterior probability of codebook of each local

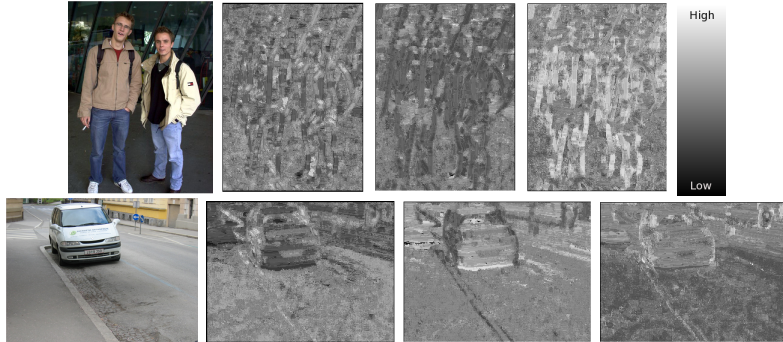


Fig. 1. Probability images (codebook plausibility): Original image, probability of BIKE(left), probability of CAR(middle), probability of PEOPLE(right). This probability image is calculated by 2 pixel interval ($p=2$), for easy understanding the original images are resized to the same size to probability images. The actual size of the original images are larger than the probability images by $p \times p$ pixels. Local features and codebook are the same as those used in experiment.

regions as a probability image. Examples of probability images are shown in Fig. 1. White color shows the high probability. The data are comes from IG02 used in the following experiment. The number of categories is 3 (BIKE, CAR and PEOPLE). It is noticed the human-like contours are appeared in PEOPLE probability.

3.2 PHLAC

We call HLAC features [6] on this probability images as PHLAC. The definition of Nth order PHLAC is as follows.

$$R(c, \mathbf{a}_1, \dots, \mathbf{a}_N) = \int_{I_p} P(c|V_{VQ}(\mathbf{r}))P(c|V_{VQ}(\mathbf{r} + \mathbf{a}_1)) \cdots P(c|V_{VQ}(\mathbf{r} + \mathbf{a}_N))d\mathbf{r}. \quad (6)$$

In practice, Eq.(6) can take so many forms by varying the parameters N and \mathbf{a}_n . In this paper, these are restricted to the following subset: $N \in \{0, 1, 2\}$ and $a_{nx}, a_{ny} \in \{\pm \Delta r \times p, 0\}$. By eliminating duplicates which arise from shifts, the mask patterns of PHLAC becomes as shown in Fig. 2. This mask pattern is the same as 35 HLAC mask patterns [6]. Thus, PHLAC inherits the desirable properties of HLAC for object recognition, namely shift-invariance and additivity. Although PHLAC does not have scale-invariance, we can deal with scale changes by using several size of mask patterns.

By calculating correlations in local regions, PHLAC becomes to robust against small spatial difference and noise. There are several alternatives of preprocessing of these local regions such as {max, average, median}. We found average is the

Algorithm1. PHLAC computation

Training Image :

- 1) Create codebook by local features and clustering algorithm (e.g. SIFT + K means).
- 2) Configure posterior probability of each codebook {plausibility, uncertainly, SVM}.

Training and Test Image :

- 3) Create C posterior probability images by p pixel interval.
 - 4) Preprocessing posterior probability images (local averaging).
 - 5) Calculate HLAC on posterior probability images by sliding HLAC mask patterns.
-

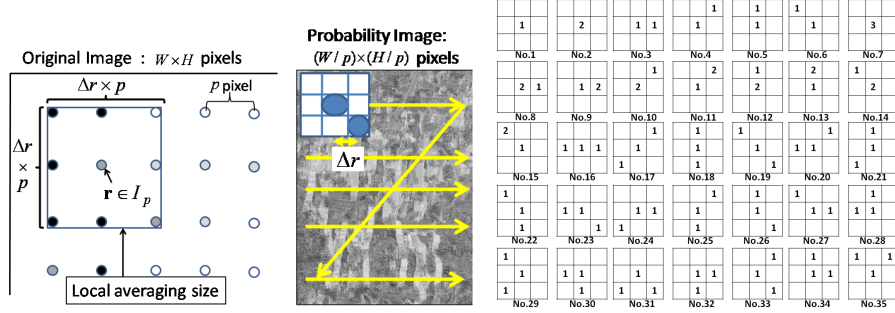


Fig. 2. PHLAC: local averaging size(left), extracting process(middle) and mask patterns(right).The number $\{1,2,3\}$ of mask patterns show the frequency for which its pixel value is used for product in Eq.(6).

best for other settings. Thus the practical formulation of PHLAC is given by

$$\begin{aligned}
 \mathbf{0}^{\text{th}} \text{ order} \quad R_{N=0}(c) &= \sum_{\mathbf{r} \in I_p} L_a(P(c|V_{VQ}(\mathbf{r}))) & (7) \\
 \mathbf{1}^{\text{st}} \text{ order} \quad R_{N=1}(c, \mathbf{a}_1) &= \sum_{\mathbf{r} \in I_p} L_a(P(c|V_{VQ}(\mathbf{r})))L_a(P(c|V_{VQ}(\mathbf{r} + \mathbf{a}_1))) \\
 \mathbf{2}^{\text{nd}} \text{ order} \quad R_{N=2}(c, \mathbf{a}_1, \mathbf{a}_2) &= \sum_{\mathbf{r} \in I_p} L_a(P(c|V_{VQ}(\mathbf{r})))L_a(P(c|V_{VQ}(\mathbf{r} + \mathbf{a}_1))) \\
 &\quad L_a(P(c|V_{VQ}(\mathbf{r} + \mathbf{a}_2))),
 \end{aligned}$$

where L_a means local averaging on a $(\Delta r \times p) \times (\Delta r \times p)$ region centered on \mathbf{r} (Fig. 2). Actually, PHLAC are obtained by HLAC calculation on local averaged probability image (see Algorithm.1.). PHLAC are extracted from probability images of all categories, thus the total number of features of PHLAC becomes $35 \times C$. There are two possibilities of classification using PHLAC image representations. One is the classification using all PHLAC of all categories (PHLAC_{ALL}) and the other is using one categories PHLAC for each one-against-all classifiers (PHLAC_{CLASSWISE}). We compare these methods in the following experiments.

3.3 Interpretation of PHLAC

Bag-of-features(0th) + local auto-correlations(1st + 2nd) : If we use SVM weights as pseudo probabilities, then 0-th order of PHLAC becomes the same as the classification by the standard bag-of-features using linear-SVM. Because \mathbf{H} is a histogram (see Eq.(3)), Eq.(4) is rewritten as follows.

$$\arg \max_{c \in C} \left\{ \sum_{k=1}^K \alpha_{c, VQ(r_k)} + b_c \right\} \quad (8)$$

$$= \arg \max_{c \in C} \left\{ \sum_{k=1}^K (\alpha_{c, VQ(r_k)} - \min\{\alpha_c\}) + K \min\{\alpha_c\} + b_c \right\} \quad (9)$$

$$= \arg \max_{c \in C} \{A_c R_{N=0}(c) + B_c\}, \quad (10)$$

where $A_c = \sum_{m=1}^M (\alpha_{c,m} - \min\{\alpha_c\})$, $B_c = K \min\{\alpha_c\} + b_c$. (In this transformation from Eq.(9) to Eq.(10), the relationship $R_{N=0}(c) = \sum_{k=1}^K \frac{\alpha_{c, VQ(r_k)} - \min\{\alpha_c\}}{A_c}$ is used.) This equation shows that the classification by the standard bag-of-features is possible by using only 0-th order of PHLAC and the learned parameters A_c and B_c . (Exactly, this was assumed no-preprocessing in the calculation of PHLAC). This is the case that SVM weight is used as pseudo probability, but it is expected other probabilities have also similar property. Because the histogram of the standard bag-of-feature is created by not utilizing local co-occurrences, the 0th order of PHLAC is thought as almost the one-against-all bag-of-features classifications. Higher order features of PHLAC have richer information of probability images (e.g. the shape of local probability distributions). Thus, if any commonly existed patterns are contained in the specific classes, this representation can be expected to achieve better classification performance than the standard bag-of-features.

The relationship of the standard bag-of-features and PHLAC classification is shown in Fig.3. In our PHLAC classification, we train additional classifier using 0th order PHLAC $\{R_{N=0}(1), \dots, R_{N=0}(C)\}$ and higher order PHLAC as feature vector. Thus, the only 0-th order PHLAC_{SVM} can achieve better performance than the standard bag-of-features.

Synonymy invariance : The synonymous codebooks are the codebooks which have similar posterior probabilities [18]. PHLAC calculates directly on the probability images, the same features can be extracted even a local appearance is exchanged to other appearances whose posterior probabilities are same. This synonymy invariance is important for creating compact image representations [19].

4 Experiment

We compared the classification performances of the standard bag-of-features and PHLAC using two commonly used image datasets: IG02[8] and fifteen natural scene categories [7].

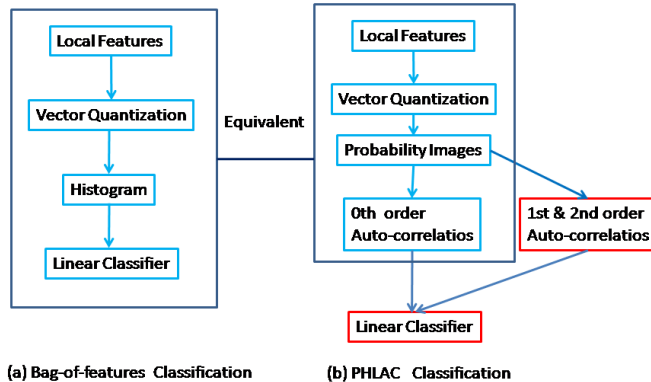


Fig. 3. Schematic comparison of the standard bag-of-features classification with our proposed PHLAC classification.

To obtain reliable results, we repeated the experiments 10 times. Ten random subsets were selected from the data to create 10 pairs of training and test data. For each of these pairs a codebook was created by using k-means clustering on training set. For classification, a linear SVM was used by one-against-all. As implementation of SVM, we used LIBSVM. Five-fold cross-validation on the training set was used to tune parameters of SVM. The classification rate we report is the average of the per-class recognition rates which in turn are averaged over the 10 random test sets.

As local features, we used a SIFT descriptor [2] sampled on a regular grid. The modification by the dominant orientation was not used and computed on 16×16 pixel patch sampled every 8 pixels ($p = 8$). In the codebook creation process, all features sampled every 16 pixel on all training images were used for k-means clustering. As normalization method, we used L2-norm normalization for both the standard bag-of-features and PHLAC. In PHLAC, the features were L2 normalized by each auto-correlations order. Below we denote the classification of PHLAC using probability by codebook plausibility as PHLAC_{Plau} , PHLAC using pseudo probability by codebook uncertainty as PHLAC_{Unc} and SVM weight as PHLAC_{SVM} . Note that although the SVM of standard bag-of-features is used for Eq. (5) of PHLAC_{SVM} , the result of 0th order PHLAC_{SVM} is different from the result of standard bag-of-features from the reason mentioned in Sec 3.3.

4.1 Result of IG02

At first, we used IG02 [8](INRIA Annotations for Granz-02) dataset which contains large variations of target size. The classification task is to classify the test images to 3 categories, CAR, BIKE and PEOPLE. The number of training im-

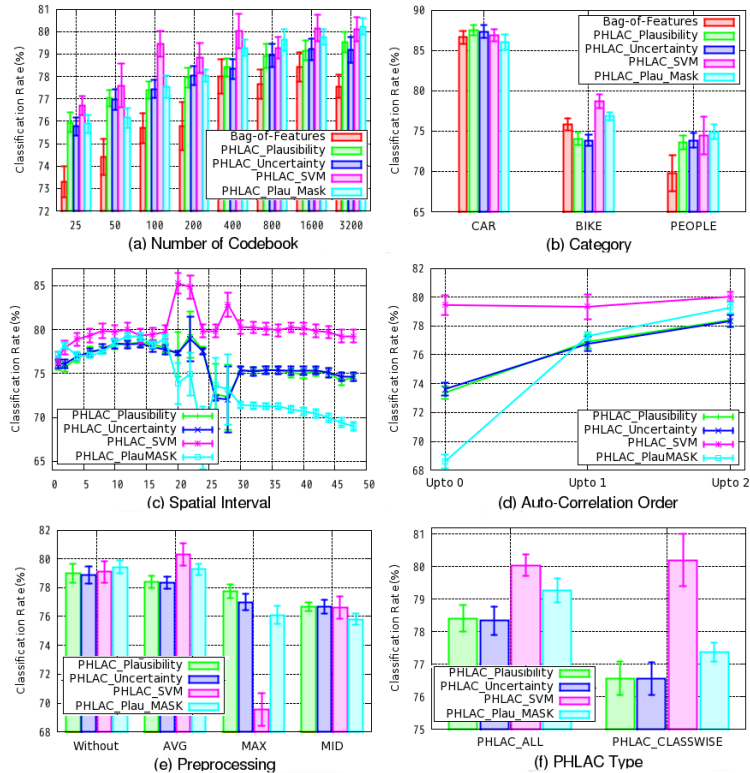


Fig. 4. Recognition rates of IG02. The basic setting is codebook size = 400 ((b)-(f)), Spatial Interval $\Delta r=12$ ((a),(b),(d)-(f)) and PHLAC_{ALL} ((a)-(e)).

ages of each category is 162 for CAR, 177 for BIKE and 140 for PEOPLE. The number of test images is same as training images. We resampled 10 sets of training and test sets from all images. Image size is 640×480 pixels or 480×640 pixels. Maraszalek et al prepared mask images which indicates target object locations. We also attempted to estimate probability of Eq.(1) by using only target object region’s local features. We denote this PHLAC features as PHLAC_{Plau-MASK}. The experimental results are shown in Fig. 4.

Overall performance: As basic settings we used spatial interval $\Delta r=12$ and PHLAC_{ALL}. In all codebook size, all types of PHLAC achieves higher classification performances than the standard bag-of-features (Fig.4(a)). PHLAC_{SVM} achieves higher classification rates than PHLAC_{Plau} and PHLAC_{Unc}. By using mask images for estimating probability, the performance of PHLAC_{Plau} becomes better when the codebook size is larger than 400.

Recognition rates per category: The classification rates of PHLAC becomes higher than the standard bag-of-features almost all cases (Fig.4(b)). Especially,

the classification rates of PEOPLE are higher than the standard bag-of-features in any settings of PHLAC. This is because human-like contours which are shown in Fig.1 are appeared in human’s regions and not existed in other images.

Spatial interval: The spatial interval seems to be better near $\Delta r=12$ ($12 \times 8 = 96$ pixel) in all settings except for PHLAC_{SVM} (Fig.4(c)). The classification rates of PHLAC_{Plau} and PHLAC_{Unc} become lower as to increase the spatial interval. In the case of PHLAC_{SVM} , classification rates is still high when the spatial interval becomes large and the peak of classification rates is appeared near $\Delta r=20$. But the classification rates in $\Delta r=20$, PHLAC_{Plau} and PHLAC_{Unc} become to be low, so we set the spatial interval as to $\Delta r=12$ as basic settings. In practice, multi-scale spatial interval is more useful than single spatial interval, because there are several optimal spatial intervals.

Auto-correlation order: In the case of PHLAC_{Plau} and PHLAC_{Unc} , the classification rates become higher as to increase auto-correlation order (Fig.4(d)). PHLAC_{SVM} is higher classification performance than other PHLAC only 0-th order auto-correlations. This is the reason of high classification rates of PHLAC_{SVM} in the large spatial intervals. Using up to 2nd order auto-correlations, PHLAC_{SVM} also can achieve the best classification performance. Especially in the optimal spatial interval of $\text{PHLAC}_{SVM}(\Delta r=20)$, the 2nd order auto-correlation of PHLAC_{SVM} were 5.01% better than 0th order (Fig.4(c)).

Preprocessing: In local averaging and no preprocessing seems to be comparable in Fig.4(e). But when we tried another codebook size and spatial intervals, the local averaging were often outperformed no preprocessing cases. Thus, we recommend to using local averaging for preprocessing.

PHLAC type: PHLAC_{ALL} are better performance than $\text{PHLAC}_{CLASSWISE}$ in PHLAC_{Plau} and PHLAC_{Unc} (Fig.4(f)). On the other hand PHLAC_{SVM} are better in the case of using $\text{PHLAC}_{CLASSWISE}$. This indicates the dimension for training of each SVM can be reduced to 35 dimension when using PHLAC_{SVM} .

4.2 Result of scene-15

Next we performed experiments on Scene-15 dataset [7]. The Scene-15 dataset consists of 4485 images spread over 15 categories. The fifteen categories contain 200 to 400 images each and range from natural scene like mountains and forest to man-made environments like kitchens and office. We selected 100 random images per categories as a training set and the remaining images as the test set. We used PHLAC_{ALL} and experimentally set spatial interval as to $\Delta r = 8$. Some examples of dataset images and probability images are shown in Fig.5. Recognition rates of scene 15 are shown in Fig.6. In Scene-15, PHLAC achieves higher recognition performances than the standard bag-of-features classification in all categories and all number of codebook. In this dataset, PHLAC_{Plau} and PHLAC_{Unc} indicates higher accuracy than PHLAC_{SVM} . In the case of codebook size is 200, PHLAC_{Plau} gives more than 15% higher recognition rate.

In our experimental settings, classification rates of the standard bag-of-features using histogram intersection kernel [7] is $66.31(\pm 0.15)\%$ in codebook size 200 and PHLAC_{Plau} achieves $69.48(\pm 0.27)\%$ by using linear SVM. While

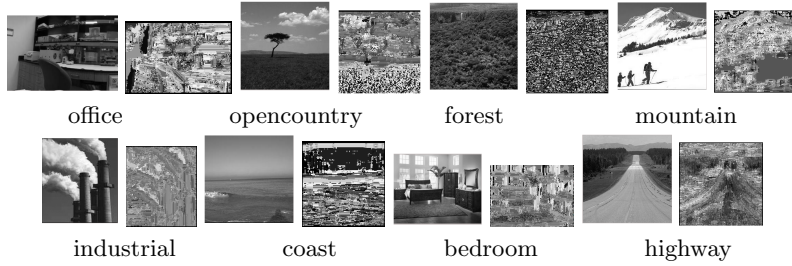


Fig. 5. Example of Scene15. Probability image shows probabilities of own category.

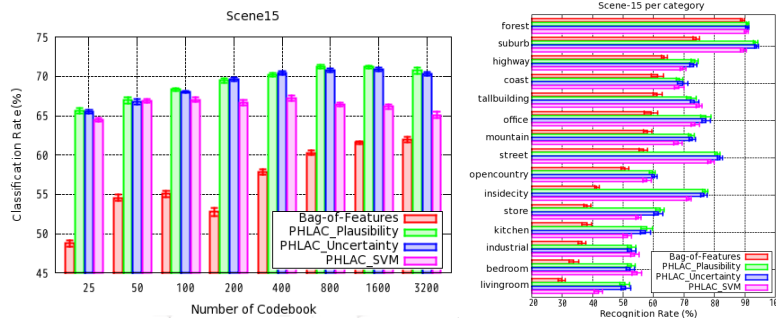


Fig. 6. Recognition Rates of Scene15: per codebook size (left) and per category when codebook size is 200(right).

Lazabnik reported $72.2(\pm 0.6)$ % on the standard bag-of-features, this difference is caused by the difference of implementations such as feature extractions and codebook creations. The proposed method and the standard bag-of-features use the same codebook and features through in our experiments.

5 Conclusion

In this paper, we proposed an image description method using higher-order local auto-correlations on probability images called "Probability Higher-order Auto Correlations(PHLAC)". This method is regarded as an extension of the standard bag-of-features for improving the limitation of spatial information by utilizing co-occurrence of local spatial pattern in posterior probabilities. This method has shift-invariance and additivity as in HLAC [6]. Experimental results show the proposed method achieved higher classification performance than the standard bag-of-features in average 2 % and 15 % in the case of IG02 and Fifteen Scene Dataset respectively using 200 codebooks. We think combinations with other method (e.g. spatial binning and correlation features) probably improve the performance by the proposed probability auto-correlations scheme.

References

1. Agarwal, A., Triggs, B.: Multilevel Image Coding with Hyperfeatures. *International Journal of Computer Vision*. **78** (2008) 15–27
2. Lowe, D., G.: Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*. **60** (2004) 91–110
3. Csurka, G., Dance, C., R., Fan, L., Willamowski, J., Bray, C.: Visual Categorization with Bag of Keypoints. *European conference on computer vision 2004 workshop on Statistical Learning in Computer Vision*. (2004) 59-74.
4. Vapnik, V.: *Statistical Learning Theory*. John Wiley & Sons, 1998.
5. Kobayashi, T., Otsu, N.: Image Feature Extraction Using Gradient Local Auto-Correlations. *European conference on computer vision, 2008, Part I, LNCS 5302*, pp.346-358 (2008).
6. Otsu, N., Kurita, T.: A new scheme for practical flexible and intelligent vision systems. *IAPR Workshop on Computer Vision* (1988).
7. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories, In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp.2169-2178 (2006).
8. Marszalek, M., Schmid, C.: Spatial Weighting for Bag-of-Features, In: *IEEE Conference on Computer Vision and Pattern Recognition, Vol.2*, pp.2118-2125 (2006).
9. Jurie, F., Triggs, B.: Creating efficient codebooks for visual recognition, In: *IEEE International Conference on Computer Vision, Vol.1*, pp.604-610 (2005).
10. Nowak, E., Jurie, F., Triggs, B.: Sampling strategies for bag-of-features image classification. *European conference on computer vision, (2006)*.
11. Gemert, J.C.V., Geusebroek, J.-M., Veenman, C.J.: Kernel Codebooks for Scene Classification. *European conference on computer vision, 2008, Part III, LNCS 5304*, pp.696-709 (2008).
12. Shechtman, E., Irani, M.: Matching local self-similarities across images and videos. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp.511-518 (2007).
13. Snavely, S., Winn, J., Criminisi, A.: Discriminative Object Class Models of Appearance and Shape by Correlations. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp.2033-2040 (2006).
14. Bosch, A., Zisserman, A., Munoz, X.: Image classification using random forests and ferns. In: *IEEE International Conference on Computer Vision*, pp.1-8 (2007).
15. Shotton, J., Johnson, M., Cipolla, R.: Semantic texton forests for image categorization and segmentation, In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp.1-8 (2008).
16. Wang, X., Grimson, E.: Spatial Latent Dirichlet Allocation. In: *Proceedings of Neural Information Processing Systems Conference (NIPS) 2007*.
17. Quack, T., Ferrari, V., Leibe, B., and L. Van-Gool: Efficient mining of frequent and distinctive feature configurations. In: *IEEE International Conference on Computer Vision, (2007)*.
18. Yuan, J., Wu, Y., Yang, M.: Discovery of Collocation Patterns: from Visual Words to Visual Phrases. In: *IEEE Conference on Computer Vision and Pattern Recognition, (2007)*.
19. Zheng, Y.-T., Zhao, M., Neo, S.-Y., Chua, T.-S., Tian, Q.: Visual Synset: Towards a Higher-level Visual Representation. In: *IEEE Conference on Computer Vision and Pattern Recognition, (2008)*.