



## Nearest-Neighbor based Metric Functions for indoor scene recognition

Fatih Cakir, Uğur Güdükbay\*, Özgür Ulusoy

Bilkent University, Department of Computer Engineering, 06800 Bilkent, Ankara, Turkey

### ARTICLE INFO

#### Article history:

Received 7 February 2011

Accepted 29 July 2011

Available online 5 August 2011

#### Keywords:

Scene classification

Indoor scene recognition

Nearest Neighbor classifier

Bag-of-visual words

### ABSTRACT

Indoor scene recognition is a challenging problem in the classical scene recognition domain due to the severe intra-class variations and inter-class similarities of man-made indoor structures. State-of-the-art scene recognition techniques such as capturing holistic representations of an image demonstrate low performance on indoor scenes. Other methods that introduce intermediate steps such as identifying objects and associating them with scenes have the handicap of successfully localizing and recognizing the objects in a highly cluttered and sophisticated environment.

We propose a classification method that can handle such difficulties of the problem domain by employing a metric function based on the Nearest-Neighbor classification procedure using the bag-of-visual words scheme, the so-called *codebooks*. Considering the codebook construction as a Voronoi tessellation of the feature space, we have observed that, given an image, a learned weighted distance of the extracted feature vectors to the center of the Voronoi cells gives a strong indication of the image's category. Our method outperforms state-of-the-art approaches on an indoor scene recognition benchmark and achieves competitive results on a general scene dataset, using a single type of descriptor.

© 2011 Elsevier Inc. All rights reserved.

### 1. Introduction

Scene classification is an active research area in the computer vision community. Many classification methods have been proposed that aim to solve different aspects of the problem such as topological localization, indoor–outdoor classification and scene categorization [1–9]. In scene categorization the problem is to associate a semantic label to a scene image. Although categorization methods address the problem of categorizing any type of a scene, they usually only perform well on outdoors [10]. In contrast, classifying indoor images has remained a further challenging task due to the more difficult nature of the problem. The intra-class variations and inter-class similarities of indoor scenes are the biggest barriers for many recognition algorithms to achieve satisfactory performance on images that have never been seen, i.e., test data. Moreover, recognizing indoor scenes is very important for many fields. For example, in robotics, the perceptual capability of a robot for identifying its surroundings is a highly crucial ability.

Earlier works on scene recognition are based on extracting low-level features of the image such as color, texture and shape properties [1,3,5]. Such simple global descriptors are not powerful enough to perform well on large datasets with sophisticated environmental settings. Olivia and Torralba [4] introduce a more compact and robust global descriptor, the so-called *gist*, which

captures the holistic representation of an image using spectral analysis. Their descriptor performs well on categorizing outdoor images such as forests, mountains and suburban environments but has difficulties recognizing indoor scenes.

Borrowing ideas from the human perceptual system, recent work on indoor scene recognition focuses on classifying images by using representations of both global and local image properties and integrating intermediate steps such as object detection [10,11]. This is not surprising since indoor scenes are usually characterized by the objects they contain. Consequently, indoor scene recognition can be mainly considered as a problem of first identifying objects and then classifying the scene accordingly. Intuitively, this idea seems reasonable but it is unlikely that even state-of-the-art object recognition methods [12–14], can successfully localize and identify unknown number of objects in cluttered and sophisticated indoor images. Hence, classifying a particular scene via objects becomes yet a more challenging issue.

A solution to this problem is to classify an indoor image by implicitly modeling objects with densely sampled local cues. These cues will then give indirect evidence of a presence of an object. Although this solution seems contrary to the methodology of recognizing indoor scenes by the human visual system, i.e., explicitly identifying objects and associating them with scenes, it provides a successful alternative by bypassing the drawbacks of trying to localize objects in highly intricate environments.

The most successful and popular descriptor that captures the crucial information of an image region is the Scale-Invariant Feature Transform (SIFT) [15,16]. This proposes the idea that SIFT-like

\* Corresponding author.

E-mail addresses: [fcakir@cs.bilkent.edu.tr](mailto:fcakir@cs.bilkent.edu.tr) (F. Cakir), [gudukbay@cs.bilkent.edu.tr](mailto:gudukbay@cs.bilkent.edu.tr) (U. Güdükbay), [oulusoy@cs.bilkent.edu.tr](mailto:oulusoy@cs.bilkent.edu.tr) (Ö. Ulusoy).

features extracted from images of a certain class may have more similarities in some manner than those extracted from images of irrelevant classes. This similarity measure can be achieved by first defining a set of categorical words (the so-called *visual words*) for each class and then using a learned metric function to measure the distance between local cues and these visual words.

Therefore, we introduce a novel non-parametric weighted metric function with a spatial extension based on the approach described in [17]. In their work, Bolman et al. show that a Nearest-Neighbor (NN) based classifier which computes direct image-to-class distances without any quantization step achieves performance rates among the top leading learning-based classifiers. We show that a NN-based classifier is also well suited for categorizing indoor scenes because: (i) It incorporates image-to-class distances which is extremely crucial for classes with high variability. (ii) Considering the insufficient performance of state-of-the-art recognition algorithms on a large object dataset [12], it successfully allows classifying indoor scenes directly from local cues without incorporating any intermediate steps such as categorizing via objects. (iii) Given a query image, it allows ranked results and thus can be employed for a preprocessing step to successfully narrow down the size of possible categories for subsequent analyses.

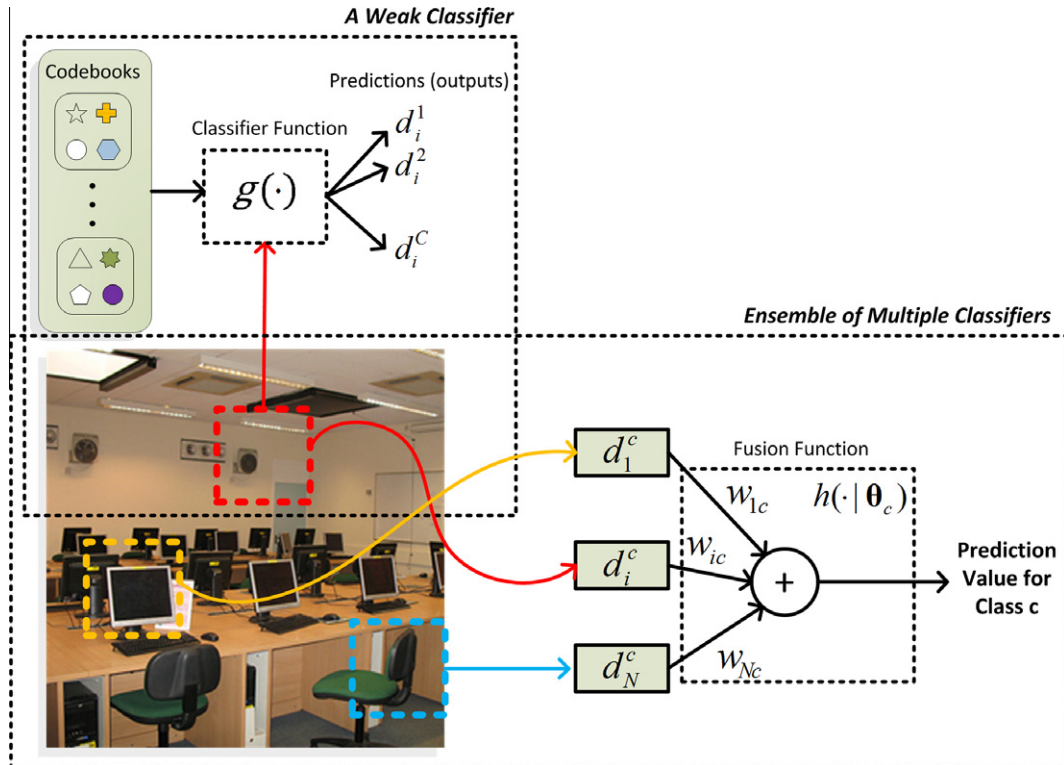
Bolman et al. also show that a descriptor quantization step, i.e., codebook generation, severely degrades the performance of the classifier by causing information loss in the feature space. They argue that a non-parametric method such as the Nearest-Neighbor classifier has no training phase as the learning-based methods do to compensate for this loss of information. They evaluate their approach on Caltech101 [18] and Caltech256 datasets [19], where each image contains only one object and maintains a common position, and on the Graz-01 dataset [20], which has three classes (bikes, persons and a background class) with a basic class vs. no-class classification task. On the other hand, for a multi-category recognition task of scenes where multiple objects co-exist in a highly cluttered, varied and complicated form, we observe that

our NN-based classifier with a descriptor quantization step outperforms the state-of-the-art learning-based methods. The additional quantization step allows us to incorporate spatial information of the quantized vectors, and more importantly, it significantly reduces the performance gap between our method and other learning-based approaches. It is computationally inefficient for a straightforward NN-based method without a quantization step to perform classification, considering the datasets with large number of training images.

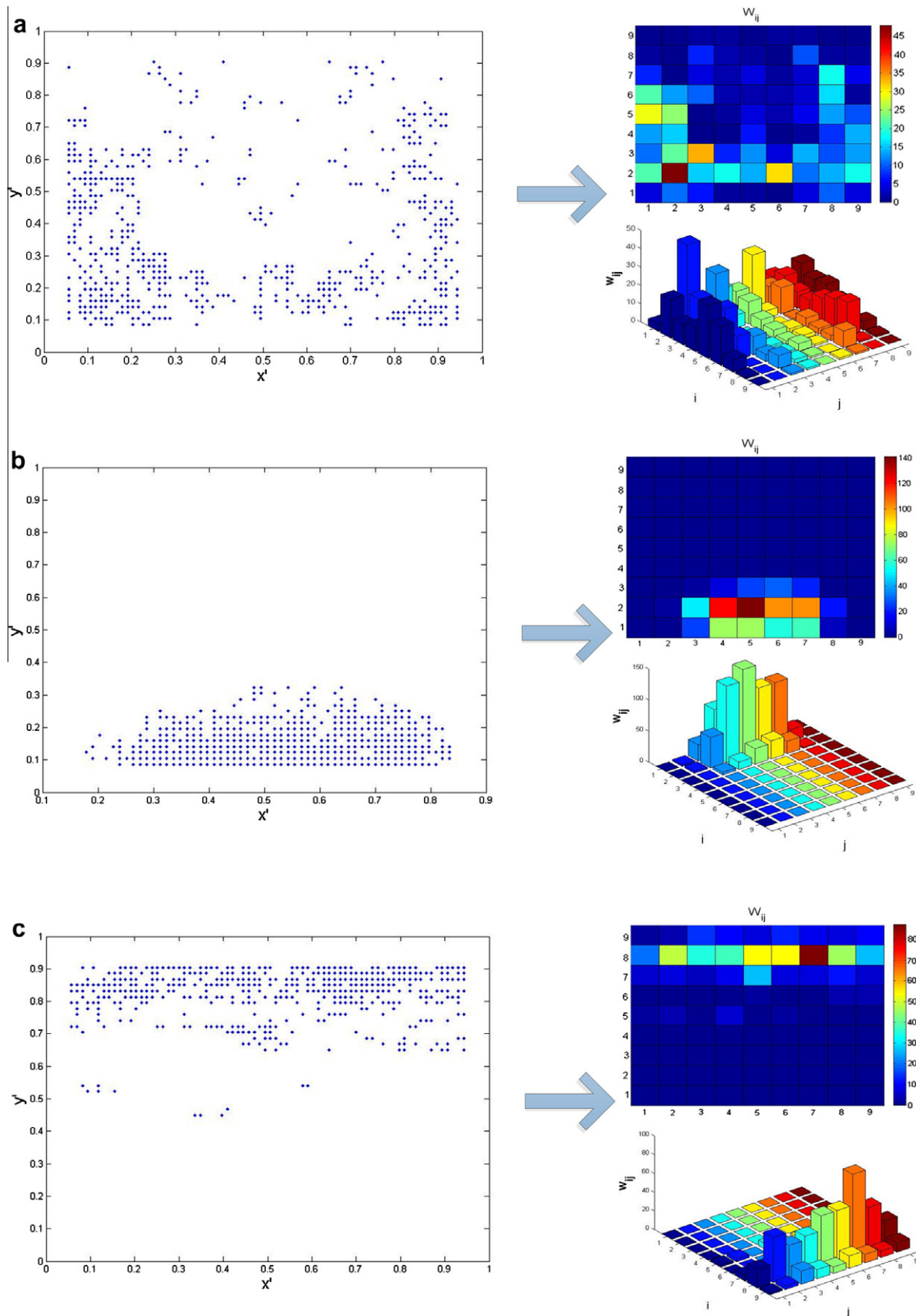
The rest of this paper is organized as follows: Section 2 discusses related work. In Section 3 we describe the framework of our proposed method. We present experimental results and evaluate the performance in Section 4. Section 5 gives conclusions and future work.

## 2. Related work

Earlier works on scene classification are based on extracting low-level features of the image such as color, texture and shape properties. Szummer and Picard [1] use such features to determine whether an image is an outdoor or an indoor scene. Vailaya et al. [3] use color and edge properties for the city vs. landscape classification problem. Ulrich and Nourbakhsh [5] employ color-based histograms for mobile robot localization. Such simple global features are not discriminative enough to perform well on a difficult classification problem, such as recognizing scene images. To overcome this limitation, Olivia and Torralba [4] introduce the gist descriptor, a technique that attempts to categorize scenes by capturing its spatial structure properties, such as the degree of openness, roughness, naturalness, using spectral analysis. Although a significant improvement over earlier basic descriptors, it has been shown in [10] that this technique performs poorly in recognizing indoor images. One other popular descriptor is SIFT [16]. Due to its strong discriminative power even under severe image transformations,



**Fig. 1.** The Nearest-Neighbor based metric function as an ensemble of multiple classifiers based on the local cues of a query image. Each local cue can be considered as a weak classifier that outputs a numeric prediction value for each class. The combination of these predictions can then be used to classify the image.



**Fig. 2.** Spatial layouts and weight matrix calculation for three different visual words. The left sides of (a), (b) and (c) represent the spatial layouts of the visual words that themselves represent the relative positions of the extracted descriptors to their image boundaries. These layouts are then geometrically partitioned into  $M \times M$  bins and a weight matrix  $W$  is computed as shown on the right sides of (a), (b) and (c).

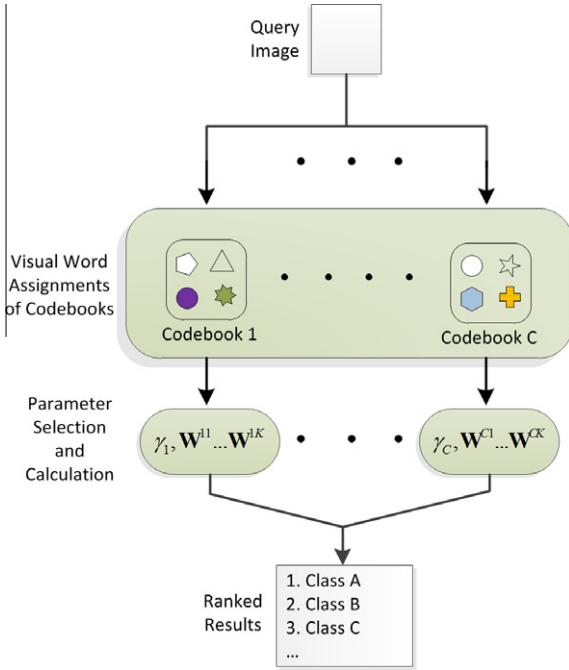


Fig. 3. Flow chart of the testing phase of our method.

noise and illumination changes, it has been the most preferred visual descriptor in many scene recognition algorithms [6,7,21–23].

Such local descriptors have been successfully used with the bag-of-visual words scheme for constructing codebooks. This concept has been proven to provide good results in scene categorization [23]. Fei-Fei and Perona [22] represent each category with such a codebook and classify scene images using Bayesian hierarchical models. Lazebnik et al. [7] use the same concept with spatial extensions. They hierarchically divide an image into sub-regions, which they call the spatial pyramid, and compute histograms based on quantized SIFT vectors over these regions. A histogram intersection kernel is then used to compute a matching score for each quantized vector. The final spatial pyramid kernel is implemented as concatenating weighted histograms of all features at all sub-regions. The traditional bag-of-visual words scheme discards any spatial information; hence many methods utilizing this concept also introduce different spatial extensions [7,24].

Bosch et al. [25] present a review of the most common scene recognition methods. However, recognizing indoor scenes is a more challenging task than recognizing outdoor scenes, owing to severe intra-class variations and inter-class similarities of man-made in-

door structures. Consequently, this task has been investigated separately within the general scene classification problem. Quattoni and Torralba [10] brought attention to this issue by introducing a large indoor scene dataset consisting of 67 categories. They argue that together with the global structure of a scene which they capture via the gist descriptor, the presences of certain objects described by local features are strong indications of its category. Espinace et al. [11] suggest using objects as an intermediate step for classifying a scene. Such approaches are coherent with the human vision system since we identify and characterize scenes by the objects they contain. However, with the state-of-the-art object recognition methods [12–14,26], it is very unlikely to successfully identify multiple objects in a cluttered and sophisticated environmental setting. Instead of explicitly modeling the objects, we can use local cues as indirect evidence for their presence and thus bypass the drawbacks of having to successfully recognize them, which is a very difficult problem considering the intricate nature of indoor scenes.

### 3. Nearest-Neighbor based Metric Functions (NNbMF)

#### 3.1. Baseline problem formulation

The popular bag-of-visual words paradigm introduced in [27] has become commonplace in various image analysis tasks. It has been proven to provide powerful image representations for image classification and object/scene detection. To summarize the procedure, consider  $\mathbf{X}$  to be a set of feature descriptors in  $D$ -dimensional space, i.e.,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L]^T \in \mathbb{R}^{L \times D}$ . A vector quantization or a codebook formation step involves the Voronoi tessellation of the feature space by applying  $K$ -means clustering to set  $\mathbf{X}$  to minimize the cost function

$$J = \sum_{i=1}^K \sum_{l=1}^L \|\mathbf{x}_l - \mathbf{v}_i\|^2 \quad (1)$$

where the vectors in  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K]^T$  correspond to the centers of the Voronoi cells, i.e., the visual words of codebook  $\mathbf{V}$ , and  $\|\cdot\|$  denotes the  $L_2$ -norm.

After forming a codebook for each class using Eq. (1), a set  $\mathbf{X}_q = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$  denoting the extracted feature descriptors from a query image can be categorized to class  $c$  by employing the Nearest-Neighbor classification function  $y: \mathbb{R}^{N \times D} \rightarrow \{1, \dots, C\}$  given as

$$y(\mathbf{X}_q) = \underset{c=1, \dots, C}{\operatorname{argmin}} \left[ \underbrace{\sum_{n=1}^N \|\mathbf{x}_n - \operatorname{NN}_c(\mathbf{x}_n)\|}_{h(\cdot|\theta_c)} \right] \quad (2)$$

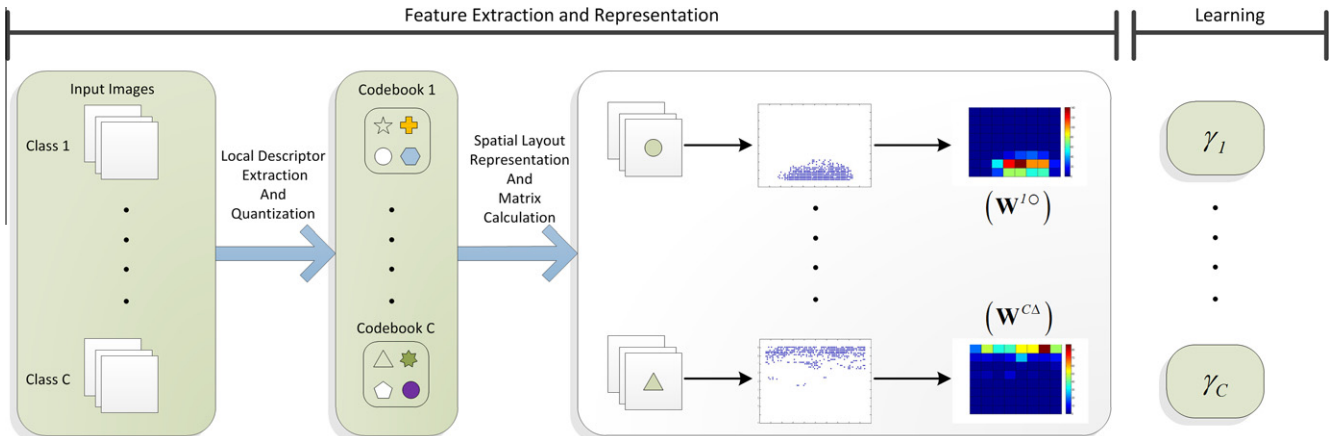


Fig. 4. The flow chart for the training phase of our method.



where  $NN_c(\mathbf{x})$  denotes the nearest visual word of  $\mathbf{x}$ , i.e., the nearest Voronoi cell center, in the Voronoi diagram of class  $c$ ,  $y_i \in \{1, \dots, C\}$  refers to class labels and  $h(\cdot | \theta_c)$  denotes a combination function with the parameter vector  $\theta_c$  associated with class  $c$ . Intuitively, Eq. (2) can be considered as an ensemble of multiple experts based on the extracted descriptor set  $\mathbf{X}_q$ . In this ensemble learning scheme there are  $|\mathbf{X}_q|$  weak-classifiers and  $h: \mathbb{R}^N \rightarrow \mathbb{R}$  is a fusion function to combine the outputs of such experts. This large ensemble scheme is very suitable for the particular problem domain where each scene object, implicitly modeled by local cues, provides little discriminative power in the classification objective but in combination they significantly increase the predictive performance.

From this perspective, given a query image, assume  $N$  base-classifiers corresponding to the extracted descriptor set  $\mathbf{X}_q = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ . Let  $\mathbf{V}_c = [\mathbf{v}_{c1}, \mathbf{v}_{c2}, \dots, \mathbf{v}_{cN}]^T$  and  $d_i^c$  be the codebook and the prediction of base classifier  $g(\mathbf{x}_i, \mathbf{V}_c) = \|\mathbf{x}_i - NN_c(\mathbf{x}_i)\|$  for class  $c$ , respectively. Taking  $d_1^c = g(\mathbf{x}_i, \mathbf{V}_c)$ , the final prediction value for the particular class is then

$$h(d_1^c, d_2^c, \dots, d_N^c | \theta_c) = \sum_{n=1}^N \omega_{nc} d_n^c \quad (3)$$

where  $\theta_c = [\omega_{1c}, \dots, \omega_{Nc}]^T$  denote the parameters of the fusion function associated with class  $c$ . Note that  $\theta_c = 1, \forall c \in \{1, \dots, C\}$  in Eq. (2). In the next section, we will use spatial information of the extracted descriptors to determine the parameter vector set  $\theta = \{\theta_1, \dots, \theta_C\}$ . Fig. 1 illustrates this concept. It should be noted that Eq. (2) does not take into account unquantized descriptors, as in [17]. There is a trade-off between information loss and computational efficiency because of the quantization of the feature space.

### 3.2. Incorporating spatial information

The classic bag-of-visual words approach does not take into account spatial information and thus loses crucial data about the distribution of the feature descriptors within an image. Hence, this is an important aspect to consider when working to achieve satisfactory results in a classification framework. We incorporate spatial information as follows. Given extracted descriptors in  $D$ -dimensional space,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L]^T \in \mathbb{R}^{L \times D}$  and their spatial locations  $\mathbf{S} = [(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)]$ , during the codebook generation step we also calculate their relative position with respect to the corresponding image boundaries from which they are extracted. Hence their relative locations are  $\mathbf{S}' = [(x'_1, y'_1), (x'_2, y'_2), \dots, (x'_L, y'_L)] = \left[ \left( \frac{x_1}{w_1}, \frac{y_1}{h_1} \right), \left( \frac{x_2}{w_2}, \frac{y_2}{h_2} \right), \dots, \left( \frac{x_L}{w_L}, \frac{y_L}{h_L} \right) \right]$ , where the  $(w_1, h_1), (w_2, h_2), \dots, (w_L, h_L)$  pairs represent the width and height values of the corresponding images. After applying clustering to the set  $\mathbf{X}$ , we obtain the visual word set  $\mathbf{V}$  as described in the previous section. Since similar feature descriptors of  $\mathbf{X}$  are expected to be assigned to the same visual word, their corresponding coordinate values described in set  $\mathbf{S}'$  should have similar values. Fig. 2 shows the spatial layout of the descriptors assigned to several visual words.

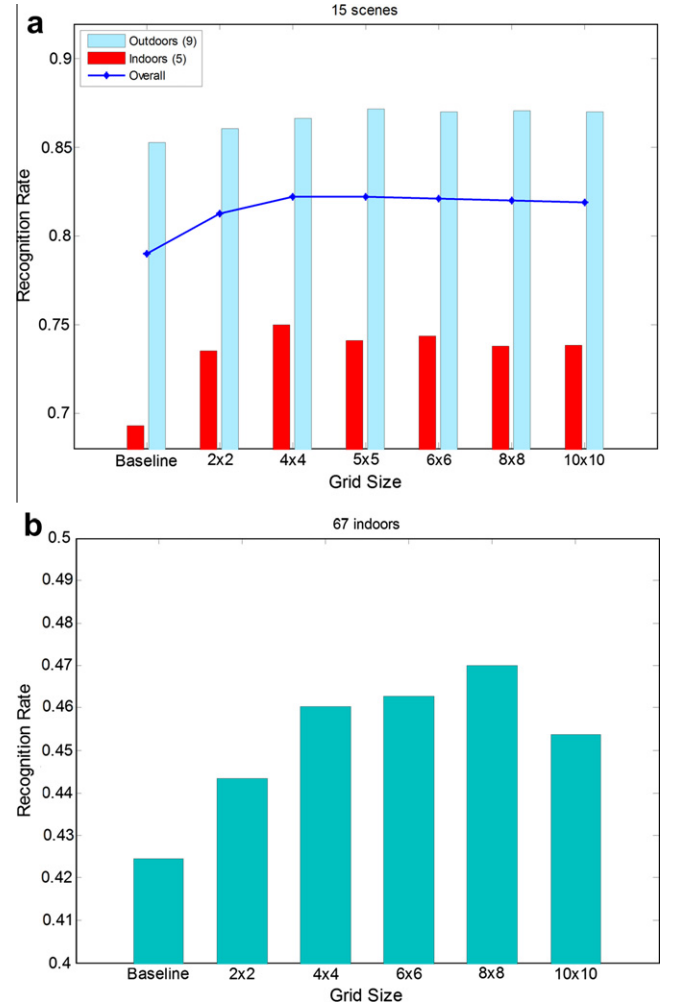
To incorporate this information into Eq. (2), we consider the density estimation methods which are generally used for determining unknown probabilistic density functions. It should be noted that we do not consider a probabilistic model; thus obtaining and using a legitimate density function is irrelevant in our case.

**Table 1**

Performance comparison with different  $\gamma$  settings.

	Baseline	$\gamma_{LP} \in \mathbb{R}^C$	$\gamma_{QP} \in \mathbb{R}^C$	$\gamma_{LP} \in \mathbb{R}$	Baseline <sub>full</sub>	$\gamma_m = \gamma_{LP}$	$\gamma_m^*$
15-Scenes	78.93	79.60	79.83	81.17	78.99	81.04	<b>82.08</b>
67-Indoor scenes	40.75	35.15	35.15	43.13	42.46	45.22	<b>47.01</b>

C refers to the number of categories in a dataset and Baseline refers to the method when Eq. (2) is used. Subscripts LP and QP stand for linear and quadratic programming, respectively. They refer to the optimization model with different  $n$  settings in Eq. (8).  $\gamma_m$  refers to the manual selection of the scale parameter.



**Fig. 5.** Recognition rates based on different grid size settings.

We can assign weights for each grid on the spatial layout of every visual word using a histogram counting technique (cf. Fig. 2). Suppose we geometrically partition this spatial layout into  $M \times M$  grids. Then for the  $j$ th visual word of class  $c$ ,  $\mathbf{v}_{cj}$ , the weight of a grid can be calculated as

$$\mathbf{W}^{cf} = [w_{ij}^{cf}] = \frac{k}{N} \quad (4)$$

where  $k$  is the number of descriptors assigned to  $\mathbf{v}_{cj}$  that fall into that particular grid and  $N$  is the total number of descriptors assigned to  $\mathbf{v}_{cj}$ . During the classification of a query image, the indices  $i, j$  correspond to the respective grid location of an extracted feature descriptor. An alternative way for defining weights is to first consider  $\mathbf{W}^{cf} = [w_{ij}^{cf}] = k$  then scale this matrix as

$$\mathbf{W}^{cf'} = \frac{[w_{ij}^{cf}]}{\max(\mathbf{W}^{cf})} \quad (5)$$

**Table 2**  
Performance comparison with state-of-the-art methods.

Methods	Descriptor	67 indoor scenes classification rate	15 scenes classification rate
Morioka et al. [26]	SIFT( $D = 36$ )	$39.63 \pm 0.69$	$83.40 \pm 0.58$
Quattoni and Torralba [10]	SIFT( $D = 128$ ) GIST( $D = 384$ )	$\sim 28$	–
Zhou et al. [29]	PCA-SIFT( $D = 64$ )	–	<b>85.20</b>
Yang et al. [13]	SIFT( $D = 128$ )	–	$80.28 \pm 0.93$
Lazebnik et al. [7]	SIFT( $D = 128$ )	–	$81.40 \pm 0.50$
NNbMF	SIFT (2 scales, $D = 256$ )	<b>47.01</b>	82.08

$D$  refers to the dimension of the descriptor.

where  $\max(\cdot)$  describes the largest element. Eq. (5) does not provide weight consistency of the visual words throughout a codebook. It assigns larger weights to visual words that have a sparse distribution in the spatial layout while attenuating the weights of the visual words that are more spatially compact. The choice of a weight matrix assignment is directly related to the problem domain; as we have found Eq. (4) more suitable for the 67-indoor benchmark and Eq. (5) suitable for the 15-scenes benchmark.

We calculate the weight matrices for all visual words of every codebook. The function  $h(\cdot | \theta_c)$  described in Eq. (2) now can be improved as

$$\sum_{n=1}^N \left(1 - \gamma_c \mathbf{W}_{ij}^{cf}\right) \times \|\mathbf{x}_n - NN_c(\mathbf{x}_n)\| \quad (6)$$

where  $NN_c(\mathbf{x}_n) \equiv \mathbf{v}_{cf}$ . The parameter set now includes the weight matrices associated with each visual word of a codebook, i.e.,  $\theta_c = [\mathbf{W}^c_1, \mathbf{W}^c_2, \dots, \mathbf{W}^c_K]$ . Obviously  $\gamma_c$  functions as a scale operator for a particular class, e.g., if  $\gamma_c = 0$  then the spatial location for class  $c$  is entirely omitted when classifying an image, i.e., only the sum of the descriptors' Euclidean distance to their closest visual words is considered.

This scale operator can be determined manually or by using an optimization model. Now, given codebook  $c$ , assume a vector  $\mathbf{d}^c \in \mathbb{R}^N$  that holds the predictions of every extracted descriptor  $\mathbf{x}_n$  of a query image as its elements; i.e.,  $d_n^c = g(\mathbf{x}_n, \mathbf{V}_c) = \|\mathbf{x}_n - NN_c(\mathbf{x}_n)\|$ , where  $n \in \{1, \dots, N\}$  corresponds to extracted descriptor indices and  $NN_c(\mathbf{x}_n)$  refers to the nearest visual word to  $\mathbf{x}^n$  ( $NN_c(\mathbf{x}_n) \equiv \mathbf{v}_{cf}$ ).  $\alpha_n^c$  denotes the corresponding spatial weights assigned to  $d_n^c$ ; i.e.,  $\alpha_n^c = \gamma_c \mathbf{W}_{ij}^{cf}$ . Referring to the vector of these spatial weights as  $\alpha^c \in \mathbb{R}^N$ , Eq. (6) can now be redefined as  $(\mathbf{1} - \alpha^c) \cdot \mathbf{d}^c$  and an image can be classified to class  $c$  by using the function

$$y(\mathbf{X}_q) = \underset{c=1, \dots, C}{\operatorname{argmin}} \left[ \frac{(\mathbf{1} - \alpha^c) \cdot \mathbf{d}^c}{h(\cdot | \theta_c)} \right] \quad (7)$$

Consider an image  $i$  that belongs to class  $j$  with an irrelevant class  $k$ . We would like to satisfy the inequalities  $(\mathbf{1} - \alpha_i^j)^T \mathbf{d}_i^j < (\mathbf{1} - \alpha_i^k)^T \mathbf{d}_i^k$ . Given  $i$  training images and  $j$  classes, we specify a set of  $S = i \times j \times (j - 1)$  inequality constraints where  $k = j - 1$ . Since we will not be able to find a scale vector that satisfies all such constraints, we introduce slack variables,  $\xi_{ijk}$ , and try to minimize the sum of slacks allowed. We also aim to select a scale vector  $\gamma$  so that Eq. (6) remains as close to Eq. (2) as possible. Hence we minimize the  $L_n$ -norm of  $\gamma$ . Consequently, finding the scale vector  $\gamma = [\gamma_1, \dots, \gamma_j]$  can now be modeled as an optimization problem as follows:

$$\begin{aligned} \min \quad & \|\gamma\|_n + \varphi \sum_{i,j,k} \xi_{ijk} \\ \text{subject to} \quad & \forall (i, j, k) \in S : \\ & (-\alpha_i^j)^T \mathbf{d}_i^j + (\alpha_i^k)^T \mathbf{d}_i^k < \mathbf{d}_i^k - \mathbf{d}_i^j + \xi_{ijk} \\ & \xi_{ijk} \geq 0, \gamma \geq \mathbf{0} \end{aligned} \quad (8)$$

where  $\varphi$  is a penalizing factor. We choose  $n$  from  $\{1, 2\}$ , resulting in linear and quadratic programming problems, respectively. One may prefer the  $L_2$ -norm, since sparsity is not desirable in our case due to the fact that sparse solutions may heavily bias categories associated with large scale weights. An alternative model is to define one weight value associated with all categories. This model is less flexible but it prevents a possible degradation in recognition performance caused by sparsity. The scale vector can also be manually chosen. Figs. 3 and 4 depict the testing and training phase of the proposed method, respectively.

## 4. Experimental setup and results

### 4.1. Training data and parameter selections

This section presents the training setup of our NN-based metric function on the 15 scenes [7] and 67 indoor scenes datasets [10]. The 15-scenes dataset contains 4485 images spread over 15 indoor and outdoor categories containing 200–400 images each. We use the same experimental setup as in [7] and randomly choose 100 images per class for training, i.e., for codebook generation and learning the scale vector  $\gamma$ , and use the remaining images for testing.

The 67-indoor scenes dataset contains images solely from indoor scenes with very high intra-class variations and inter-class similarities. We use the same experimental setup, as in [10] and [28]. Approximately 20 images per class are used for testing and 80 images per class for training.

We use two different scales of SIFT descriptors for evaluation. For the 15-scenes dataset, patches with bin sizes of 6 and 12 pixels are used, and for the 67-indoor scenes dataset, the bin sizes are selected as 8 and 16 pixels. The SIFT descriptors are sampled and concatenated at every four pixels and are constructed from  $4 \times 4$  grids with eight orientation bins (256 dimension in total). The training images are first resized to speed the computation and to provide scale consistency. The aspect ratio is maintained, but all images are scaled down so their largest resolution does not exceed 500 and 300 pixels and the feature space is clustered using  $K$ -means into 500 and 800 visual words, for the 67-indoor scenes and 15-scenes datasets, respectively. We use 100 K SIFT descriptors extracted from random patches to construct a codebook.

The spatial layout of each visual word from each category is geometrically partitioned into  $M \times M$  bins and a weight matrix is formed for each visual word from Eqs. (4) and (5). Several settings are used to determine the scale vector  $\gamma$ . We first consider assigning different weights to all categories ( $\gamma \in \mathbb{R}^C$ ). We find the optimal scale vector by setting  $n = \{1, 2\}$  in Eq. (8) and solving the corresponding optimization problem. We also use another setting for the optimization model where we assign the same weight to all categories ( $\gamma \in \mathbb{R}$ ). Alternatively, we select the scale parameter manually.

The constraints in Eq. (8) are formed as described in the previous section with 10 training images per class. The rest of the training set is used for codebook construction. The subset of the training images used for parameter learning is also employed as the validation set when manually tuning the scale parameter to find its optimal value. The value that yields the highest performance for this validation set is then selected for our method.

The performance rate is calculated by the ratio of correctly classified test images within each class. The final recognition rate is the total number of correctly classified images divided by the total number of test images used in the evaluation.

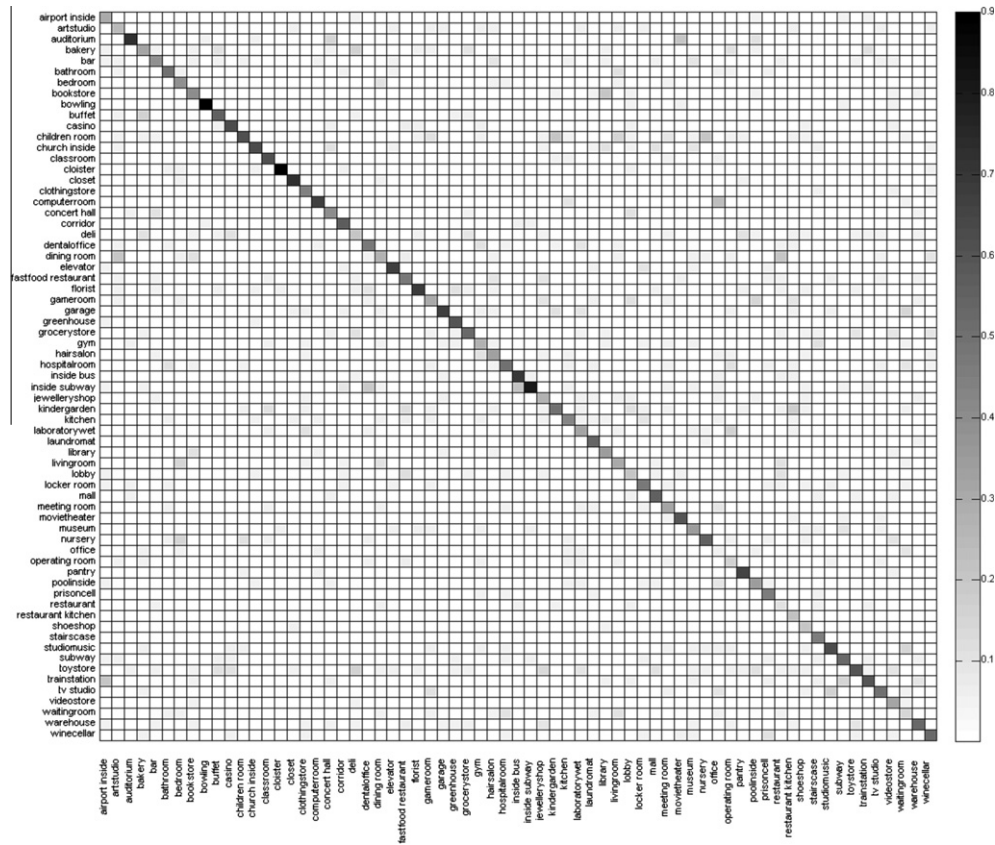


Fig. 6. Confusion matrix for the 67-indoor scenes dataset. The horizontal and vertical axes correspond to the true and predicted classes, respectively.

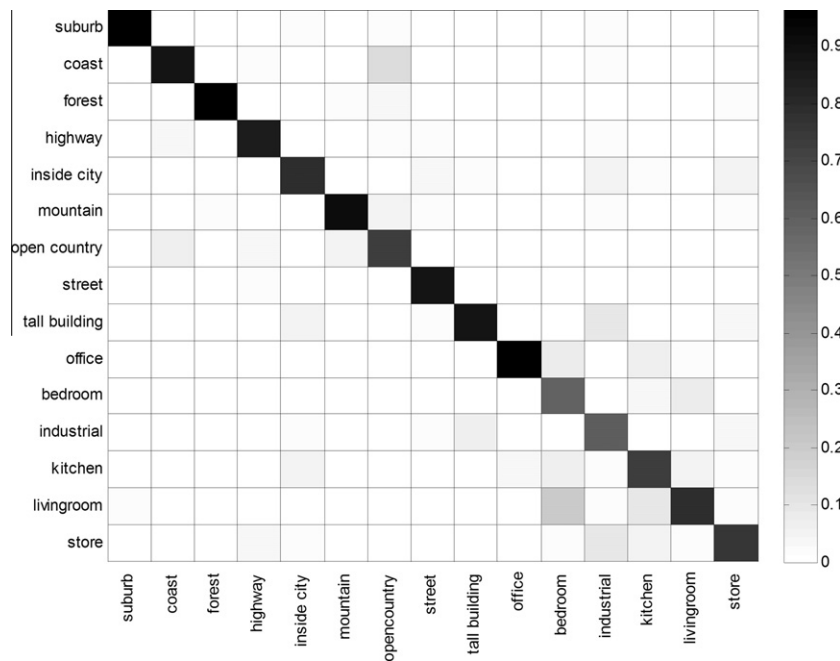


Fig. 7. Confusion matrix for the 15-scenes dataset. The columns and rows denote the true and predicted classes, respectively.

#### 4.2. Results and discussion

Table 1 shows recognition rates for both datasets with different scale vector settings. *Baseline* and *Baseline<sub>full</sub>* refer to the method when Eq. (2) is used (no spatial information is incorporated). The

difference is that *Baseline<sub>full</sub>* uses all available training images for codebook generation while leaves 10 images per class for scale parameter learning. In Table 1, the settings to the right of the baselines use the corresponding codebook setup. Observe the positive correlation between the number of training images used for

constructing codebooks and the general recognition rate. This impact is clearly visible on the 67-indoors dataset. When we generate codebooks using all available training data the recognition rate increases by 2%. The 15-scenes dataset has little intra-class variations with respect to the 67-indoors dataset, hence increasing the number of training images for codebooks generation yields only a slight increase in the performance.

The results where a scale parameter is assigned to every category ( $\gamma = [\gamma_1, \gamma_2, \dots, \gamma_c] \in \mathbb{R}^c$ ) are slightly better than the baseline implementation in the 15-scenes benchmark. In spite of an insignificant increase, we observe that setting  $n = 2$  in Eq. (8) gives a higher recognition rate compared to that with  $n = 1$ . This confirms our previous assertion that dense solutions increase the performance. This effect is clearly observed when we assign the same scaling parameter  $\gamma$  to all 15 categories. On the other hand, assigning a different scale parameter for each category in the 67-indoor scenes dataset decreases the performance values for both the LP and QP programming models. In fact we observed that the solutions to these models are identical for our setting. This situation can be avoided and the overall performance value can be increased by using more training images, however this results in the reduction of the number of available training images for codebook construction which also degrades the recognition rate.

Another solution is to assign the same scale parameter to all categories. This positively affects the performance, resulting in a 43% and 45% recognition rate with the two corresponding codebook setups when a LP optimization model is used to determine the scale parameter. One can easily expect that this effect will be much stronger in a problem domain where spatial distributions of the visual words are more ordered and compact. The last two columns in Table 1 shows the recognition rate when the scale parameter is manually tuned. As the initial selection for the parameter we used the value determined by the LP model. The performance rate of this initial selection is also included in Table 1 ( $\gamma_m = \gamma_{LP}$ ). The heuristic optimal value  $\gamma_m^*$  is then found by a simple numerical search.

Although the learned value of the scale parameter increases the accuracy of the method, manually tuning the parameter with respect to a validation set provides the highest accuracy in our setting. A more robust learning scheme can be constructed by introducing further constraints to the optimization model in Eq. (8).

Fig. 5 shows the recognition rates with different weight matrix ( $\mathbf{W}$ ) sizes. Geometrically partitioning the spatial layout into  $5 \times 5$  and  $8 \times 8$  grids yields the best results for the 15-scenes and 67-indoor scenes datasets, respectively. The 15-scenes dataset can be separated into five indoor and nine outdoor categories. We ignore the *industrial* category since it contains both indoor and outdoor images. Observe that incorporating spatial information improves the performance rate of the outdoor categories by 2% only. The performance rate for the indoor categories is improved by up to 6%. This difference can be explained by the more orderly form of the descriptors extracted from the indoor images. This improvement is  $\sim 4.5\%$  for the 67-indoor scenes dataset due to further difficulty and intra-class variations.

Table 2 compares our method with the state-of-the-art scene recognition algorithms. Our method achieves more than 7% improvement over the best published result in the 67-indoor benchmark [26] and shows competitive performance in the 15-scenes dataset. Figs. 6 and 7 show the confusion matrix for the 67 indoor scenes and 15 scenes datasets, respectively.

Our method also induces rankings that could naturally be used as a pre-processing step in another recognition algorithm. As shown in Fig. 8a and b, our method returns the correct category within the top ten results by ranking the categories for a query image with 82% overall accuracy in the 67-indoor scenes benchmark.

This rate is near 100% considering the returned top three results in the 15-scenes dataset (cf. Fig. 8a). Hence one can utilize this aspect of our algorithm to narrow down category choices, consequently increasing their final recognition rate by analyzing other information channels of the query image with different complementary descriptors or classification methods. Fig. 9 shows a set of classified images.

#### 4.3. Runtime performance

Compared to learning-based methods such as the popular Support Vector Machines (SVM), the Nearest-Neighbor classifier has a slow classification time, especially when the dataset is too large and the dimension of the feature vectors is too high. Several approximation techniques have been proposed to increase the efficiency of this method, such as [30] and [31]. These techniques involve pre-processing the search space using data structures, such as KD-trees or BD-trees. These trees are hierarchically structured so that only a subset of the data points in the search space is considered for a query point. We utilize the Approximate Nearest Neighbors library (ANN) [30]. For the 67 indoor scenes benchmark, it takes approximately 0.9 s to form a tree structure of a category codebook and about 2.0 s to search all query points of an image in a tree structure, using an Intel Centrino Duo 2.2 GHz CPU.

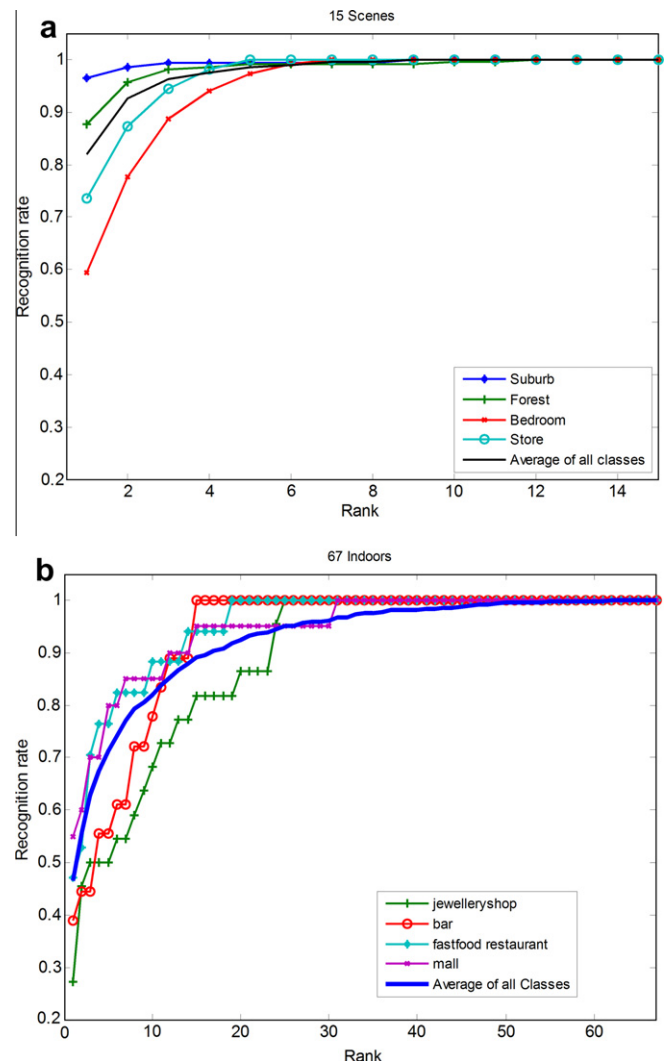
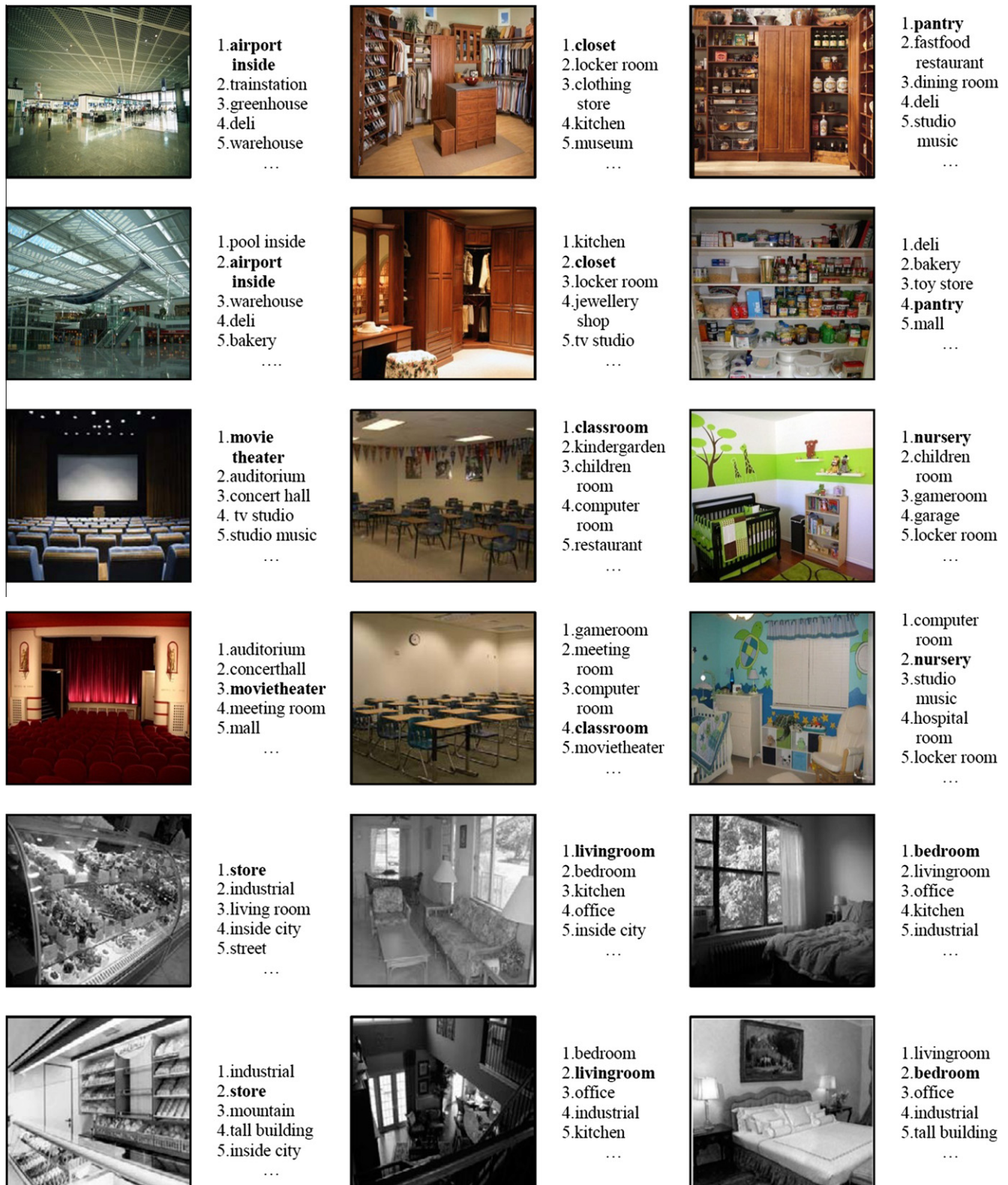


Fig. 8. Recognition rates based on rankings. Given a query image, if the true category is returned in the top- $k$  results, it is considered a correct classification.





**Fig. 9.** Classified images for a subset of indoor scene images. Images from the first four rows are taken from the 67-indoor scenes and the last two rows are from the indoor categories of the 15-scenes dataset. For every query image the list of ranked categories is shown on the right side. The bold name denotes the true category.

Without quantizing, it takes about 100 s to search all the query points. For the 15-scenes benchmark, it takes about 1.5 s to construct a search tree and 4.0 s to search all query points in it. Without quantizing, it takes approximately 200 s to search all the query points.

The CUDA implementation of the  $K$ -nearest neighbor method [32] further increases the efficiency by parallelizing the search process. We observed  $\sim 0.2$  s per class needed to search the query points extracted from an image using a NVIDIA Geforce 310 M graphics card.

## 5. Conclusion

We propose a simple, yet effective Nearest-Neighbor based metric function for recognizing indoor scene images. In addition, given an image our method also induces rankings of categories for a possible pre-processing step for further classification analyses. Our method also incorporates the spatial layout of the visual words formed by clustering the feature space. Experimental results show that the proposed method effectively classifies indoor scene images compared to state-of-the-art methods.

We are currently investigating how to further improve the spatial extension part of our method by using other estimation techniques to better capture and model the layout of the formed visual words. We are also investigating how to apply the proposed method to other problem domains such as auto-annotation of images.

## Acknowledgments

We thank to Muhammet Bastan for various discussions. We are grateful to Rana Nelson for proofreading and suggestions.

## References

- [1] M. Szummer, R.W. Picard, Indoor-outdoor image classification, in: Proceedings of the International Workshop on Content-based Access of Image and Video Databases (CAIVD '98), Washington, DC, USA, 1998, p. 42.
- [2] A. Torralba, K. Murphy, W. Freeman, M. Rubin, 2003. Context-based vision system for place and object recognition, in: Proceedings of the International Conference on Computer Vision.
- [3] A. Vailaya, A. Jain, H. Zhang, On image classification: city vs. landscapes, *Pattern Recogn.* 31 (1998) 1921–1935.
- [4] A. Oliva, A. Torralba, Modeling the shape of the scene: a holistic representation of the spatial envelope, *Int. J. Comput. Vis.* 42 (2001) 145–175.
- [5] I. Ulrich, I. Nourbakhsh, Appearance-based place recognition for topological localization, in: Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), 2000.
- [6] A. Bosch, A. Zisserman, X. Munoz, Scene classification using a hybrid generative/discriminative approach, *IEEE Trans. Pattern Anal. Mach. Intell.* 30 (4) (2008) 712–727.
- [7] S. Lazebnik, C. Schmid, J. Ponce, Beyond bags of features: spatial pyramid matching for recognizing natural scene categories, in: Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2006.
- [8] S. Se, D.G. Lowe, J.J. Little, Vision-based mobile robot localization and mapping using scale-invariant features, in: Proceedings of the International Conference on Robotics and Automation, 2001, pp. 2051–2058.
- [9] A. Pronobis, B. Caputo, P. Jensfelt, H.I. Christensen, A discriminative approach to robust visual place recognition, in: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems, 2006.
- [10] A. Quattoni, A. Torralba, Recognizing indoor scenes, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2009.
- [11] P. Espinace, T. Kollar, A. Soto, N. Roy, Indoor scene recognition through object detection, in: Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), 2010.
- [12] P. Gehler, S. Nowozin, On feature combination for multiclass object classification, in: Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2009, pp. 221–228.
- [13] J. Yang, K. Yu, Y. Gong, T. Huang, Linear spatial pyramid matching using sparse coding for image classification, in: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2009, pp. 1794–1801.
- [14] J. Wu, J.M. Rehg, Beyond the Euclidean distance: creating effective visual codebooks using the histogram intersection kernel, in: Proceedings of the Twelfth IEEE International Conference on Computer Vision (ICCV), 2009.
- [15] K. Mikolajczyk, C. Schmid, A performance evaluation of local descriptors, *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (10) (2005) 1615–1630.
- [16] D. Lowe, Distinctive image features from scale-invariant keypoints, *Int. J. Comput. Vis.* 60 (2) (2004) 91–110.
- [17] O. Bolman, E. Shechtman, M. Irani, In defense of Nearest-Neighbor based image classification, in: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2008, pp. 1–8.
- [18] L. Fei-Fei, R. Fergus, P. Perona, Learning generative visual models from few training examples: an incremental Bayesian approach tested on 101 object categories, in: Proceedings of the CVPR Workshop on Generative-Model Based Vision, 2004.
- [19] G. Grin, A. Holub, P. Perona, Caltech 256 Object Category Dataset, Technical Report, UCB/CSD-04-1366, California Institute of Technology, 2006.
- [20] A. Opelt, M. Fussenegger, A. Pinz, P. Auer, Weak hypotheses and boosting for generic object detection and recognition, in: Proceedings of the Eighth European Conference on Computer Vision (ECCV), Lecture Notes in Computer Science, vol. 2, 2004, pp. 71–84.
- [21] J. Vogel, B. Schiele, A semantic typicality measure for natural scene categorization, in: Proceedings of 26th Pattern Recognition Symposium (DAGM), Lecture Notes in Computer Science, vol. 3175, 2004, pp. 195–203.
- [22] L. Fei-Fei, P. Perona, A Bayesian hierarchical model for learning natural scene categories, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), vol. II, 2005, pp. 524–531.
- [23] P. Quelhas, F. Monay, J.-M. Odobez, D. Gatica-Perez, T. Tuytelaars, A thousand words in a scene, *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (9) (2007) 1575–1589.
- [24] V. Viitaniemi, J. Laaksonen, Spatial extensions to bag of visual words, in: Proceedings of the 8th ACM International Conference on Image and Video Retrieval (CVIR), 2009.
- [25] A. Bosch, X. Muñoz, R. Martí, A review: which is the best way to organize/classify images by content?, *Image Vis. Comput.* 25 (6) (2007) 778–791.
- [26] N. Morioka, S. Satoh, Building compact local pairwise codebook with joint feature space clustering, in: Proceedings of the 11th European Conference on Computer Vision (ECCV), 2010, pp. 692–705.
- [27] J. Sivic, A. Zisserman, Video Google: a text retrieval approach to object matching in videos, in: Proceedings of the Ninth International Conference on Computer Vision, 2003, pp. 1470–1478.
- [28] A. Torralba, Indoor Scene Recognition. <<http://web.mit.edu/torralba/www/indoor.html>> (accessed May 2011).
- [29] X. Zhou, N. Cui, Z. Li, F. Liang, T.S. Huang, Hierarchical Gaussianization for image classification, in: Proceedings of the International Conference on Computer Vision (ICCV), 2009.
- [30] D. Mount, S. Arya, ANN: A library for approximate nearest neighbor searching, in: Proceedings of the 2nd Annual Fall Workshop on Computational Geometry, 1997.
- [31] A. Andoni, P. Indyk, Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions, in: Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS '06), 2006, pp. 459–468.
- [32] V. Garcia, E. Debreuve, M. Barlaud, Fast k-nearest neighbor search using GPU, in: Proceedings of the CVPR Workshop on Computer Vision on GPU, 2008.