



Texture representations using subspace embeddings

Xiaodong Yang, YingLi Tian*

Department of Electrical Engineering, The City College of New York, CUNY, USA

ARTICLE INFO

Article history:

Received 2 August 2012

Available online 26 March 2013

Communicated by G. Borgefors

Keywords:

Texture representation

Texture classification

Subspace embedding

ABSTRACT

In this paper, we propose a texture representation framework to map local texture patches into a low-dimensional texture subspace. In natural texture images, textons are entangled with multiple factors, such as rotation, scaling, viewpoint variation, illumination change, and non-rigid surface deformation. Mapping local texture patches into a low-dimensional subspace can alleviate or eliminate these undesired variation factors resulting from both geometric and photometric transformations. We observe that texture representations based on subspace embeddings have strong resistance to image deformations, meanwhile, are more distinctive and more compact than traditional representations. We investigate both linear and non-linear embedding methods including Principle Component Analysis (PCA), Linear Discriminant Analysis (LDA), and Locality Preserving Projections (LPP) to compute the essential texture subspace. The experiments in the context of texture classification on benchmark datasets demonstrate that the proposed subspace embedding representations achieve the state-of-the-art results while with much fewer feature dimensions.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

The automated analysis of texture is widely applied in a number of real-world applications, e.g., image and video retrieval, object recognition and segmentation, and natural scene classification (Bouguila, 2012; Caputo et al., 2010; Mailing and Cernuschi-Frias, 1982; Nguyen et al., 2012; Yang et al., 2011). However, it is a challenging problem to represent texture images due to scaling changes, affine deformations, and lighting variations. A desired texture representation is thus supposed to handle both geometric and photometric variations. There has been extensive research in the literature on designs of robust texture representations. Early work for modeling texture includes filter banks (Randen and Husoy, 1999) and co-occurrence features (Haralick. Statistical et al., 1979). They mainly concentrate on global 2D transformations such as rotation and scaling. Most recent work further captures the effects of 3D transformations such as viewpoint change and non-rigid surface deformation. The representation methods based on fractal analysis (Varma and Garg, 2007; Xu et al., 2009) have also been proposed to model spatial distribution properties of textons with impressive recognition performance. In addition, most recent state-of-the-art results in texture recognition are obtained by using histograms of local image features as distributions of textons (Csurka et al., 2004; Lazebnik et al., 2005; Zhang et al.,

2007). Therefore, the effective computations of textons are crucial for robust texture representations.

It is common to define texture as a visual pattern with the repetition of a set of basic primitives named textons. Accordingly, a histogram or distribution of textons can be used as an effective representation of texture images. For nature textures, textons can be approximated by the prototypes from clustering local texture patches. However, natural texture images are generated from interaction of multiple factors related to rotation, scaling, lighting, viewpoint, and non-rigid surface deformation, as illustrated in Fig. 1. The multiple factor variations result in severe difficulties for accurately capturing the essential factor, i.e., textons. In this paper, we propose to employ both linear and non-linear embedding approaches to map normalized local texture patches into a texture subspace for analysis. In our framework, the low-dimensional structures hidden in high-dimensional texture observations correspond to the essential factor for texture representation. In this way, the unwanted variation modes resulting from geometric and photometric transformations can be reduced or removed from the essential factor.

The approaches of subspace embedding have been demonstrated effectiveness in mining meaningful low-dimensional structures hidden in original high-dimensional feature space (Roweis et al., 2000; Tenenbaum et al., 2000). They are based upon the biological observation that human brain extracts a manageably small amount of perceptually relevant features from high-dimensional sensory inputs (about 10^4 auditory nerve fibers or 10^6 optic nerve fibers) (Tenenbaum et al., 2000). On the other hand, it also has been explored to transfer the design of local image descriptors to

* Corresponding author. Tel.: +1 212 650 7046.

E-mail addresses: xyang02@ccny.cuny.edu (X. Yang), ytian@ccny.cuny.edu (Y. Tian).

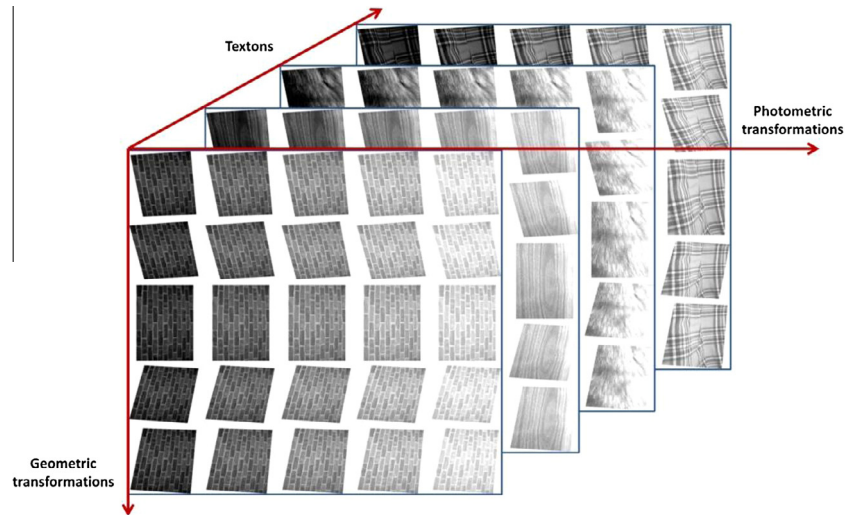


Fig. 1. Texture images are generated from interactions of multiple factors including textons, geometric transformations, and photometric transformations, among which, textons are the essential factor for texture representation.

a dimensionality reduction problem in the context of image matching (Hua et al., 2007; Ke and Sukthankar, 2004). Our proposed method of subspace embedding textons is mainly inspired by the appearance-based face recognition. As discussed in Belhumeur et al. (1997), He et al. (2005), Martinez and Kak (2001), Turk and Pentland (1991), face images varying in rotation, pose, illumination, and expression reside in a manifold of original data space. Mapping face images into a face subspace is able to conserve the essential factors of person identity but suppresses other factor variations. *Eigenfaces*, *Fisherfaces*, and *Laplacianfaces* are the state-of-the-art embedding algorithms in face recognition literature (Belhumeur et al., 1997; He et al., 2005; Turk and Pentland, 1991). *Eigenfaces* and *Fisherfaces* are linear methods which are used to effectively model the Euclidean structure of original feature space. *Laplacianfaces* is a non-linear approach that is able to preserve local data relationships and to discover the subspace of essential factor. Motivated by the success of subspace embedding methods in face recognition, in this paper we explore texture subspaces detected by PCA, LDA, and LPP, and then evaluate our approach in the context of texture classification. Following the conventions in face recognition, we name textons embedded by PCA, LDA, and LPP as *EigenTextons*, *FisherTextons*, and *LaplacianTextons*.

The remainder of this paper is organized as follows. Section 2 reviews existing approaches for texture representation and subspace embeddings. Section 3 describes PCA, LDA, and LPP methods. In Section 4, we provide the detailed procedures of representing texture images using the proposed methods. A variety of experimental results and discussions are presented in Section 5. Finally, Section 6 summarizes the remarks of this paper.

2. Related work

A major challenge of texture representation is to achieve invariance under a wide range of geometric and photometric variations. Early research work (Haralick. Statistical et al., 1979; Randen and Husoy, 1999) in this domain mainly focused on the analysis of global 2D image transformations including in-plane rotation and scaling. Because of lacking invariance to general geometric transformations, these approaches however cannot effectively model texture images with large 3D transformations such as viewpoint change and non-rigid surface deformation. Multi-fractal analysis

has recently been proposed and achieves good resilience to 3D deformations (Varma and Garg, 2007; Xu et al., 2006; Xu et al., 2009). Texture representations based on this method benefit from the invariance of fractal dimension to geometric transformations. For example, MFS proposed by Xu et al. (2006, 2009) combined fractal dimensions of pixel sets grouped by density functions and orientation templates.

In order to make texture representations more robust to 3D image transformations (e.g., viewpoint change and non-rigid surface deformation) as well as illumination variations, most of recent methods on texture representation rely on extracting local features by local image detectors and descriptors (Csurka et al., 2004; Lazebnik et al., 2005; Zhang et al., 2007). A textons dictionary is then generated by clustering the extracted local features. For example, Lazebnik et al. (2005) proposed a texture representation method based on affine-invariant detectors (Harris and Laplacian) and descriptors (RIFT and SPIN). Zhang et al. (2007) represented textures by combining multiple local image features (SIFT, RIFT, and SPIN). Extensive experiments in texture classification and retrieval have demonstrated that histogram of local image feature is well adapted for texture representation. This is mainly because these textons inherit the resistance to geometric and photometric transformations of local image detectors and descriptors. However, computations of most local image descriptors (Ambai and Yoshida, 2011; Bay et al., 2008; Fan et al., 2011; Lowe, 2004; Rublee et al., 2011; Winder et al., 2009) are complicated and some choices behind their specific designs are not clear. Most descriptors are carefully crafted by hand with many parameters to be manually tuned, such as the number of orientation bins, the number of grids in each support region, and grids sampling schemes (e.g., Cartesian or log-polar). Another limitation is their high dimensions that result in expensive computations in the clustering process to generate textons dictionary. Instead of using manually designed local image descriptors, some research work attempted to employ the data-driven approach to compute texture representations. You et al. (2009) applied a family of PCA subspace decompositions to recognize the very specified wood texture. A null-space based LDA in Liao and Chung (2010) was used in the frequency domain to perform texture recognition.

Similar to texture images resulting from multiple factors of geometric and photometric transformations, facial images are also formed by interactions of multiple modes related to facial geometry (e.g., person identity and expression), pose, and illumination. In

order to disentangle and extract the essential factor, i.e., person identity, for robust and fast face recognition, face images are usually mapped into a face manifold by subspace embedding techniques. Turk and Pentland (1991) proposed to use PCA to represent face images. Belhumeur et al. (1997) used LDA with the class specific linear projection to compute a face subspace. Both PCA and LDA are linear embedding methods. A number of research efforts have developed to discover the non-linear structure hidden in original image space, e.g., Isomap (Tenenbaum et al., 2000), Local Linear Embedding (LLE) Roweis et al. (2000), and Laplacian Eigenmap (Belkin and Niyogi, 2001). However, these non-linear approaches suffer the out-of-sample problem, i.e., a subspace yielded by such techniques is only defined on training data but is not able to extend to new testing data. LPP proposed by He et al. (2005) explicitly addressed this problem. LPP models a subspace by a nearest-neighbor graph where the local structure of original image space is preserved.

Motivated by the similarity of image formation between texture images and face images, we propose to use subspace embedding methods to map texture images into a texture subspace. This enables us to disentangle and extract essential factors of texture images. Compared to representations of local image descriptors, the data-driven textons based upon subspace embeddings are more distinctive, more compact, and with less parameters to tune.

3. Subspace embedding methods

We investigate both linear and non-linear embedding methods to compute a texture subspace. PCA effectively models the Euclidean structure and the variance of entire data. LDA incorporates class specific information and finds the projection that actively discriminates between different categories. LPP preserves intrinsic local structure and detects a non-linear subspace hidden in original data space.

As an illustration, Fig. 2 shows the distributions of normalized local texture patches mapped into a texture subspace with the top three dimensions. In this figure, (a)–(c) correspond to subspaces obtained by PCA, LDA, and LPP, respectively. The mapping of PCA tends to spread data to capture the factor of the maximum variance. The projection of LDA is based on the factor of texture identities, i.e., to cluster texture patches from the same class close while to separate the ones from different classes far from each other. The embedding of LPP also forms reasonably separated clusters. It maintains the similarities of local patches in the texture subspace and in the original data space.

Let us consider a set of n d -dimensional local texture patches $X = x_1, x_2, \dots, x_n$ belonging to l classes. $a \in \mathbb{R}^{d \times k}$ represents the

embedding that maps original data to a new $k \ll d$ -dimensional texture subspace, where new data $y_i \in \mathbb{R}^k$ are defined by $y_i = a^T x_i$, $i = 1, 2, \dots, n$.

3.1. EigenTextons of PCA

PCA is an eigenvector approach to model linear variations in the data with high dimensions. The goal of PCA is to construct a series of mutually orthogonal basis that are able to capture the maximum variance directions. It performs embedding by projecting original feature vectors with d dimensions to a k -dimensional linear subspace spanned by k leading eigenvectors of the covariance matrix. The objective function $J(a)$ is defined as following:

$$J(a) = \sum_{i=1}^n (y_i - \bar{y})^2, \quad \bar{y} = \sum_{i=1}^n y_i, \quad (1)$$

$$a^* = \arg \max_a J(a). \quad (2)$$

The optimal embedding a^* in Eq. (2) is the EigenTextons, which correspond to the basis that maximizes the above objective function.

3.2. FisherTextons of LDA

LDA is a supervised linear subspace embedding algorithm. By encoding class specific information, LDA seeks a projection basis on which data points of different classes are separated far from each other while simultaneously clustering feature points of the same class close to each other. Therefore, the subspace yielded by LDA is efficient for discrimination. The objective function of LDA is:

$$J(a) = \frac{a^T S_B a}{a^T S_W a}, \quad (3)$$

$$S_B = \sum_{i=1}^l n_i (m_i - m)(m_i - m)^T, \quad (4)$$

$$S_W = \sum_{i=1}^l \left(\sum_{j=1}^{n_i} (x_i^j - m_i)(x_i^j - m_i)^T \right), \quad (5)$$

where m is the mean vector of local texture patches in training set, m_i is the average feature vector of the i th class, n_i is the number of local texture patches in the i th class, x_i^j is the j th local texture patch in the i th class, l is the number of classes. S_B and S_W are between-class scatter matrix and within-class scatter matrix, where the class specific information is incorporated. The optimal mapping basis a^*

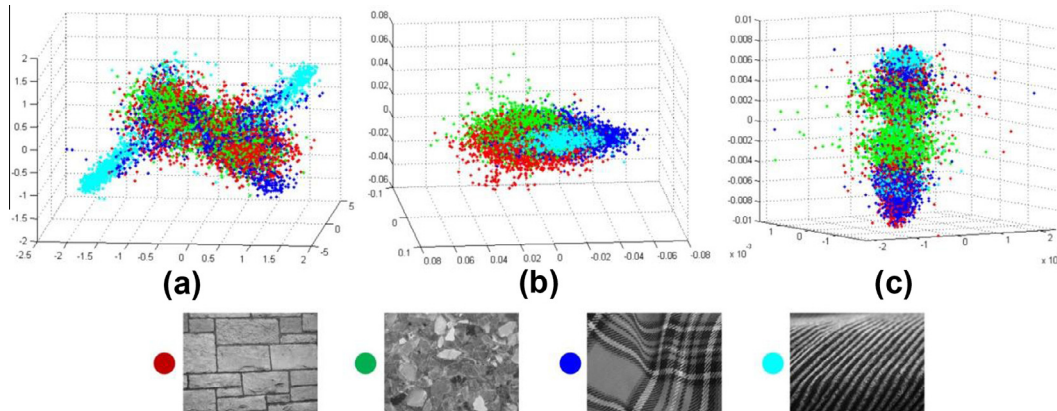


Fig. 2. Visualization of normalized texture patches in texture subspaces with top three dimensions computed by PCA (a), LDA (b), and LPP (c). Each color is encoded according to one texture class. For figure clarity, four texture classes from UIUC Texture dataset are visualized.

is the maximizer of the objective function in Eq. (3). It can be solved by a generalized eigenvalue problem defined in Eq. (6). Note the upper bound of reduced dimension $k = l - 1$ as there are at most $l - 1$ nonzero generalized eigenvalues.

$$S_B a^* = \lambda S_W a^*. \quad (6)$$

In the case of face recognition, S_W usually becomes singular. This stems from the fact that the rank of S_W is less than or equal to $n - l$, but the number of training images n is much smaller than the number of pixels d in each image. In texture representation, this difficulty however can be avoided. In our framework, n is the number of local patches in texture images of training set. This number is much larger (10^3) than the amount of images. In addition, the dimension d of each local texture patch is far smaller than the dimension of the entire image.

It was observed in Cai et al. (2007) that the coefficients of a^* used to map spatially smooth features (e.g., local texture patches) tend to become spatially rough. We take the scheme of spatially smooth regularization in Cai et al. (2007) to smooth and stabilize the mapping coefficients. Spatially smooth regularization takes advantage of the spatial relationships between pixels within each local texture patch and makes the embedding coefficients smoother and more stable. The regularized optimal embedding a^* is the FisherTextons.

3.3. LaplacianTextons of LPP

LPP is a non-linear subspace embedding approach that aims to preserve intrinsic geometry of original data space. It concentrates on discovering the manifold structure hidden in original space by modeling an adjacency graph. LPP addresses the out-of-sample problem of most non-linear embedding techniques. The embedding obtained by LPP is defined on both training and testing data. The objective function of LPP is defined as:

$$J(a) = \sum_{ij} (y_i - y_j)^2 S_{ij}, \quad (7)$$

$$S_{ij} = \begin{cases} \exp(-\|x_i - x_j\|^2 / \gamma), & \|x_i - x_j\|^2 < \varepsilon \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

where S is the adjacency matrix that measures the similarity between each pair of local texture patches (x_i, x_j) . ε defines the range of local neighborhood. γ is a constant scalar value. The intrinsic geometry of original data space is captured by S . The objective function incurs heavy penalty if a pair of neighboring texture patches are mapped far apart. So it seeks to obtain such an embedding that, if x_i and x_j are close, they will be mapped to a subspace where y_i and y_j are close as well. The optimal embedding a^* that minimizes Eq. (7) can be solved by the generalized eigenvalue problem:

$$XLX^T a = \lambda XDX^T a \quad (9)$$

where D is a diagonal matrix with $D_{ii} = \sum_j S_{ji}$. $L = D - S$ is the Laplacian matrix. The minimum eigenvalue solution a^* of Eq. (9) is the LaplacianTextons.

4. Texture representation framework

Our proposed framework to build effective texture representations is described in this section. A keypoint detector is first used to localize texture regions. We then normalize detected regions to make local texture patches invariant to scaling and rotation. The normalized texture patches are then mapped to a texture subspace using the embedding approaches described in Section 3. A textons dictionary generated from training set is employed to quantize embedded normalized texture patches. A texture image is in the end represented as a histogram of textons.

4.1. Region detection and normalization

We begin with a keypoint detector to search salient local image structures. The keypoint detector provides support regions of local texture patches. In this paper we adopt Harris–Laplace detector and Hessian–Laplace detector (Mikolajczyk and Schmid, 2005) as keypoint detectors. Both of them are rotation and scaling invariant. Harris–Laplace detector responses to corner-like structures, and Hessian–Laplace detects blob-like structures. They provide salient, complimentary, and sufficient local texture regions.

The support regions are three times larger than the detected regions in order to include more signal changes. All the support regions are first smoothed to reduce noise and aliasing and then normalized to a fixed patch size of 41×41 that provides sufficient resolution. A similar patch size was used in Mikolajczyk and Schmid (2005). As most state-of-the-art local image descriptors (Bay et al., 2008; Fan et al., 2011; Ke and Sukthankar, 2004; Lowe, 2004), a dominant orientation of a patch is computed based on gradient information. The dominant orientation corresponds to the largest bin of a histogram of gradient orientation weighted by gradient magnitudes and smoothed by a Gaussian window. A patch is then rotated to align its dominant orientation to a canonical direction. This normalization process simplifies the subspace modeling problem for embedding algorithms as variations of rotation and scaling are significantly suppressed.

4.2. Offline computation of embeddings

We compute embeddings a^* using the algorithms described in Section 3. For texture recognition, the embeddings can be pre-computed once and stored. It is important to note that the embeddings are computed based upon normalized local image patches rather than the entire images as used in face recognition. We explore two channels of normalized texture patches to compute embeddings: (1) image channel, i.e., local image patch with $41 \times 41 = 1681$ dimensions; (2) gradient channel, i.e., horizontal and vertical gradients with $2 \times 39 \times 39 = 3042$ dimensions. By using two channels of training set, we learn three embeddings: *EigenTextons*, *FisherTextons*, and *LaplacianTextons*. As discussed in Section 3, the upper bound of reduced dimension of LDA is $l - 1$. We make this number as the reduced dimension for LDA. To keep good performance and consistency with LDA, we also use the first $l - 1$ dimensions of PCA and LPP.

5. Experiments and discussions

The proposed texture representation approaches are evaluated in the context of texture classification. As discussed in Sections 3 and 4, we have three embedding methods and two feature channels. So there are 6 different combinations of texture representations that are investigated in our experiments as shown in Table 1. We extensively compare the performances of our proposed methods with the existing state-of-the-arts. They are tested on two public available datasets: UIUC Texture (Lazebnik et al., 2005) and UMD Texture (Xu et al., 2009). In addition to in-plane rotation and scaling change presented in traditional datasets

Table 1

Texture representations based upon different combinations of embeddings and feature channels.

Embeddings	Image channel	Gradient channel
EigenTextons	PCA-Img	PCA-Grad
FisherTextons	LDA-Img	LDA-Grad
LaplacianTextons	LPP-Img	LPP-Grad

(Brodatz, 1996; Caputo et al., 2010; Varma and Zisserman, 2009), the two datasets as shown in Fig. 3 capture more challenging variations including viewpoint, illumination, and non-rigid surface deformation.

5.1. Experimental setup

The UIUC dataset includes 25 texture classes and 40 images with the resolution of 640×480 in each class. These images present strong rotation, scaling, viewpoint variation, non-rigid surface deformation, and lighting change. The UMD dataset consists of 1000 uncalibrated and unregistered images with the resolution of 1280×960 pixels. It contains 25 texture categories with 40 images for each class. These images are also taken under significant geometric and photometric transformations. We downsample original images of UMD dataset to the resolution of 640×480 .

In order to facilitate a fair comparison, we follow the standard experimental setting to randomly select a portion of images from each class as the training set. The remaining images are used as the testing set. The training process is based on each corresponding

randomly generated training set. The reported recognition accuracy rates in the following experiments are the average results over 50 runs by the random generated training and testing sets. K -means clustering ($K = 100$) is employed to build the textons dictionary. We employ Support Vector Machines (SVMs) with RBF kernels as the classifier. The optimal parameters of RBF kernels are obtained by 5-fold cross-validation. The SVMs classifier in essence finds the hyperplane that separates two-class data with maximal margin. In order to apply SVMs for multi-class problem, we take the *one-versus-one* strategy.

We compare the proposed methods against the state-of-the-art approaches including VG (Varma and Garg, 2007), MFS (Xu et al., 2009), Lazebnik et al. (2005), (Zhang et al. (2007), SIFT (Lowe, 2004), SURF (Bay et al., 2008), DAISY (Winder et al., 2009), ORB (Rublee et al., 2011), CARD (Ambai and Yoshida, 2011), and MROGH (Caputo et al., 2010). VG makes use of local density function properties of a set of image measurements. MFS combines the fractal dimensions of pixel sets grouped by three local density functions. Lazebnik extracts local image features by RIFT and SPIN from affine regions. Zhang combines local features by multiple

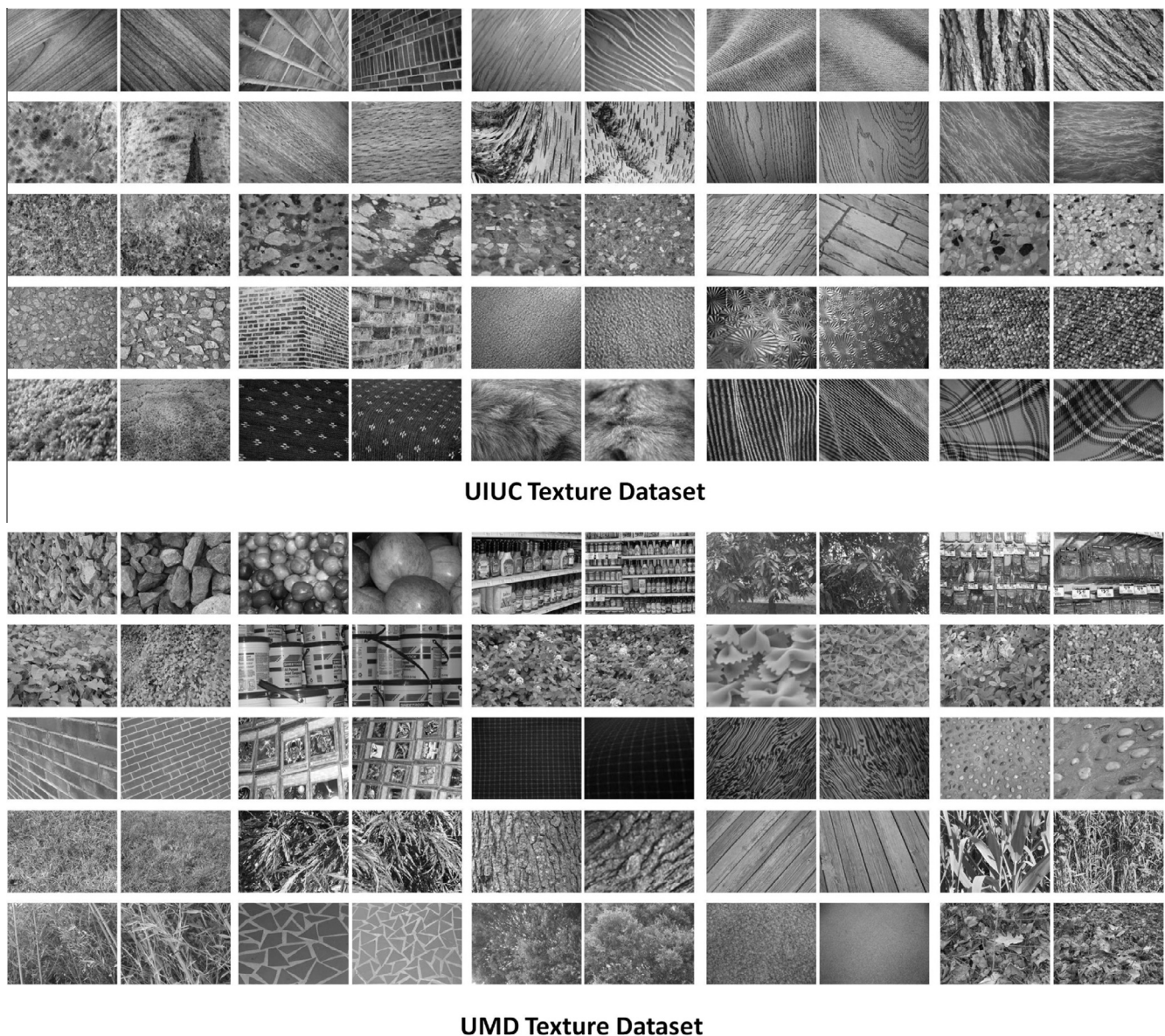


Fig. 3. Two sample images of 25 texture categories in UIUC and UMD Texture Datasets.

local image detectors and descriptors. SIFT, SURF, DAISY, ORB, CARD, and MROGH are the most recent developed local image descriptors. They achieve the state-of-the-art performances in the context of object, texture, and scene classification.

5.2. Evaluations of different combinations of embeddings and feature channels

The classification accuracies for different combinations of embeddings and feature channels on UIUC dataset and UMD dataset are shown in Fig. 4. The numbers of training images for each class are from 1 to 20 and the rest images are used for testing. Similar conclusions can be drawn from experiments on the two datasets.

As shown in the two figures, for each embedding method, the performances based on gradient channel outperform the ones based on image channel. This is probably because gradient is more resistant to lighting variations and preserves relative changes in intensity. Thus, gradient feature simplifies embedding methods to model the essential factor of texture images. The difference between two channels of LPP is more evident. As shown in Eq. (8), LPP measures pair-wise similarities of local texture patches. So the adjacency matrix S is more sensitive to illumination, which results in that LPP is relatively less robust to handle lighting change.

LDA-Grad dominates the recognition rates when the training number is small. By explicitly encoding class labels in computing the texture subspace, LDA is forced to concentrate on the essential factor of texture identities, i.e., textons. PCA-Grad and LPP-Grad also demonstrate impressive performances. When the training images are sufficient (e.g., >10), the performances of PCA-Grad and LPP-Grad are comparable to LDA-Grad. But LPP-Grad is more sensitive to the number of training images because non-linear methods require denser sampling of a manifold to reasonably recover the intrinsic structure.

5.3. Comparisons with the state-of-the-arts

Based on the evaluation results from different combinations of embedding methods and feature channels, we choose LDA-Grad to compare with the state-of-the-art methods for texture recognition on UIUC dataset and UMD dataset. The experimental results are shown in Table 2. N_t denotes the number of training images

in each class. The best recognition rates of various training numbers are the numbers in bold. We can also obtain similar conclusions on both datasets.

The results in the two tables show that our proposed method outperforms the state-of-the-art approaches in most cases. For example, our method significantly and consistently outperforms ORB and CARD, both of which are the most recent state-of-the-art local image descriptors. The performance of our method is also much better than texture representations based on fractal analysis, i.e., VG and MFS. In most cases, our approach achieves better performances than the remaining methods that are based on the state-of-the-art local image descriptors. The impressive performances based on sophisticated descriptors originate from the resistance to photometric and geometric transformations of local image descriptors. Compared to local image descriptors that are carefully crafted by hand, our methods are totally data-driven. It is based on the construction of a texture subspace where the essential factor (textons) is manifested but unwanted variation factors are reduced or removed. Our method is inferior to SIFT when $N_t = 5$. This is probably due to the fact that only 5 training images cannot provide sufficiently dense sampling of a texture subspace. The computation of embeddings is therefore biased by the rough sampling.

5.4. Computational cost of textons

The experimental results have demonstrated that textons-based methods are well-adapted for texture representation. In natural texture images, textons can be generated by clustering local texture features. However, the clustering process is always time consuming. If the clustering problem is exactly solved, the computational cost of K -means is $O(n^{d+1} \log n)$ (Inaba et al., 1994), where n is the number of local texture features to be clustered; c is the number of centers; and d is the dimension of feature. So when n and c are fixed, feature with fewer dimensions are able to reduce the computational cost and speed up clustering process.

The local image features computed by most descriptors are with high dimensions which result in expensive computations. As discussed in Section 3.2, the upper bound of reduced dimension of FisherTextons is $l - 1$, where l is the number of classes. Both UIUC dataset and UMD dataset contains 25 classes. So we use 24 as the reduced dimension of textons. Fig. 5 compares the running time in

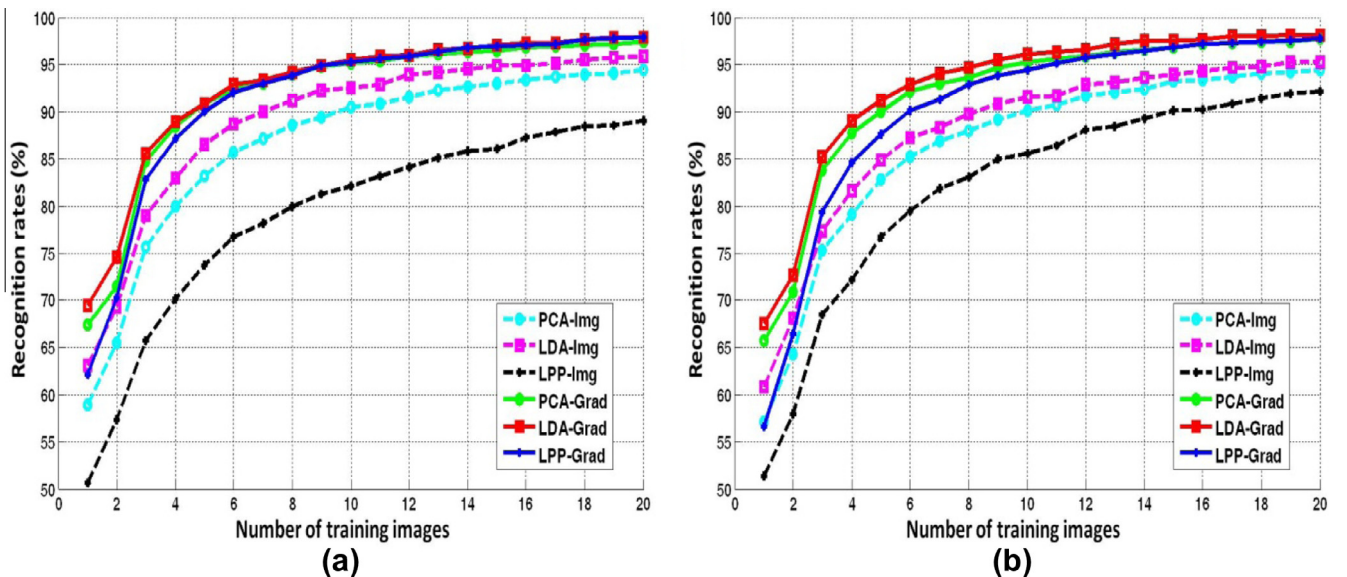


Fig. 4. Recognition accuracy for different combinations of embedding methods and feature channels on (a) UIUC dataset and (b) UMD dataset.

Table 2

Recognition rates of proposed method compared with the state-of-the-art approaches on (top) UIUC dataset and (bottom) UMD dataset. N_t is the number of training images in each class.

N_t	VG	MFS	Lazebnik	Zhang	SIFT	SURF	DAISY	ORB	CARD	MROGH	Our method
5	82.86	82.24	91.12	88.62	91.96	90.73	86.80	79.03	73.99	88.76	90.86
10	87.85	88.36	94.42	93.17	95.42	95.15	92.54	86.26	83.00	94.13	95.55
15	90.62	91.38	96.64	95.33	96.87	96.14	94.16	89.40	87.18	95.93	97.07
20	92.31	92.74	97.02	96.67	97.84	96.75	95.21	90.73	89.69	96.82	97.91
5	90.92	85.63	90.71	91.56	91.68	90.41	90.81	81.85	84.38	90.39	91.23
10	94.09	90.82	94.54	96.00	96.01	94.49	94.92	87.87	90.41	94.54	96.06
15	96.22	92.67	96.29	96.79	97.21	96.13	96.47	90.87	93.05	96.01	97.59
20	96.36	93.93	96.95	97.62	97.64	96.98	97.58	92.84	94.23	97.03	98.20

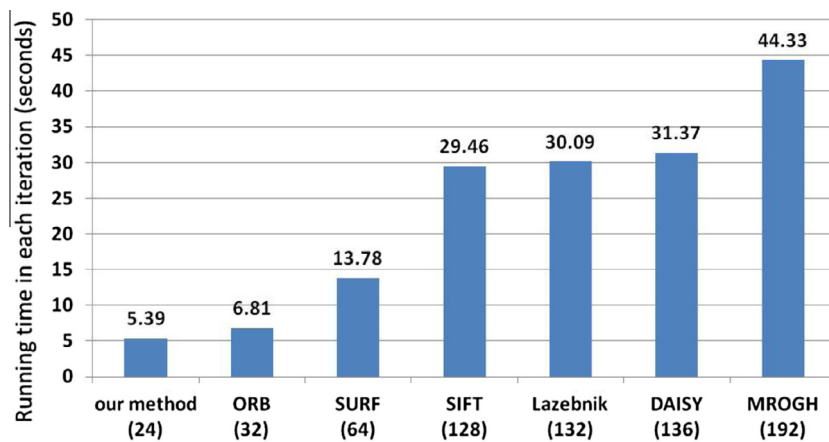


Fig. 5. The running times of different methods in each iteration of K -means clustering. The numbers in parenthesis under each method denote corresponding feature dimensions.

each iteration of clustering on UIUC dataset on an Intel Core2 CPU 2.13 GHz computer. We use 20 images of each class as the training set and extract approximately 3000 local patches from each image. We set the number of clustering centroids $c = 100$. In the experiments, our method significantly reduces the running time compared with most state-of-the-art descriptors. Note the running time difference can become huge when a clustering needs a large number (e.g., 1000) of iterations.

5.5. Discussions

The experimental results in the context of texture recognition have validated the effectiveness of our proposed texture representation methods. They also reveal a number of interesting points:

First, in all embedding methods, gradient channel consistently performs better than image channel, especially for LPP. These experimental results demonstrate that gradient channel is more suitable for embedding approaches to model the texture subspace as gradients suppress lighting variation but preserve relative intensity change.

Second, both of the linear embedding methods with gradient channel achieve the state-of-the-art classification results. PCA provides the benefits of capturing the maximum variance of original data space but reducing noisy variations. This also confirms that the major factors, i.e., the leading eigenvectors, of texture images correspond to texture identities even though significant variations are presented. LDA encodes the class specific information in the texture subspace which enables the mapping actively discriminates between different texture classes.

Third, the non-linear embedding method, i.e., LPP-Grad, also achieves impressive performance on both datasets but is inferior to LDA-Grad when the training samples are insufficient. This is different from the observation in face recognition domain where non-linear methods consistently outperform linear ones. This might be explained by the difference of data sampling. Facial images are always densely sampled, which forms smooth variations in terms of pose and expression. The smooth changes make the Euclidean distance $\|x_i - x_j\|$ in Eq. (8) more accurate as they are small enough to preserve the geodesic distance or the intrinsic geometry hidden in a facial manifold. In contrast, images in texture datasets always present drastic changes which are not smooth enough to capture local structures in original texture manifold.

6. Conclusion

In this paper, we have proposed several texture representations by subspace embeddings. To the best of our knowledge, this is the first work on texture representation that systematically and explicitly considers the texture subspace using both linear and non-linear embedding algorithms. The experimental results on benchmark texture datasets have demonstrated the texture subspace computed by embedding methods is effective to disentangle and extract the essential factor of texture images from the interactions of multiple factors resulting from geometric and photometric transformations. The experimental results also show that the state-of-the-art performances on existing texture classification datasets are now near ceiling (e.g., >97%). But in addition to classification accuracy, our methods significantly improve the computational

costs and are totally data-driven with much fewer parameters to tune. The experiments have validated that textures mapped into a texture subspace have strong resistance to image deformations, meanwhile, are more distinctive and more compact. The future work will focus on effective combinations (e.g., through Multiple Kernel Learning) of texture representations computed from different embedding methods.

Acknowledgement

This work was supported in part by NSF grant IIS-0957016, EFRI-1137172, NIH 1R21EY020990, ARO grant W911NF-09-1-0565, and FHWA DTFH61-12-H-00002.

References

- Ambai, M., Yoshida, Y., 2011. CARD: Compact and Real-time Descriptors. In Proc. IEEE International Conference on Computer Vision, 97–104.
- Bay, H., Ess, A., Gool, L., 2008. SURF: speed up robust features. *Comput. Vision Image Understand* 110 (3), 346–359.
- Belhumeur, P., Hespanha, J., Kriegman, D., 1997. Eigenfaces vs. Fisherfaces: recognition using class specific linear projection. *IEEE Trans. Pattern Anal. Mach. Intell.* 19 (7), 711–720.
- Belkin, M., Niyogi, P., 2001. Laplacian Eigenmaps and Spectral Techniques for Embedding and Clustering. In Proc. *Advances in Neural Information Processing Systems*, 585–591.
- Bouguila, N., 2012. Infinite Liouville mixture models with application to text and texture categorization. *Pattern Recognit. Lett.* 33 (2), 103–110.
- Brodatz, P., 1996. *Texture: A Photographic Album for Artists and Designers*, 66. Dover, New York.
- Cai, D., He, X., Hu, Y., Han, J., Huang, T., 2007. Learning A Spatially Smooth Subspace for Face Recognition. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 1–7.
- Caputo, B., Hayman, E., Fritz, M., Eklundh, J., 2010. Classifying materials in the real world. *Image Vision Comput.* 28 (1), 150–163.
- Csurka, G., Bray, C., Dance, C., Fan, L., 2004. Visual Categorization with Bags of Keypoints. In Proc. *European Conference on Computer Vision Workshop*.
- Fan, B., Wu, F., Hu, Z., 2011. Aggregating Gradient Distributions into Intensity Orders: A Novel Local Image Descriptor. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 2377–2384.
- Haralick, R., 1979. Statistical and Structural Approaches to Texture. In Proc. IEEE 67 (5), 786–804.
- He, X., Yan, S., Hu, Y., Niyogi, P., Zhang, H., 2005. Face recognition using Laplacianfaces. *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (3), 328–340.
- Hua, G., Brown, M., Winder, S., 2007. Discriminant Embedding for Local Image Descriptors. In Proc. IEEE International Conference on Computer Vision, 1–8.
- Inaba, M., Katoh, N., Imai, H., 1994. Application of Weighted Voronoi Diagrams and Randomization to Variance-based K-clustering. In Proc. *ACM Symposium on Computational Geometry*, 332–339.
- Ke, Y., Sukthankar, R., 2004. PCA-SIFT: A More Distinctive Representation for Local Image Descriptors. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 506–513.
- Lazebnik, S., Schmid, C., Ponce, J., 2005. A sparse texture representation using local affine regions. *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (8), 1265–1278.
- Liao, S., Chung, C., 2010. A New Subspace Learning Method in Fourier Domain for Texture Classification. In Proc. IEEE International Conference on Image Processing, 4589–4592.
- Lowe, D., 2004. Distinctive Image Features from Scale-invariant Keypoints. *International Journal of Computer Vision* 60 (2), 91–110.
- Mailing, A., Cernuschi-Frias, B., 1982. A method for mixed states texture segmentation with simultaneous parameter estimation. *Pattern Recognit. Lett.* 21 (15), 1982–1989.
- Martinez, A., Kak, A., 2001. PCA vs. LDA. *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (2), 228–233.
- Mikolajczyk, K., Schmid, C., 2005. A performance evaluation of local descriptors. *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (10), 1615–1630.
- Nguyen, K., Sabata, B., Jain, A., 2012. Prostate cancer grading: gland segmentation and structural features. *Pattern Recognit. Lett.* 33 (7), 951–961.
- Randen, T., Husoy, J., 1999. Filtering for texture classification: a comparative study. *IEEE Trans. Pattern Anal. Mach. Intell.* 21 (4), 291–310.
- Roweis, S., Saul, L., 2000. Nonlinear Dimensionality Reduction by Locally Linear Embedding. *Science* 290 (5500), 2323–2326.
- Rublee, E., Rabaud, V., Konolige, K., Bradski, G., 2011. ORB: An Efficient Alternative to SIFT and SURF. In Proc. IEEE International Conference on Computer Vision, 2564–2571.
- Yang, X., Yuan, S., Tian, Y., 2011. Recognizing clothes patterns for blind people by confidence margin based feature combination. In: *ACM Multimedia*.
- Tenenbaum, J., Silva, V., Langford, J., 2000. A global geometric framework for nonlinear dimensionality reduction. *Science* 290.
- Turk, M., Pentland, A., 1991. Face Recognition Using Eigenfaces. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 586–591.
- Varma, M., Zisserman, A., 2009. A statistical approach to material classification using image patches. *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (11), 2032–2047.
- Varma, M., Garg, R., 2007. Locally Invariant Fractal Features for Statistical Texture Classification. In Proc. IEEE International Conference on Computer Vision, 1–8.
- Winder, S., Hua, G., Brown, M., 2009. Picking the Best DAISY. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 178–185.
- Xu, Y., Ji, H., Fermuller, C., 2006. A Projective Invariant for Textures. In Proc. IEEE Conference on Computer Vision and Pattern Recognition, 1932–1939.
- Xu, Y., Ji, H., Fermuller, C., 2009. Viewpoint invariant texture description using fractal analysis. *Int. J. Comput. Vision* 83 (1), 85–100.
- You, M., Cai, C., 2009. Wood Classification based on PCA, 2DPCA, (2D)2PCA, and LDA. In Proc. IEEE International Symposium on Knowledge Acquisition and Modeling, 371–374.
- Zhang, J., Marszałek, M., Lazebnik, S., Schmid, C., 2007. Local features and kernels for classification of texture and object categories: a comprehensive study. *Int. J. Comput. Vision* 73 (2), 213–238.