

JOINT KERNEL DICTIONARY AND CLASSIFIER LEARNING FOR SPARSE CODING VIA LOCALITY PRESERVING K-SVD

Weyang Liu[†], Zhiding Yu[‡], Meng Yang^{§*}, Lijia Lu[†], and Yuexian Zou[†]

[†]School of ECE, Peking University [‡]Dept. of ECE, Carnegie Mellon University

[§]College of Computer Science & Software Engineering, Shenzhen University

ABSTRACT

We present a locality preserving K-SVD (LP-KSVD) algorithm for joint dictionary and classifier learning, and further incorporate kernel into our framework. In LP-KSVD, we construct a locality preserving term based on the relations between input samples and dictionary atoms, and introduce the locality via nearest neighborhood to enforce the locality of representation. Motivated by the fact that locality-related methods works better in a more discriminative and separable space, we map the original feature space to the kernel space, where samples of different classes become more separable. Experimental results show the proposed approach has strong discrimination power and is comparable or outperforms some state-of-the-art approaches on public databases.

Index Terms— Discriminative Dictionary Learning, Locality Preserving K-SVD, Kernel Space.

1. INTRODUCTION

Sparse coding has served an important role in a wide variety of vision problems, ranging from image restoration [1], image denoising [2] to image classification [3], etc. The technique linearly represents a query image \mathbf{y} by a few atoms from a dictionary \mathbf{D} , i.e., $\mathbf{y} = \mathbf{D}\mathbf{x}$ or $\mathbf{y} \approx \mathbf{D}\mathbf{x}$, where \mathbf{x} is sparse coefficients. In image classification, one classifies a query based on the sparse codes. Thus, the discrimination power of the dictionary is of considerable importance. Wright et al. [3] employ the entire set of training samples as the dictionary and achieve promising results on face recognition. However, the method is time-consuming, since it needs to solve the l_1 minimization problem, which reduces its scalability for large scale databases. Alternatively, some off-the-shelf bases (e.g. Fourier, wavelets) can also be used as dictionaries [4], yet these dictionaries might not be the optimal choices in certain tasks like image classification or face recognition. Learning the desired dictionary from training samples has been reported to bring additional performance gain to classification

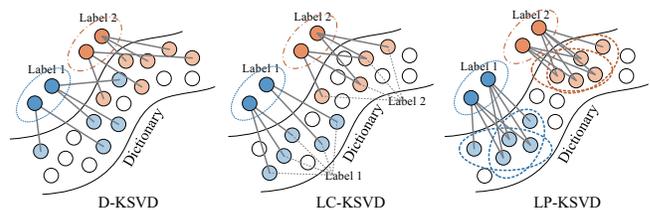


Fig. 1. Intuitive comparison among D-KSVD, LC-KSVD and LP-KSVD. The selected atoms for representation are highlighted.

tasks [5–11]. Dictionary learning methods boosts the performance from perspective of coding coefficients, dictionary, or both. Group sparse coding [9, 10] regularizes the coding coefficients to be similar within the same class. Graph-regularized sparse coding [11] combines geometrical information into coding coefficients and dictionary in order to preserve the local manifold structure. Moreover, dictionaries with labels are learned to exploit the class-specific representation residual for classification [3, 7, 8]. For instance, [8] learns a structured dictionary via Fisher criteria and enforce coding coefficients to have small within-class scatter but big between-class scatter.

Of great interest is a series of recently proposed dictionary learning approaches [5, 6, 12] via K-SVD algorithm [13]. Since K-SVD algorithm focuses on the representation power without considering the discrimination power, [12] proposes to jointly learn a classifier based on coding coefficients for classification tasks. Pham et al. [12] combine the classification error of the linear classifier, the representation error of data as well as the regularization terms into the objective function, minimizing it by iteratively updating the variables while preserving the sparsity. However, this approach involves multiple additional optimizations, which is inefficient and easy to be trapped into local minima. To address this, discriminative K-SVD (D-KSVD) [5] eliminates the representation error term of unlabeled data and formulates the objective function into the K-SVD framework. In order to further enhance the discrimination power of the learned dictionary, label consistent K-SVD (LC-KSVD) [6] constructs a label consistent term and applies similar method to formulate the optimization into the K-SVD framework. The label consistent term supervises each class has similar sparse codes, exhibiting discrimination.

Inspired by recent progress obtained via enforcing locality [14, 15], we propose the locality preserving K-SVD (LP-

Corresponding Author {yang.meng@szu.edu.cn; wylu@pku.edu.cn}.

This work is partially supported by the National Natural Science Foundation of China under Grant No. 61402289, and Shenzhen Scientific Research and Development Funding Program under Grant JCYJ20140509172609171.

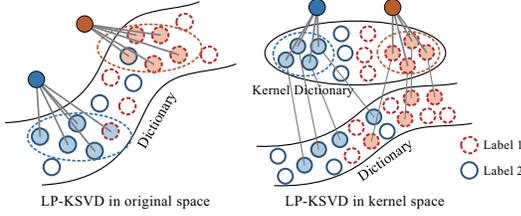


Fig. 2. Comparison between LP-KSVD in original space and LP-KSVD in kernel space. For LP-KSVD, the label information (label) in dictionary is not necessarily a prior. We show the label just to imply the hidden unsupervised label consistency when applying LP-KSVD in kernel space.

KSVD) approach in which dictionary and classifier are jointly learned similar to [5, 6]. Specifically, we first construct a penalty term that enforces the representation by local dictionary atoms similar to the input sample and penalises the representation by nonlocal (dissimilar) atoms. The classification error term, representation error term and locality penalty term are then jointly optimized by the K-SVD algorithm. An intuitive comparison among D-KSVD, LC-KSVD and LP-KSVD is provided in Fig. 1. We can see that the locality-preserving term could keep the relationship between input samples and dictionary atoms in dictionary learning, making the nearby samples (usually with the same label) represented by similar dictionary atoms. In fact, Laplacian sparse coding [16] also uses the local information to construct graph Laplacian for feature quantization. It preserves the consistence in sparse representation of similar local features with graph Laplacian and apply one-vs-all SVM as classifier, while we enforce locality into coding coefficients and dictionary as suggested in [14] and jointly optimize a classifier via K-SVD.

The label consistent term in LC-KSVD serves a similar role to the classification error term. They both penalize the dictionary representation from classes different from the input class. On the other hand, different from classification error, the locality preserving term incorporates the neighborhood information into dictionary and coding coefficients, thus retaining more crucial information. We believe the locality preserving term is endowed with more discriminative power than label consistent term. The representation obtained via the proposed dictionary learning can be further enhanced by exploiting the nonlinear structure within the data [17]. [18] also shows the benefit of representation power from kernel dictionaries in classification tasks. Moreover, the same benefit is reported in [19, 20]. Different from [18], however, we specifically aim at the image classification by jointly learning the optimal linear classifier and simultaneously taking locality into consideration. Further, mapping the original feature space into kernel space implies more separable representation between different classes can be achieved as Fig. 2 illustrates, since the locality preserving term is imposed in kernel space.

2. PRELIMINARIES

Given an input image set matrix $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k\} \in \mathbb{R}^{n \times k}$ in which each column vector \mathbf{y}_i represents a n -

dimensional input sample, we let $\mathbf{D} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_m\} \in \mathbb{R}^{n \times m}$ be the dictionary that is composed of m training samples, and $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\} \in \mathbb{R}^{m \times k}$ be the coefficient matrix where each column is a sparse code for an input sample. Learning a dictionary to sparsely represent \mathbf{Y} can be accomplished by solving the minimization problem:

$$\langle \mathbf{D}, \mathbf{X} \rangle = \arg \min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2 \quad \text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c \quad (1)$$

where $\|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2$ denotes the representation error, and c is a sparsity constraint parameter which enforces each sparse code should have fewer than c nonzero values. This objective function is exactly what K-SVD algorithm optimizes with, and it emphasizes more on representation power instead of discrimination power. Therefore, [5, 6] jointly learn a classifier at the same time and combine the classification error in the objective function. Towards this end, a classifier $g(\mathbf{x}, \mathbf{W})$ is constructed with parameter matrix $\mathbf{W} \in \mathbb{R}^{m \times m}$. An optimal classifier should satisfy

$$\mathbf{W} = \arg \min_{\mathbf{W}} \sum_i \mathcal{L}\{l_i, g(\mathbf{x}_i, \mathbf{W})\} + \lambda \|\mathbf{W}\|_F^2 \quad (2)$$

where \mathcal{L} is the classification loss function that is usually defined as logistic loss function or quadratic loss function, l_i is the label of \mathbf{y}_i , and λ is a regularization parameter. Combining Eq.(1) and Eq.(2) can jointly learn the dictionary and classifier as in [5, 6, 12] by defining objective function:

$$\begin{aligned} \langle \mathbf{D}, \mathbf{W}, \mathbf{X} \rangle = \arg \min_{\mathbf{D}, \mathbf{W}, \mathbf{X}} & \\ \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2 + \sum_i \mathcal{L}\{l_i, g(\mathbf{x}_i, \mathbf{W})\} + \lambda \|\mathbf{W}\|_F^2. & \quad (3) \\ \text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c & \end{aligned}$$

When using the l_2 norm as the loss function, the case becomes [12] without the unlabeled data therein, and it is exactly the same case in D-KSVD [5]. For LC-KSVD [6], it further enforces a label consistent term into the objective function, making the coding coefficients concentrate on its own class.

3. KERNEL DICTIONARY AND CLASSIFIER LEARNING VIA LOCALITY PRESERVING K-SVD

3.1. Locality Preserving K-SVD

As suggested in locality-constrained linear coding (LLC) [14] and local coordinate coding (LCC) [15], locality plays a significant role in classification. We therefore incorporate locality preserving constraint into the objective function:

$$\begin{aligned} \langle \mathbf{D}, \mathbf{W}, \mathbf{T}, \mathbf{X} \rangle = \arg \min_{\mathbf{D}, \mathbf{W}, \mathbf{T}, \mathbf{X}} & \\ \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2 + \alpha \|\mathbf{H} - \mathbf{W}\mathbf{X}\|_2^2 & \\ + \beta \|\mathbf{P} - \mathbf{T}\mathbf{X}\|_2^2 + \lambda \|\mathbf{W}\|_F^2 & \quad (4) \\ \text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c & \end{aligned}$$

where α, β are scaling parameters that control the contribution of classification error term and locality preserving term respectively, and $\|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2$ and $\|\mathbf{H} - \mathbf{W}\mathbf{X}\|_2^2$ are representation error term and classification error term respectively. $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_k\} \in \mathbb{R}^{v \times k}$ denotes the class labels of the input images \mathbf{Y} . Each column of \mathbf{H} is a label vector, formulated as $\mathbf{h}_i = \{0, 0, \dots, 1, \dots, 0, 0\}$ in which 1 indicates the positive label in the according category. \mathbf{W} denotes a set of parameters for a linear classifier. Most importantly, locality preserving term $\|\mathbf{P} - \mathbf{T}\mathbf{X}\|_2^2$ is added to penalize the representation which is not locally concentrated. To design \mathbf{P} , we first calculate the q_i nearest neighbors of each input image \mathbf{y}_i in dictionary, and then use $\omega(q, \mathbf{y}_i)$ to denote the positions of these q nearest neighbors in dictionary, where nonzero elements stand for the positions of nearest neighbors in dictionary. For example, in $\omega(q, \mathbf{y}_i) = \{1, 0, \dots, 0, 1, 0, \dots, 1, 1\}^T$, 1 represents the positions of nearest neighbors. Therefore, we construct the matrix \mathbf{P} via

$$\mathbf{P} = \{\omega(q, \mathbf{y}_1), \omega(q, \mathbf{y}_2), \dots, \omega(q, \mathbf{y}_k)\}. \quad (5)$$

Consider an example that \mathbf{D} has 4 samples from 2 categories, each category with 2 samples. Given \mathbf{Y} of 4 input samples, 2-nearest neighbors of $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3$ and \mathbf{y}_4 are $(\mathbf{d}_1, \mathbf{d}_3), (\mathbf{d}_1, \mathbf{d}_2), (\mathbf{d}_3, \mathbf{d}_4)$ and $(\mathbf{d}_2, \mathbf{d}_4)$ respectively. Then \mathbf{P} is constructed as

$$\mathbf{P} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}. \quad (6)$$

However, the desired \mathbf{D} is not given in the first place, so we optimize a \mathbf{D} with the locality preserving term removed and use it to initialize \mathbf{P} . \mathbf{P} contains all the locality information for the input samples, which is essential to preserve locality. \mathbf{T} is a linear transformation matrix, which transforms the original sparse codes \mathbf{x} to be more locally concentrated. Particularly, we construct the locality term in a similar form to the classification error so as to conveniently formulate the objective function for K-SVD algorithm. Furthermore, we believe the representation codes should still be sparse since locality must lead to sparsity but not necessary vice versa.

In order to apply the K-SVD algorithm to efficiently find the optimal solution, we drop the regularization term and further reformulate the objective function in Eq.(4) as

$$\begin{aligned} \langle \mathbf{D}, \mathbf{W}, \mathbf{T}, \mathbf{X} \rangle &= \arg \min_{\mathbf{D}, \mathbf{W}, \mathbf{T}, \mathbf{X}} \\ &\left\| \begin{pmatrix} \mathbf{Y} \\ \alpha^{\frac{1}{2}} \mathbf{H} \\ \beta^{\frac{1}{2}} \mathbf{P} \end{pmatrix} - \begin{pmatrix} \mathbf{D} \\ \alpha^{\frac{1}{2}} \mathbf{W} \\ \beta^{\frac{1}{2}} \mathbf{T} \end{pmatrix} \mathbf{X} \right\|_2^2 \\ &\text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c \end{aligned} \quad (7)$$

which is the generalized form optimized in [13]. We can directly use the K-SVD algorithm to solve the problem.

3.2. Objective Function

As formulated in Eq.(4), the objective function for dictionary learning contains four penalty terms, which are the representation error $\|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_2^2$, classification error $\|\mathbf{H} - \mathbf{W}\mathbf{X}\|_2^2$, locality preserving term $\|\mathbf{P} - \mathbf{T}\mathbf{X}\|_2^2$ and parameter regularizer $\|\mathbf{W}\|_F^2$ respectively.

Representation error: It is a fundamental constraint in sparse coding, which is used to keep the estimation close to the input sample in order to keep the representation power.

Classification error: As the main contribution to the discrimination power, this constraint penalize misclassifications. It makes sparse codes more discriminative and also learns a linear classifier for classification in a supervised way, avoiding to use dictionary classifier [3], since K-SVD algorithm may disable such classifier.

Locality preserving term: Motivated by the fact that the samples in the same class usually have similar local dictionary atoms, we construct this term to preserve the locality of representation and bring more discrimination power. The distance metric used to find neighbors of input samples is flexible, but we use simple l_1 distance in experiments. The direct intuition of locality preserving term is that the sparse codes become locally concentrated after a linear invertible transformation. In fact, the locality information is preserved within sparse codes via such constraint. We let $\mathbf{X}^N = \mathbf{T}\mathbf{X}$ and further rewrite $\mathbf{D}^N = \mathbf{D}\mathbf{T}^{-1}$, $\mathbf{W}^N = \mathbf{W}\mathbf{T}^{-1}$. Thus the penalty terms become $\|\mathbf{Y} - \mathbf{D}^N \mathbf{X}^N\|_2^2$, $\|\mathbf{H} - \mathbf{W}^N \mathbf{X}^N\|_2^2$ and $\|\mathbf{P} - \mathbf{X}^N\|_2^2$ respectively. From this formulation, we can clearly see the representation error and classification error stay the same form as in Eq.(4), while the locality preserving term make the sparse codes represented by local atoms.

3.3. Learning in Kernel Space

We can easily generalize a linear classifier into a nonlinear one via kernel function [21, 22], enhancing the discrimination power of the original linear classifier. To this end, we use kernel function to create a nonlinear mapping mechanism $\mathbf{y} \in \mathbb{R} \mapsto \phi(\mathbf{y}) \in \mathbb{H}$ where \mathbb{H} is a unique associated reproducing kernel Hilbert space (RKHS). We map every sample into higher dimensional kernel space via transformation ϕ and use the kernel function $K(\mathbf{y}', \mathbf{y}'') = \phi(\mathbf{y}')^T \phi(\mathbf{y}'')$, where $\mathbf{y}', \mathbf{y}''$ are different samples and ϕ denotes the implicit nonlinear mapping associated with the kernel function $K(\mathbf{y}', \mathbf{y}'')$. The kernel dictionary is written as

$$\phi(\mathbf{D}) = \{\phi(\mathbf{d}_1), \phi(\mathbf{d}_2), \dots, \phi(\mathbf{d}_m)\} \in \mathbb{R}^{s \times m} \quad (8)$$

where s is the dimension of the kernel space that is much larger than the dimension of the original feature space. Similarly, the input image set \mathbf{Y} is kernelized to $\phi(\mathbf{Y})$. Thus, the representation error term is rewritten as $\|\phi(\mathbf{Y}) - \phi(\mathbf{D})\mathbf{X}\|$.

However, the high dimensional kernel space may lead to high computational complexity, so we will perform dimen-

sionality reduction in kernel space, similar to [19,20]. Specifically, the projection matrix $\mathbf{R} \in \mathbb{R}^{s \times u}$ ($u < s$) is constructed following the similar approach in KPCA [22]. Combining the projection matrix \mathbf{R} , we derive the new representation error term $\|\mathbf{R}^T \phi(\mathbf{Y}) - \mathbf{R}^T \phi(\mathbf{D})\mathbf{X}\|$. Further, we assume that each column vector in \mathbf{R} is a linear combination of samples in kernel space, and decompose \mathbf{R} as

$$\begin{aligned} \mathbf{R} &= \{\mathbf{R}_1, \dots, \mathbf{R}_s\} = \phi(\mathbf{D}) \cdot \Psi \\ &= \{\phi(\mathbf{d}_1), \dots, \phi(\mathbf{d}_m)\} \cdot \{\psi_1, \dots, \psi_u\} \end{aligned} \quad (9)$$

in which ψ_i is a m -dimensional column vector that is also linear projection coefficients of columns in $\phi(\mathbf{D})$, satisfying $\mathbf{R}_i = \phi(\mathbf{D}) \cdot \psi_i$. Moreover, Ψ is also called pseudo-transformation matrix [19]. Putting Eq.(9) into the representation error, we can derive $\|\Psi^T K(\mathbf{D}, \mathbf{Y}) - \Psi^T \mathbf{G}\mathbf{X}\|$ where

$$\begin{aligned} K(\mathbf{D}, \mathbf{Y}) &= \{K(\mathbf{D}, \mathbf{y}_1), \dots, K(\mathbf{D}, \mathbf{y}_k)\} \\ &= \begin{Bmatrix} K(\mathbf{d}_1, \mathbf{y}_1) & \dots & K(\mathbf{d}_1, \mathbf{y}_k) \\ \vdots & \ddots & \vdots \\ K(\mathbf{d}_m, \mathbf{y}_1) & \dots & K(\mathbf{d}_m, \mathbf{y}_k) \end{Bmatrix}, \end{aligned} \quad (10)$$

and \mathbf{G} ($G_{ij} = K(\mathbf{d}_i, \mathbf{d}_j)$), also equal to $\phi(\mathbf{D})^T \cdot \phi(\mathbf{D})$. Particularly, \mathbf{G} is defined as the kernel gram matrix that should be symmetric and positive semi-definite based on Mercer's theorem. Since \mathbf{G} and $K(\mathbf{D}, \mathbf{Y})$ are known as a prior, dimensionality reduction now focuses on finding Ψ instead of \mathbf{R} . Various methods to find Ψ are provided in [23]¹. Therefore, the kernel dictionary learning can be formulated as

$$\begin{aligned} \langle \mathbf{G}, \mathbf{W}, \mathbf{T}, \mathbf{X} \rangle &= \arg \min_{\mathbf{G}, \mathbf{W}, \mathbf{T}, \mathbf{X}} \\ &\left\| \begin{Bmatrix} \Psi^T K(\mathbf{D}, \mathbf{Y}) \\ \alpha^{\frac{1}{2}} \mathbf{H} \\ \beta^{\frac{1}{2}} \mathbf{P} \end{Bmatrix} - \begin{Bmatrix} \Psi^T \mathbf{G} \\ \alpha^{\frac{1}{2}} \mathbf{W} \\ \beta^{\frac{1}{2}} \mathbf{T} \end{Bmatrix} \mathbf{X} \right\|_2^2 \\ &\text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c \end{aligned} \quad (11)$$

from which we can learn a locality-preserving kernel dictionary and a corresponding optimal linear classifier. The dictionary learned this way addresses the desirable locality property and is able to adapt to the underlying structure of the input samples for better representation and classification.

3.4. Optimization

Following the basic protocol in the original K-SVD algorithm, we can efficiently find the optimal solution for Eq.(11). For conciseness, we let $\mathbf{Y}' = \{\Psi^T K(\mathbf{D}, \mathbf{Y}), \alpha^{\frac{1}{2}} \mathbf{H}, \beta^{\frac{1}{2}} \mathbf{P}\}^T$ and $\mathbf{D}' = \{\Psi^T \mathbf{G}, \alpha^{\frac{1}{2}} \mathbf{W}, \beta^{\frac{1}{2}} \mathbf{T}\}^T$, rewriting Eq.(11) as

$$\langle \mathbf{D}', \mathbf{X} \rangle = \arg \min_{\mathbf{D}', \mathbf{X}} \|\mathbf{Y}' - \mathbf{D}'\mathbf{X}\|_2^2 \quad \text{s.t. } \forall i, \|\mathbf{x}_i\|_0 \leq c \quad (12)$$

¹We find no advantages of using complex method to construct Ψ in our experiments, so we simply use identity matrix as Ψ in order to retain the intuitive interpretation.

which is identical to the problem that K-SVD algorithm solves. K-SVD algorithm updates the dictionary atom by atom until convergence by solving the following problem:

$$\langle \mathbf{d}'_k, \mathbf{x}_T^k \rangle = \arg \min_{\mathbf{d}'_k, \mathbf{x}_T^k} \|\mathbf{E}_k - \mathbf{d}'_k \mathbf{x}_T^k\|_F^2 \quad (13)$$

where \mathbf{x}_T^k is the k th row in \mathbf{X} (\mathbf{x}_k is the k th column), and $\mathbf{E}_k = \mathbf{Y} - \sum_{j \neq k} (\mathbf{d}'_j \mathbf{x}_T^j)$. Note that, \mathbf{E}_k denotes the error for all input samples while the k th atom is removed. Then we discard the zero entries in \mathbf{x}_T^k and obtain the row vector \mathbf{x}_R^k . We also define \mathbf{E}_k^R to stand for the selection of error column vectors that correspond to samples that use the atom \mathbf{d}_k . Thus Eq.(13) is equivalent to the following minimization:

$$\langle \mathbf{d}'_k, \mathbf{x}_R^k \rangle = \arg \min_{\mathbf{d}'_k, \mathbf{x}_R^k} \|\mathbf{E}_k^R - \mathbf{d}'_k \mathbf{x}_R^k\|_F^2. \quad (14)$$

Singular value decomposition (SVD) can be utilized to solve Eq.(14). For the error matrix \mathbf{E}_k^R , SVD decomposes it to $\mathbf{E}_k^R = \mathbf{U}\Sigma\mathbf{V}^T$. According to the decomposition results, the solution for \mathbf{d}'_k and the representation coefficients \mathbf{x}_R^k is $\mathbf{d}'_k = \mathbf{U}(:, 1)$, $\mathbf{x}_R^k = \Sigma(1, 1) \cdot \mathbf{V}(:, 1)$. Then the nonzero entries in \mathbf{x}_T^k are replaced by \mathbf{x}_R^k . The updating is done iteratively until convergence. The algorithm is given in Algorithm 1.

Algorithm 1 Kernel Dictionary and Classifier Learning via LP-KSVD

Input: $\mathbf{Y}, \mathbf{P}^{(0)}, \mathbf{H}, \Psi, K(\cdot), \alpha, \beta, c, r, i=0$

Output: $\mathbf{D}^{ke} (\mathbf{D}^{ke} = \Psi^T \mathbf{G}), \mathbf{T}, \mathbf{W}$

- 1: Obtain $\mathbf{D}^{(0)}$ by constructing the dictionary with training atoms. It does not matter whether atoms of the same label are grouped together.
 - 2: Kernelize the dictionary by computing $\mathbf{D}^{ke} = \Psi^T \mathbf{G}$ where $\mathbf{G} = \phi(\mathbf{D}^{(0)})^T \cdot \phi(\mathbf{D}^{(0)})$. In the kernel dictionary, K-SVD algorithm is used to initialize $\mathbf{D}^{ke(0)}$.
 - 3: Kernelize the input samples \mathbf{Y} via $\Psi^T K(\mathbf{D}, \mathbf{Y})$.
 - 4: Compute sparse codes $\mathbf{X}^{(0)}$ for $\Psi^T K(\mathbf{D}, \mathbf{Y})$ via Eq.(1) with $\mathbf{D}^{ke(0)}$ as dictionary.
 - 5: Initialize \mathbf{T}, \mathbf{W} with $\mathbf{T}^{(0)}, \mathbf{W}^{(0)}$ via Eq.(15).
 - 6: Initialize \mathbf{Y}', \mathbf{D}' which are defined in Eq.(14).
 - 7: Update each column and the corresponding representation coefficient by iteratively solving Eq.(14) via K-SVD algorithm.
 - 8: Obtain $\Psi^T \mathbf{G}, \mathbf{T}, \mathbf{W}$ via Eq.(16).
 - 9: Let $i \leftarrow i + 1$ and use the newly learned dictionary to recompute $\mathbf{P}^{(i)}$.
 - 10: If $i = r$, output $\Psi^T \mathbf{G}, \mathbf{T}, \mathbf{W}$; Else, go to Step 1;
-

For initialization, we first compute $\mathbf{P}^{(0)}$ via the \mathbf{D}^* obtained from Eq(11) with $\beta=0$. It is equivalent to removing the locality preserving term to obtain \mathbf{D}^* . Then we use \mathbf{D}^* to compute $\mathbf{P}^{(0)}$. For $\mathbf{W}^{(0)}, \mathbf{T}^{(0)}$, we use multivariate regression to initialize:

$$\begin{aligned} \mathbf{W}^{(0)} &= \mathbf{H}(\mathbf{X}^{(0)})^T ((\mathbf{X}^{(0)})(\mathbf{X}^{(0)})^T + \mu_1 \mathbf{I})^{-1} \\ \mathbf{T}^{(0)} &= \mathbf{P}(\mathbf{X}^{(0)})^T ((\mathbf{X}^{(0)})(\mathbf{X}^{(0)})^T + \mu_1 \mathbf{I})^{-1} \end{aligned} \quad (15)$$

where H, P are given a prior, μ_1, μ_2 are two small parameters, and I represents an identical matrix.

Since P is updated iteratively, the stopping criteria needs to be discussed. We perform an experiment to evaluate the convergence of P , as shown in Fig 3(a). The recognition rate corresponding to $P^{(i)}$ in each iteration is also shown in Fig 3(a). From Fig. 3(a), we can see the matrix P becomes stable after approximately 5 times iterations. In fact, we have tested on multiple database and found out 5 times iterations are enough for an appropriate P . Moreover, we sum up P from each iteration (totally 100 iterations), and obtain $P^\Sigma = \sum_i P^{(i)}$. We binarize P^Σ and draw it in Fig. 3(b) where black is 0 and white is 1. It shows that locality preserving term has hidden label consistent property, since label consistent matrix [6] is a block-diagonal matrix. LP-KSVD is adaptive to the local relations between atoms and training samples. Interestingly, in an ideal scenario that samples of each label are perfectly separated, LP-KSVD becomes LC-KSVD.

3.5. Classification Strategy

We compute the sparse codes of the query with the learned kernel dictionary, and then use the linear classifier to get the label. The original K-SVD algorithm requires D' to be normalized column-wise, so W is learned with the unnormalized $D^{ke} = \Psi^T G = \{d_1^{ke}, \dots, d_m^{ke}\}$, making it impossible to directly use dictionary D^{ke} and classifier parameters $W = \{w_1, \dots, w_m\}$ for classification. We do the following operations to obtain \hat{D}^{ke} and \hat{W} :

$$\begin{aligned} \hat{D}^{ke} = \{\hat{d}_1^{ke}, \dots, \hat{d}_m^{ke}\} &= \left\{ \frac{d_1^{ke}}{\|d_1^{ke}\|_2}, \dots, \frac{d_m^{ke}}{\|d_m^{ke}\|_2} \right\} \\ \hat{W} = \{\hat{w}_1, \dots, \hat{w}_m\} &= \left\{ \frac{w_1}{\|d_1^{ke}\|_2}, \dots, \frac{w_m}{\|d_m^{ke}\|_2} \right\}. \end{aligned} \quad (16)$$

The relation between the desired (\hat{D}^{ker}, \hat{W}) and the obtained (D^{ker}, W) is shown as follows (\bar{T}, T are left out):

$$\begin{aligned} y &\approx D^{ke} x = \sum x_j \|d_j^{ke}\|_2 \frac{d_j^{ke}}{\|d_j^{ke}\|_2} = \sum x_j^* \hat{d}_j^{ke} = \hat{D}^{ke} x^* \\ h &\approx W x = \sum x_j \|d_j^{ke}\|_2 \frac{w_j}{\|d_j^{ke}\|_2} = \sum x_j^* \hat{w}_j = \hat{W} x^* \end{aligned} \quad (17)$$

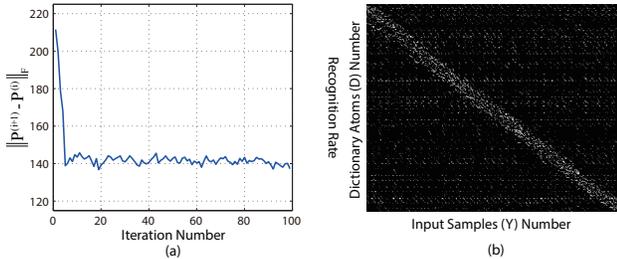


Fig. 3. (a) An example of the convergence of P on extended YaleB database. (b) Binary Representation of matrix P^Σ . This experiment is performed on extended YaleB database and under the same settings in [6].

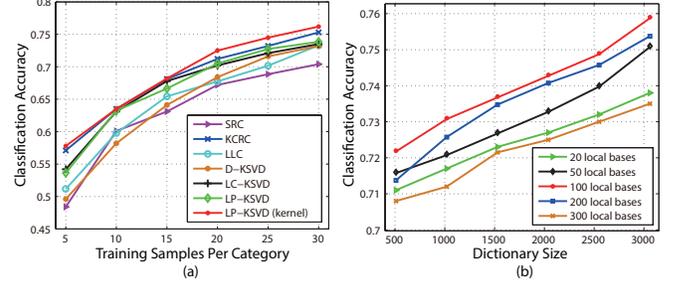


Fig. 4. (a) Classification accuracy on Caltech101 database. (b) Evaluation of different local bases under varying dictionary size.

where x^* denotes a new set of sparse codes, y is a query sample and h is the label vector. We obtain the label by first coding the query over dictionary \hat{D}^{ker} , and the label of the query refers to the index of the largest element in $h = \hat{W} x^*$.

4. EXPERIMENTS AND RESULTS

This section uses $\alpha=4, \beta=4$ and adopts RBF kernel with parameter 0.5 in all experiments. Note that, The matrix P is obtained by 5 iterations. α and β are determined by fivefold cross validation on the training dataset.

4.1. Caltech101 Database

Caltech101 database [24] contains 9144 images from 102 classes with significant variance in shape. We train on 5, 10, 15, 20, 25, 30 samples per category and test on the rest. Dictionary size is set as 510, 1020, 1530, 2040, 2550 and 3060 for 5, 10, 15, 20, 25 and 30 samples per category respectively. Note that, we adopt the spatial pyramid features. The results are averaged for 10 times and shown in Fig. 4(a). We also evaluate the classification performance with different local bases under varying dictionary size trained in LP-KSVD. Note that, we use 30 samples per category to train the dictionary. Fig. 4(a) shows the proposed LP-KSVD (kernel) performs well in image classification tasks and surpasses most competitive approaches. LP-KSVD is 0.3% to 0.4% better than LC-KSVD in 25 and 30 training samples per category. LP-KSVD with kernel outperforms LP-KSVD and LC-KSVD by approximate 2.5%. Fig. 4(b) shows the proposed method with 100 local atoms has the best performance, indicating the number of local bases has an optimal range in which the classification accuracy reaches its best.

4.2. Extended YaleB Database

Extended YaleB database [25] contains 2414 frontal face under different illumination conditions. We randomly select 32 images as training samples, the rest as testing samples, and adopt the 504 dimensional randomface features. For K-SVD based algorithms, We use 32 images per person to train a dictionary of size 570. The number of local bases for LP-KSVD

is set as 40. We run the experiment for 10 times and obtain the average recognition rate in Table 1. The running time² for dictionary training and query testing is shown in Table 2. We can see the proposed LP-KSVD outperforms LC-KSVD in kernel space and other competitive approaches. As for the efficiency of LP-KSVD, the time to classify a query is approximate 0.3ms, which is fast enough for applications.

Table 1. Recognition results (%) on extended YaleB database. For SRC, we use a dictionary of 570 atoms, same as D-KSVD, LC-KSVD and LP-KSVD.

Method	Accuracy	Method	Accuracy
SRC [3]	81.7	LLC [14] (25 bases)	83.1
D-KSVD [5]	93.5	LLC [14] (50 bases)	89.7
LC-KSVD [6]	95.0	LC-KSVD [6] (kernel)	94.4
LP-KSVD	94.8	LP-KSVD (kernel)	94.9

Table 2. Average dictionary training time and running time for classifying a testing image on extended YaleB database.

Method	Testing Time (ms)
SRC [3]	19.72
LC-KSVD [6]	0.298
LP-KSVD	0.328
LP-KSVD (kernel)	0.536

4.3. 15 Scenes Database

15 scenes categories database [26] contains 200 to 400 images per category, including kitchen, mountain and store. We randomly select 100 images per category for training and the rest for testing. The size of learned dictionary is set as 450. Note that, we adopt the spatial pyramid features. Results are shown in Table 3. Our proposed method achieves better performance than SRC, LLC, D-KSVD and LC-KSVD. We can see the kernel indeed boosts the performance for both LC-KSVD and LP-KSVD. LP-KSVD performs better than LC-KSVD and LP-KSVD with kernel achieves the best classification accuracy in some competitive methods.

Table 3. Classification results (%) on 15 scenes database.

Method	Accuracy	Method	Accuracy
SRC [3]	91.8	LLC [14] (30 bases)	90.1
D-KSVD [5]	89.2	LSC [16]	89.9
LC-KSVD [6]	92.9	LC-KSVD [6] (kernel)	93.5
LP-KSVD	93.4	LP-KSVD (kernel)	94.0

5. CONCLUDING REMARKS

This paper proposes the locality preserving K-SVD algorithm and elaborates the gain that locality brings especially in kernel space. Specifically, a locality preserving term is designed to enforce locality into dictionary and coding coefficients. In order to make locality information more discriminative, we map the original feature space to the kernel space, making samples of different classes more separable. LP-KSVD approach in kernel space and its specific algorithm are presented and comprehensively discussed. Experimental results validate the su-

²We use a PC with 3.4 GHz dual-core CPU and 16GB RAM. All experiments presented in the paper are performed in this PC.

periority of LP-KSVD in kernel space and also show that our approach achieves state-of-the-art results on extended YaleB database, Caltech101 database and 15 scenes database.

6. REFERENCES

- [1] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Non-local sparse models for image restoration," in *ICCV*, 2009. 1
- [2] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE TIP*, 2006. 1
- [3] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE TPAMI*, 2009. 1, 3.2, 1, 2, 3
- [4] K. Huang and S. Aviyente, "Sparse representation for signal classification," in *NIPS*, 2006. 1
- [5] Q. Zhang and B. Li, "Discriminative k-svd for dictionary learning in face recognition," in *CVPR*, 2010. 1, 1, 2, 2, 2, 1, 3
- [6] Z. Jiang, Z. Lin, and L. S. Davis, "Label consistent k-svd: learning a discriminative dictionary for recognition," *IEEE TPAMI*, 2013. 1, 1, 2, 2, 2, 3.4, 3, 1, 2, 3
- [7] M. Yang, D. Dai, L. Shen, and L. V. Gool, "Latent dictionary learning for sparse representation based classification," in *CVPR*, 2014. 1
- [8] M. Yang, L. Zhang, X. Feng, and D. Zhang, "Sparse representation based fisher discrimination dictionary learning for image classification," *IJCV*, 2014. 1
- [9] Y. Sun, Q. Liu, J. Tang, and D. Tao, "Learning discriminative dictionary for group sparse representation," *IEEE TIP*, 2014. 1
- [10] S. Bengio, F. Pereira, Y. Singer, and D. Strelow, "Group sparse coding," in *NIPS*, 2009. 1
- [11] M. Zheng, J. Bu, C. Chen, C. Wang, L. Zhang, G. Qiu, and D. Cai, "Graph regularized sparse coding for image representation," *IEEE TIP*, 2011. 1
- [12] D. Pham and S. Venkatesh, "Joint learning and dictionary construction for pattern recognition," in *CVPR*, 2008. 1, 2, 2
- [13] M. Aharon, M. Elad, and A. Bruckstein, "K-svd: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE TSP*, 2006. 1, 3.1
- [14] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong, "Locality-constrained linear coding for image classification," in *CVPR*, 2010. 1, 3.1, 1, 3
- [15] K. Yu, T. Zhang, and Y. Gong, "Nonlinear learning using local coordinate coding," in *NIPS*, 2009. 1, 3.1
- [16] S. Gao, I. W. Tsang, L. Chia, and P. Zhao, "Local features are not lonely-laplacian sparse coding for image classification," in *CVPR*, 2010. 1, 3
- [17] S. Gao, I. W. Tsang, and L. T. Chia, "Kernel sparse representation for image classification and face recognition," in *ECCV*, 2010. 1
- [18] H. Nguyen, V. M. Patel, N. M. Nasrabadi, and R. Chellappa, "Kernel dictionary learning," in *ICASSP*, 2012. 1
- [19] W. Liu et al., "A novel kernel collaborative representation approach for image classification," in *ICIP*, 2014. 1, 3.3, 3.3
- [20] W. Liu and Z. Yu et al., "Kcrc-lcd: Discriminative kernel collaborative representation with locality constrained dictionary for visual categorization," *arXiv:1410.4673*, 2014. 1, 3.3
- [21] B. Scholkopf and A. J. Smola, *Learning with kernels: support vector machines, regularization, optimization, and beyond*, MIT press, 2001. 3.3
- [22] B. Schölkopf, A. Smola, and K. Müller, "Kernel principal component analysis," in *ICANN*, 1997. 3.3, 3.3
- [23] L. Zhang, W. Zhou, P. Chang, J. Liu, Z. Yan, T. Wang, and F. Li, "Kernel sparse representation-based classifier," *IEEE TSP*, 2012. 3.3
- [24] Li F., R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," *CVIU*, 2007. 4.1
- [25] A. S. Georghades, P. N. Belhumeur, and D. Kriegman, "From few to many: Illumination cone models for face recognition under variable lighting and pose," *IEEE TPAMI*, 2001. 4.2
- [26] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *CVPR*, 2006. 4.3