

Adaptive Kernel Principal Component Analysis With Unsupervised Learning of Kernels

Daoqiang Zhang Zhi-Hua Zhou
National Laboratory for Novel Software Technology
Nanjing University, Nanjing 210093, China
{zhangdq, zhouzh}@lamda.nju.edu.cn

Songcan Chen
Department of Computer Science and Engineering
NUAA, Nanjing 210016, China
s.chen@nuaa.edu.cn

Abstract

Choosing an appropriate kernel is one of the key problems in kernel-based methods. Most existing kernel selection methods require that the class labels of the training examples are known. In this paper, we propose an adaptive kernel selection method for kernel principal component analysis, which can effectively learn the kernels when the class labels of the training examples are not available. By iteratively optimizing a novel criterion, the proposed method can achieve nonlinear feature extraction and unsupervised kernel learning simultaneously. Moreover, a non-iterative approximate algorithm is developed. The effectiveness of the proposed algorithms are validated on UCI datasets and the COIL-20 object recognition database.

1 Introduction

As a powerful nonlinear feature extraction method, kernel principal component analysis (KPCA) [10] has been widely used in many applications [7]. The essence of KPCA is to perform principal component analysis (PCA) in the transformed high-dimensional feature space through an implicit nonlinear mapping from the original input space to the feature space. Here it is not necessary to know the explicit form of the nonlinear mapping and only the inner products between two data points in the feature space are needed. Since those inner products in feature space can be equivalently and efficiently computed by a kernel function in the original input space, KPCA elegantly avoids the ‘curse of dimensionality’ encountered by many classical algorithms working in high-dimensional space. How-

ever, like other kernel-based methods such as support vector machines (SVM) [7] and kernel discriminant analysis (KDA) [7], the performance of KPCA is greatly affected by the choice of the kernel and parameters. In other words, the selection of the optimal kernel and parameters is crucial for KPCA to obtain good performance.

Early research mainly focus on the selection of parameters for a certain kernel [2]. There are two widely used approaches for this purpose. The first approach empirically chooses a series of candidate values for the concerned kernel parameter, executes the learning algorithm using every candidate value, and finally assigns the value corresponding to the best performance to the kernel parameter. The second approach is the well-known cross-validation [2], which is also widely used in model selection. Both approaches are time-consuming and can hardly examine a large range of parameters. A recent advances in kernel parameter selection is to use gradient-based methods for adaptive learning the kernel parameters when class labels of the training data are available [14].

Learning the optimal kernel directly from a set of kernels has attracted much attention during the past few years [9, 11]. This kind of methods are based on the observation that in practice no kernel is the best. Thus, seeking the optimal combination of a set of kernels seems more reasonable than only adjusting the parameters of a single kernel. This kind of kernel selection methods include kernel alignment [3], learning the kernel matrix directly with semi-definite programming [6], hyperkernels [9, 11], and idealized kernels [5], etc. However, most of those methods require that the class labels of the training examples are available. Since KPCA works in an unsupervised learning setting where the class labels of the training examples are not known, those

methods can hardly be applied to KPCA directly. Recently, Yang et al. [12] proposed a novel kernel selection method called fisher+kernel criterion for KDA, where the key idea is to generate the feature matrix from multiple kernels. Yet this method is still only applicable when the class labels of training examples are available.

This paper proposes Adaptive Kernel Principal Component Analysis (A-KPCA), which can effectively learn the kernels under the unsupervised learning setting. Inspired by [12], we first transform the original 1D input vectors into 2D feature matrices through a set of nonlinear mappings induced from different kernels, each corresponding to one column of the 2D feature matrix. Then, two coupled sets of projective vectors are extracted from those feature matrices using an iterative procedure. One set of projective vectors corresponds to the column direction of feature matrices and is used for nonlinear feature extraction, while the other corresponds to the row direction of feature matrices and is used for searching the optimal combination of kernels. Moreover, an efficient non-iterative algorithm is proposed to approximate A-KPCA. The effectiveness of the proposed algorithms are validated by extensive experiments.

2 Kernel Principal Component Analysis

Given a training set $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where the training examples $\mathbf{x}_i \in \mathbb{R}^n$. Let $\psi : \mathbf{x} \in X \rightarrow \psi(\mathbf{x}) \in F$ be a nonlinear mapping from the original input space X to a high-dimensional feature space F . The inner product in F is defined as a kernel $\mathcal{K}(\mathbf{x}, \mathbf{y}) = \psi(\mathbf{x})^T \psi(\mathbf{y})$ in the original input space. Denote $\Psi = (\psi(\mathbf{x}_1), \psi(\mathbf{x}_2), \dots, \psi(\mathbf{x}_n))$, and $\bar{\psi} = 1/n \sum_{i=1}^n \psi(\mathbf{x}_i)$. Without loss of generality, assume that the data are centered in F , i.e. $\bar{\psi} = 0$, then the total scatter matrix is defined as $\mathbf{S}_t = \sum_{i=1}^n (\psi(\mathbf{x}_i) - \bar{\psi})(\psi(\mathbf{x}_i) - \bar{\psi})^T = \Psi \Psi^T$.

KPCA uses the criterion shown in Eq. 1 to compute the optimal projective vector \mathbf{v} .

$$J(\mathbf{v}) = \sum_{k=1}^n \|\mathbf{v}^T \psi(\mathbf{x}_k)\|^2 = \mathbf{v}^T \mathbf{S}_t \mathbf{v} \quad (1)$$

Maximizing Eq. 1 is equivalent to solving the eigenvalue problem: To find $\lambda \geq 0$ and eigenvectors $\mathbf{v} \in F - \{0\}$ satisfying $\lambda \mathbf{v} = \mathbf{S}_t \mathbf{v}$. It is easy to verify that all solutions \mathbf{v} with $\lambda \neq 0$ lie within the span of $\{\psi(\mathbf{x}_1), \psi(\mathbf{x}_2), \dots, \psi(\mathbf{x}_n)\}$, i.e. $\mathbf{v} = \Psi \mathbf{q}$. Thus, we can consider the equivalent problem shown in Eq. 2.

$$\lambda \mathbf{q} = \mathbf{K} \mathbf{q} \quad (2)$$

Here $\mathbf{K} = \Psi^T \Psi$ is the $n \times n$ kernel matrix. Suppose $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_d$ are the solutions of Eq. 2 corresponding to the largest d eigenvalues, then $\mathbf{v}_i = \Psi \mathbf{q}_i$ is the solution to Eq. 1. To extract features for a new example \mathbf{x} with KPCA, one

simply projects the image $\psi(\mathbf{x})$ onto \mathbf{v}_i as shown in Eq. 3.

$$\mathbf{v}_i^T \psi(\mathbf{x}) = \mathbf{q}_i^T \Psi^T \psi(\mathbf{x}) = \sum_{j=1}^n \mathbf{q}_{ij} \mathcal{K}(\mathbf{x}_i, \mathbf{x}) \quad (3)$$

Eqs. 2 and 3 show that the computation of KPCA components does not involve the nonlinear mapping ψ , and only the kernel is needed. As mentioned before, the performance of KPCA is greatly affected by the selection of kernels and parameters.

3 A-KPCA

Let $\phi_i : \mathbf{x} \in X \rightarrow \phi_i(\mathbf{x}) \in \mathcal{H}_i, i \in \{1, 2, \dots, f\}$ be a set of nonlinear mappings from the original input space X to a high-dimensional feature space \mathcal{H}_i . The inner products in \mathcal{H}_i are defined as the kernels $\mathcal{K}^i(\mathbf{x}, \mathbf{y}) = \phi_i(\mathbf{x})^T \phi_i(\mathbf{y})$ respectively. From ϕ_i , we define $\hat{\phi}_i : \mathbf{x} \in X \rightarrow \hat{\phi}_i(\mathbf{x}) = (\underbrace{\mathbf{0}^T \dots \mathbf{0}^T}_{1 \dots i-1}, \phi_i(\mathbf{x})^T, \underbrace{\mathbf{0}^T \dots \mathbf{0}^T}_{i+1 \dots f})^T \in \mathcal{H}, i \in \{1, 2, \dots, f\}$. Here \mathcal{H} is the Hilbert space as the direct sum of \mathcal{H}_i and the j -th $\mathbf{0}$ vector lies in \mathcal{H}_j . The inner product in \mathcal{H} can be defined as: $\hat{\phi}_i(\mathbf{x})^T \hat{\phi}_j(\mathbf{y}) = 0$ ($i \neq j$) and $\hat{\phi}_i(\mathbf{x})^T \hat{\phi}_i(\mathbf{y}) = \phi_i(\mathbf{x})^T \phi_i(\mathbf{y}) = \mathcal{K}^i(\mathbf{x}, \mathbf{y})$.

For an original training example \mathbf{x}_k , map it into the high-dimensional feature space with all $\hat{\phi}_i$'s and align the mapped vectors column by column to form a 2D feature matrix, denoted by $\Phi_k = (\hat{\phi}_1(\mathbf{x}_k), \hat{\phi}_2(\mathbf{x}_k), \dots, \hat{\phi}_f(\mathbf{x}_k))$. Thus, the original vector-based representation of training data turns into a matrix-based representation. Denote $\bar{\Phi} = \frac{1}{n} \sum_{i=1}^n \Phi_i$, and without loss of generality we assume that Φ_k 's have zero means, i.e. $\bar{\Phi} = 0$.

For a series of 2D feature matrices $\Phi_1, \Phi_2, \dots, \Phi_n$, we want to seek two matrices \mathbf{U} and \mathbf{R} , such that the criterion shown in Eq. 4 is optimized.

$$J(\mathbf{U}, \mathbf{R}) = \sum_{k=1}^n \|\mathbf{U}^T \Phi_k \mathbf{R}\|_F^2 \quad (4)$$

Here $\|\cdot\|_F$ is the Frobenius norm of matrix, $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d)$ is the projective vectors corresponding to the column direction of Φ_k and its purpose is for nonlinear feature extraction, while $\mathbf{R} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_g)$ is the projective vectors corresponding to the row direction of Φ_k and its purpose is for kernel selection.

It is easy to verify that $\mathbf{U} \subseteq \text{span}(\hat{\phi}_i(\mathbf{x}_k)), 1 \leq i \leq f, 1 \leq k \leq n$. Denote $\Phi = (\Phi_1, \Phi_2, \dots, \Phi_n)$, then $\mathbf{U} = \Phi \mathbf{L}$. The size of the original Φ is $f \times n$, which is generally very large. To reduce the size of Φ , we replace it with $\tilde{\Phi} = (\sum_{i=1}^f \hat{\phi}_i(\mathbf{x}_1), \sum_{i=1}^f \hat{\phi}_i(\mathbf{x}_2), \dots, \sum_{i=1}^f \hat{\phi}_i(\mathbf{x}_n))$ so that $\mathbf{U} = \tilde{\Phi} \mathbf{L}, \mathbf{L} \in \mathbb{R}^{n \times d}$. Thus, Eq. 4 can be replaced

by the criterion shown in Eq. 5.

$$\begin{aligned} J(\mathbf{L}, \mathbf{R}) &= \sum_{k=1}^n \|\mathbf{L}^T \tilde{\Phi}^T \Phi_k \mathbf{R}\|_F^2 \\ &= \sum_{k=1}^n \|\mathbf{L}^T \mathbf{K}_k \mathbf{R}\|_F^2 \end{aligned} \quad (5)$$

Here we constrain $\mathbf{L} \in \mathbb{R}^{n \times d}$, $\mathbf{L}^T \mathbf{L} = \mathbf{I}$ and $\mathbf{R} \in \mathbb{R}^{f \times g}$, $\mathbf{R}^T \mathbf{R} = \mathbf{I}$, and $\mathbf{K}_k = \tilde{\Phi}^T \Phi_k$ is the $n \times f$ kernel matrix, whose (i, j) -th element is computed according to Eq. 6.

$$\begin{aligned} \mathbf{K}_k(i, j) &= \left(\sum_{l=1}^f \hat{\phi}_l(\mathbf{x}_i) \right)^T \hat{\phi}_j(\mathbf{x}_k) \\ &= \mathcal{K}^j(\mathbf{x}_i, \mathbf{x}_k) \end{aligned} \quad (6)$$

Eq. 6 indicates that the kernel matrices in Eq. 5 can be directly computed from the kernel function bank $\{\mathcal{K}^i\}$. To the best of our knowledge, there is no closed form solution to Eq. 5. Inspired by [13], the following theorem presents an iterative procedure to solve this problem.

Theorem 1 *Let \mathbf{L} and \mathbf{R} be the optimal solution to Eq. 5, then: (1) For a given \mathbf{R} , \mathbf{L} consists of the d eigenvectors of the matrix $\mathbf{M}_L = \sum_{k=1}^n \mathbf{K}_k \mathbf{R} \mathbf{R}^T \mathbf{K}_k^T$ corresponding to the largest d eigenvalues; (2) For a given \mathbf{L} , \mathbf{R} consists of the g eigenvectors of the matrix $\mathbf{M}_R = \sum_{k=1}^n \mathbf{K}_k^T \mathbf{L} \mathbf{L}^T \mathbf{K}_k$ corresponding to the largest g eigenvalues.*

Proof. Equation 5 can be rewritten as

$$\begin{aligned} \sum_{k=1}^n \|\mathbf{L}^T \mathbf{K}_k \mathbf{R}\|_F^2 &= \sum_{k=1}^n \text{trace}(\mathbf{L}^T \mathbf{K}_k \mathbf{R} \mathbf{R}^T \mathbf{K}_k^T \mathbf{L}) \\ &= \text{trace}(\mathbf{L}^T \sum_{k=1}^n (\mathbf{K}_k \mathbf{R} \mathbf{R}^T \mathbf{K}_k^T) \mathbf{L}) \\ &= \text{trace}(\mathbf{L}^T \mathbf{M}_L \mathbf{L}) \end{aligned} \quad (7)$$

Here, for a given \mathbf{R} , the maximum of Eq. 5 or Eq. 7 is obtained only if \mathbf{L} consists of the d eigenvectors of the matrix \mathbf{M}_L corresponding to the largest d eigenvalues. Similarly, for a given \mathbf{L} , the maximum of Eq. 5 or $\text{trace}(\mathbf{R}^T \mathbf{M}_R \mathbf{R})$ is obtained only if \mathbf{R} consists of the g eigenvectors of the matrix \mathbf{M}_R corresponding to the largest g eigenvalues. \square

Theorem 1 provides an iterative procedure for computing \mathbf{L} and \mathbf{R} , i.e. the A-KPCA algorithm, as shown in **Algorithm 1**.

After obtaining \mathbf{L} and \mathbf{R} by A-KPCA, we can use them to extract the nonlinear features for an unseen instance \mathbf{x} . First, construct the kernel matrix $\mathbf{K}_{test}(i, j) = \mathcal{K}^j(\mathbf{x}_i, \mathbf{x})$. Then, project \mathbf{K}_{test} according to $\mathbf{C} = \mathbf{L}^T \mathbf{K}_{test} \mathbf{R}$. After that, the nonlinear features are contained in \mathbf{C} .

Algorithm 1: A-KPCA

Input: Training set $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$.

Output: Matrices \mathbf{L} and \mathbf{R} .

1. Construct the kernel matrix \mathbf{K}_k for each \mathbf{x}_k using Eq. 6.
2. Obtain initial \mathbf{R}_0 and set $i \leftarrow 1$.
3. For given \mathbf{R}_{i-1} , Compute the d eigenvectors $\mathbf{L}_i = (l_1, l_2, \dots, l_d)$ of $\mathbf{M}_L = \sum_{k=1}^n \mathbf{K}_k \mathbf{R}_{i-1} \mathbf{R}_{i-1}^T \mathbf{K}_k^T$ corresponding to the largest d eigenvalues.
4. For given \mathbf{L}_i , Compute the g eigenvectors $\mathbf{R}_i = (r_1, r_2, \dots, r_g)$ of $\mathbf{M}_R = \sum_{k=1}^n \mathbf{K}_k^T \mathbf{L}_i \mathbf{L}_i^T \mathbf{K}_k$ corresponding to the largest g eigenvalues.
5. Set $i \leftarrow i + 1$, and goto step 3 until convergence.

4 Approximate A-KPCA

In the above section, we have derived an iterative procedure to compute the matrices \mathbf{L} and \mathbf{R} . In this section, we will develop a non-iterative algorithm to approximate A-KPCA in a more efficient way.

Given training set $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, suppose that the kernel matrices \mathbf{K}_k 's have been computed using Eq. 6. Then, instead of optimizing Eq. 5 directly, we adopt the alternative criterion shown in Eq. 8.

$$J(\mathbf{L}, \mathbf{R}) = \sum_{k=1}^n \|\mathbf{L}^T \mathbf{K}_k\|_F^2 + \sum_{k=1}^n \|\mathbf{K}_k \mathbf{R}\|_F^2 \quad (8)$$

where $\mathbf{L} \in \mathbb{R}^{n \times d}$, $\mathbf{L}^T \mathbf{L} = \mathbf{I}$ and $\mathbf{R} \in \mathbb{R}^{f \times g}$, $\mathbf{R}^T \mathbf{R} = \mathbf{I}$. Comparing Eq. 8 with Eq. 5, we can find that \mathbf{L} and \mathbf{R} are not tangled together any more, and therefore can be computed analytically and simultaneously.

Theorem 2 presents a closed form solution to Eq. 8, i.e. the approximate A-KPCA algorithm, as shown in **Algorithm 2**. The proof of the theorem is similar to that of PCA [4] and therefore we omit it due to the page limit.

Theorem 2 *The problem in Eq. 8 has closed form solution \mathbf{L} and \mathbf{R} : (1) \mathbf{L} consists of the d eigenvectors of the matrix $\mathbf{M}_L = \sum_{k=1}^n \mathbf{K}_k \mathbf{K}_k^T$ corresponding to the largest d eigenvalues; (2) \mathbf{R} consists of the g eigenvectors of the matrix $\mathbf{M}_R = \sum_{k=1}^n \mathbf{K}_k^T \mathbf{K}_k$ corresponding to the largest g eigenvalues.*

5 Experiments

Experiments are conducted to compare the proposed algorithms, A-KPCA and approximate A-KPCA, with traditional PCA and KPCA with single kernel on an toy exam-

Algorithm 2: Approximate A-KPCA**Input:** Training set $X = \{x_1, x_2, \dots, x_n\}$.**Output:** Matrices L and R .

1. Construct the kernel matrix K_k for each x_k using Eq. 6.
2. Compute the d eigenvectors $L = (l_1, l_2, \dots, l_d)$ of $M_L = \sum_{k=1}^n K_k K_k^T$ corresponding to the largest d eigenvalues.
3. Compute the g eigenvectors $R = (r_1, r_2, \dots, r_g)$ of $M_R = \sum_{k=1}^n K_k^T K_k$ corresponding to the largest g eigenvalues.

ple (see Figure 1), four UCI datasets [1] (*Ionosphere*, *Balance*, *Glass* and *Heart*), and the COIL-20 database for object recognition [8]. For each UCI dataset, half of the data are used for training, and the rest for testing. For COIL-20 database the first 12 images per object are used for training and the rest for testing.

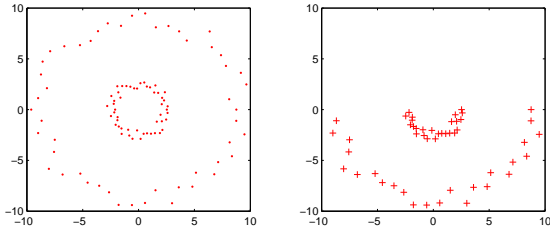


Figure 1. Left: The whole 2D toy dataset, with 2 classes (corresponding to the outer/inner circles respectively). Right: the training set.

In the experiments we adopted the Gaussian kernel function $\mathcal{K}^i(x, y) = \exp(-\|x - y\|^2 / 2\sigma_i^2)$ with width $\sigma_i = \sigma_0 \times i$, $i = 1, 2, \dots, 5$, where σ_0 is the standard variation of training data. We denote KPCA with the above 5 different kernels as KPCA-1 to KPCA-5, respectively. In A-KPCA and approximate A-KPCA, the value of g (number of columns in R) is set to 1. After feature extraction, the well-known nearest neighbor (1-NN) classifier is used for classification.

Figure 2 compares the classification accuracies of PCA, KPCA (with different kernels) and A-KPCA under different dimensions. It can be seen from Figure 2 that in most cases, A-KPCA outperforms PCA and KPCA-1 to KPCA-5 greatly. It is impressive that on almost all data sets the proposed A-KPCA consistently outperforms PCA and KPCAs no matter which dimension is considered. A-KPCA achieves the best performance on all the data sets when a relatively large dimension is used.

Figure 2 also indicates that among KPCAs with single kernel, i.e. KPCA-1 to KPCA-5, none could always achieve

the highest accuracy. For example, when large dimension is used, KPCA-1 outperforms KPCA-2 to KPCA-5 on *Ionosphere*, while on *Glass* KPCA-1 is with the lowest accuracy among KPCA-1 to KPCA-5. Thus, it leaves space for ensembling KPCA with different kernels for better performance.

Table 1 compares the average accuracies under different dimensions. From Table 1 we can find that A-KPCA outperforms PCA and all the five KPCAs on the averaged accuracies (see the bold). It is interesting to note, from Table 1, that no single KPCA always achieves the best performance (see underlined), which verifies our former claim again.

6 Conclusion

This paper proposes an algorithm for adaptive kernel principal component analysis, which is based on a new criterion enabling nonlinear feature extraction and unsupervised kernel selection be performed simultaneously. An iterative algorithm is proposed to optimize the criterion. Besides, a non-iterative approximate algorithm is presented to improve the efficiency. Experiments validate the superiority of the proposed algorithms.

In the reported experiments we only use Gaussian kernel as the base kernel for Eq. 6. We will investigate the performance of the proposed algorithms using different kinds of base kernels simultaneously in future work. Also, comparison with other kernel tuning algorithms will be carried out in the future. Moreover, it is possible to extend the proposed unsupervised kernel learning algorithms to other kernel-based methods, such as kernel k -means clustering.

Acknowledgments: Supported by the National Science Foundation of China (60505004, 60473035), the Jiangsu Science Foundation (BK2004001, BK2006521), and the National Science Fund for Distinguished Young Scholars of China (60325207).

References

- [1] C. Blake, E. Keogh, and C. J. Merz. UCI repository of machine learning databases. [http://www.ics.uci.edu/~mllearn/MLRepository.html], Department of Information and Computer Science, University of California, Irvine, CA, 1998.
- [2] O. Chapelle, V. Vapnik, O. Bousquet, and S. Mukherjee. Choosing multiple parameters for support vector machines. *Machine Learning*, 46(1):131–159, 2002.
- [3] N. Cristianini, J. Shawe-Taylor, A. Elisseeff, and J. Kandola. On kernel-target alignment. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages 367–373. MIT Press, Cambridge, MA, 2002.
- [4] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. Wiley and Sons, New York, 1973.

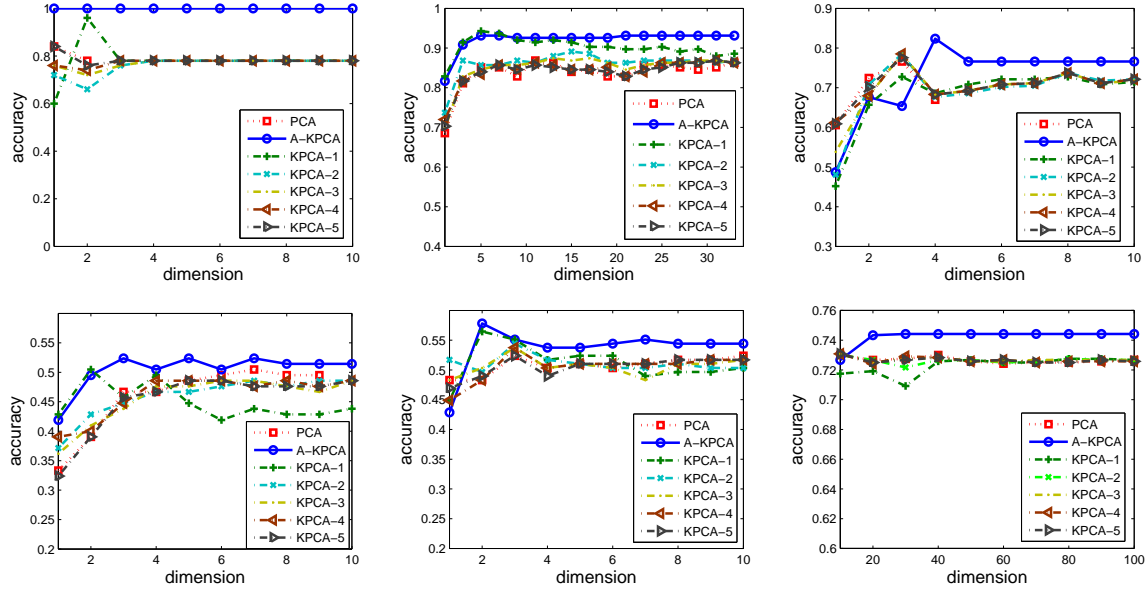


Figure 2. Comparisons on *Toy* (top left), *Ionosphere* (top middle), *Balance* (top right), *Glass* (bottom left), *Heart* (bottom middle) and *COIL-20* (bottom right).

Table 1. Comparisons of averaged accuracies (%) under different dimensions

Methods	Ionosphere	Balance	Glass	Heart	COIL-20
PCA	83.70	69.15	49.53	50.71	72.67
KPCA-1	<u>90.29</u>	68.24	44.86	50.76	72.00
KPCA-2	86.12	69.20	45.91	50.97	72.63
KPCA-3	85.04	69.65	45.62	<u>51.18</u>	<u>72.70</u>
KPCA-4	84.27	70.42	<u>46.19</u>	50.92	72.67
KPCA-5	84.03	<u>70.51</u>	45.24	50.86	72.66
A-KPCA	92.13	72.37	50.38	53.53	74.24
Approximate A-KPCA	92.20	72.27	49.24	53.58	74.22

[5] J. T. Kwok and I. W. Tsang. Learning with idealized kernels. In *Proceedings of the 20th International Conference on Machine Learning*, pages 400–407, Washington, DC, 2003.

[6] C. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, 5:27–72, 2004.

[7] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf. An introduction to kernel-based learning algorithms. *IEEE Transactions on Neural Networks*, 12(2):181–201, 2001.

[8] S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (COIL-20) [<http://www1.cs.columbia.edu/cave/research/softlib/coil-20.html>]. Technical Report CUCS-005-96, Columbia University, New York, NY, 1996.

[9] C. S. Ong, A. J. Smola, and R. C. Williamson. Hyperkernels. In S. T. S. Becker and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 478–485. MIT Press, Cambridge, MA, 2003.

[10] B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigen-value problem. *Neural Computation*, 10(5):1299–1319, 1998.

[11] I. W. Tsang and J. T. Kwok. Efficient hyperkernel learning using second-order cone programming. In *Proceedings of the 15th European Conference on Machine Learning*, pages 453–464, Pisa, Italy, 2004.

[12] S. Yang, S. Yan, D. Xu, X. Tang, and C. Zhang. Fisher + Kernel criterion for discriminant analysis. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 197–202, San Diego, CA, 2005.

[13] J. Ye. Generalized low rank approximations of matrices. In *Proceedings of the 21st International Conference on Machine Learning*, pages 887–894, Banff, Alberta, Canada, 2004.

[14] D. Q. Zhang, S. C. Chen, and Z.-H. Zhou. Learning the kernel parameters in kernel minimum distance classifier. *Pattern Recognition*, 39(1):133–135, 2006.