

1. Summary

- In solving complex visual learning tasks, we establish an approach in which **multiple kernel learning** (MKL) is incorporated into the training process of **dimensionality reduction** (DR) methods.
- Our approach
 - Based on MKL, data described by various descriptors are jointly considered, and projected into a unified space of low dimension.
 - Built on **graph embedding**, any DR methods explainable by graph embedding can be generalized by our approach.
 - Via integrating with different DR methods, MKL can address not only supervised learning problems but also semi-supervised and unsupervised ones.

2. Background

- Multiple kernel learning and graph embedding are two important components in the establishment of the approach.

2.1 Multiple Kernel Learning

- In complex vision applications, adopting multiple descriptors to characterize data is a feasible way for improving performances.
- Kernel matrix as a unified feature representation:
 - Represent the pairwise relationships among data under each descriptor by a kernel matrix.
- Multiple kernel learning: Learning an optimal ensemble kernel over a given convex set of base kernels:

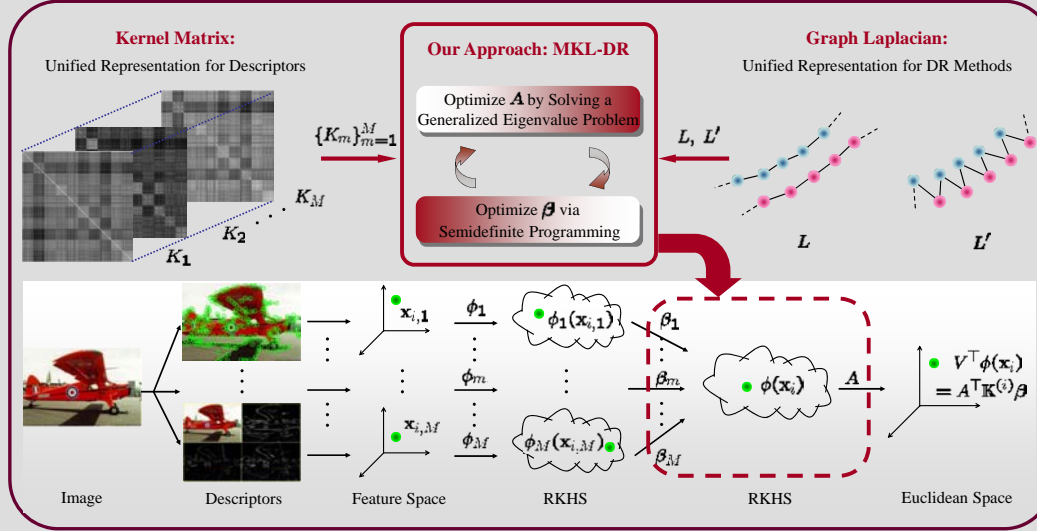
$$K = \sum_{m=1}^M \beta_m K_m, \beta_m \geq 0 \quad (1)$$

$$\text{or } k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^M \beta_m k_m(\mathbf{x}_i, \mathbf{x}_j), \beta_m \geq 0 \quad (2)$$
- Multiple kernel learning for finding optimal coefficients $\{\beta_m\}_{m=1}^M$ can be interpreted as descriptor combination.

2.2 Graph Embedding

- Many DR methods focus on modeling pairwise relationships among data points via graph structures.
- The framework of graph embedding [Yan et al. PAMI07] provides a unified formulation for a set of graph-based DR methods:

$$\mathbf{v}^* = \arg \min_{\mathbf{v}^T D D^T \mathbf{v} = 1, \text{ or } \mathbf{v}^T X L X^T \mathbf{v} = 1} \mathbf{v}^T X L X^T \mathbf{v}, \quad (3)$$
 where $X = [\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_N]$ is the data matrix, L and L' are graph Laplacian, D is a diagonal matrix, \mathbf{v}^* is the optimized linear embedding.
- By specifying particular L and L' (or L and D), a set of DR methods can be expressed by (3), such as PCA, LPP, LDA, LDE, or SDA.
- These DR methods include supervised, semi-supervised, and unsupervised ones.



3. Problem Definition

- Our goal is to make DR methods that can be expressed by (3) to consider multiple base kernels simultaneously, such that
 - Data characteristics captured in different descriptors can be jointly utilized to accomplish the objectives of these DR methods.
 - The diversity in the objectives of these DR methods enhances the applicability of MKL.
 - We could provide more prior knowledge to benefit the analysis of given data.
- With some derivation, sample i in the projected space can be expressed as

$$\mathbf{v}^T \mathbf{x}_i = \alpha^T \mathbf{K}^{(i)} \beta \in \mathbb{R},$$

where $\alpha = [\alpha_1 \alpha_2 \dots \alpha_N]^T$ is the **sample coefficient vector**,

$\beta = [\beta_1 \beta_2 \dots \beta_M]^T$ is the **kernel weight vector**,

$$\mathbf{K}^{(i)} = [K_1^{(i)} K_2^{(i)} \dots K_M^{(i)}] \in \mathbb{R}^{N \times M}.$$

- The resulting constrained optimization problem is

$$\min_{\mathbf{A}, \beta} \sum_{i,j=1}^N \|\mathbf{A}^T \mathbf{K}^{(i)} \beta - \mathbf{A}^T \mathbf{K}^{(j)} \beta\|^2 w_{ij} \quad (4)$$

$$\text{subject to } \sum_{i,j=1}^N \|\mathbf{A}^T \mathbf{K}^{(i)} \beta - \mathbf{A}^T \mathbf{K}^{(j)} \beta\|^2 w'_{ij} = 1,$$

$$\beta_m \geq 0, \text{ for } m = 1, 2, \dots, M,$$

where $W = [w_{ij}]$ is the affinity matrix of L in (3),

$W' = [w'_{ij}]$ is the affinity matrix of L' in (3),

$$\mathbf{A} = [\alpha_1 \alpha_2 \dots \alpha_P].$$

4. Optimization

- An **alternative optimization** procedure is used to solve (4).
 - Variable \mathbf{A} is optimized by solving a **generalized eigenvalue problem**.
 - Variable β is optimized via **semidefinite programming**.

4.1 On Optimizing \mathbf{A}

- By fixing β , the optimization problem (4) can be expressed as

$$\min_{\mathbf{A}} \text{trace}(\mathbf{A}^T \mathbf{S}_W^{\beta} \mathbf{A})$$

$$\text{subject to } \text{trace}(\mathbf{A}^T \mathbf{S}_{W'}^{\beta} \mathbf{A}) = 1$$

where \mathbf{S}_W^{β} and $\mathbf{S}_{W'}^{\beta}$ are two fixed square matrices.

- The columns of the optimized \mathbf{A} can be obtained by solving a generalized eigenvalue problem.

4.2 On Optimizing β

- By fixing \mathbf{A} , (4) becomes a **nonconvex QCQP** with respect to β .
- Thus, we instead consider its **SDP relaxation**:

$$\min_{\beta, B} \text{trace}(\mathbf{S}_W^{\beta} B)$$

$$\text{subject to } \text{trace}(\mathbf{S}_{W'}^{\beta} B) = 1,$$

$$\beta_m \geq 0, \text{ for } m = 1, 2, \dots, M,$$

$$\begin{bmatrix} 1 & \beta^T \\ \beta & B \end{bmatrix} \succeq 0,$$

where \mathbf{S}_W^{β} and $\mathbf{S}_{W'}^{\beta}$ are two fixed square matrices.

- The optimal β can be obtained via semidefinite programming.

5. Experimental Results

- The Caltech 101 image dataset is used in the experiments.
 - It consists of 101 object categories and one additional class of background images.
 - All 102 classes are used in the application to object recognition
 - 30 classes are chosen in the application to image clustering
- We implement seven kinds of image descriptors that result in the following seven base kernel matrices:
 - GB-1 / GB-2: Based on geometric blur descriptor.
 - SIFT-Dist / SIFT-Grid: Based on the SIFT descriptor.
 - C2-SWP / C2-ML: Based on biologically inspired features.
 - PHOG: Based on PHOG descriptor.

5.1 Supervised Application to Recognition

- We adopt two supervised DR schemes for the applications:
 - Linear discriminant analysis (LDA)**: Assume data of a class can be modeled by a Gaussian.
 - Local discriminant embedding (LDE)** [Chen et al. CVPR05]: Assume data of a class spread as a submanifold.
- To express LDA and LDE in the form of (3), we need specify their corresponding graph Laplacian matrices, i.e., L and L' .
- The recognition rates are listed as follows

Table 1: Recognition rates (mean \pm std %) for Caltech-101 dataset

kernel(s)	method	number of classes		method	number of classes	
		102	101		102	101
GB-1	KFD	57.3 \pm 2.5	57.7 \pm 0.7	KLDE	57.1 \pm 1.4	57.7 \pm 0.8
GB-2		60.0 \pm 1.5	60.6 \pm 1.5		60.9 \pm 1.4	61.3 \pm 2.1
SIFT-Dist		53.0 \pm 1.4	53.2 \pm 0.8		54.2 \pm 0.5	54.6 \pm 1.5
SIFT-Grid		48.8 \pm 1.9	49.6 \pm 0.7		49.5 \pm 1.3	50.1 \pm 0.3
C2-SWP		30.3 \pm 1.2	30.7 \pm 1.5		31.1 \pm 1.5	31.3 \pm 0.7
C2-ML		46.0 \pm 0.6	46.8 \pm 0.9		45.8 \pm 0.2	46.7 \pm 1.5
PHOG		41.8 \pm 0.6	42.1 \pm 1.3		42.2 \pm 0.6	42.6 \pm 1.3
-	KFD- <i>fixing</i>	68.4 \pm 1.5	68.9 \pm 0.3	KLDE- <i>fixing</i>	68.4 \pm 1.4	68.7 \pm 0.8
-	KFD-SAMME	71.2 \pm 1.4	72.1 \pm 0.7	KLDE-SAMME	71.1 \pm 1.9	71.3 \pm 1.2
All	MKL-LDA	74.6 \pm 2.2	75.3 \pm 1.7	MKL-LDE	75.3 \pm 1.5	75.8 \pm 1.7

5.2 Unsupervised Application to Clustering

- We apply our approach to **locality preserving projections (LPP)** [He & Niyogi NIPS03] to serve as a data preprocessing tool.
- Data in various representations are projected into a unified space.
- Consider clustering methods affinity propagation and k -means.

Table 2: Clustering performance (NMI / ERR %) on the 20-class image dataset

kernel(s)	preprocessing method	affinity propagation		k -means clustering	
		without data preprocessing	with data preprocessing	without data preprocessing	with data preprocessing
GB-1	kernel LPP	0.553 / 50.8	0.609 / 38.3	-	0.611 / 44.7
GB-2		0.577 / 48.0	0.624 / 43.7	-	0.640 / 43.0
SIFT-Dist		0.627 / 43.7	0.651 / 31.0	-	0.671 / 41.7
SIFT-Grid		0.598 / 41.3	0.631 / 45.7	-	0.642 / 43.0
C2-SWP		0.383 / 70.3	0.379 / 60.5	0.383 / 68.5	0.379 / 66.7
C2-ML		0.499 / 54.5	0.488 / 56.0	0.525 / 53.0	0.507 / 55.8
PHOG		0.455 / 57.3	0.482 / 52.7	-	0.510 / 52.7
All	MKL-LPP	-	0.714 / 25.0	-	0.758 / 25.7