

# Kernelized View Adaptive Subspace Learning for Person Re-identification

Qin Zhou<sup>12</sup>  
zhou.qin.190@sjtu.edu.cn  
Shibao Zheng<sup>12</sup>  
sbzh@sjtu.edu.cn  
Hang Su<sup>34</sup>  
suhangss@gmail.com  
Hua Yang<sup>12</sup>  
hyang@sjtu.edu.cn  
Yu Wang<sup>12</sup>  
txtxs@sjtu.edu.cn  
Shuang Wu<sup>12</sup>  
shuangwu@sjtu.edu.cn

<sup>1</sup> Department of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University  
<sup>2</sup> Shanghai Key Lab of Digital Media Processing and Transmission  
<sup>3</sup> Department of Computer Science and Technology, Tsinghua University  
<sup>4</sup> Robotics Institute, Carnegie Mellon University

Person re-identification refers to the task of recognizing the same person under different non-overlapping camera views and across different time and places. We propose a novel algorithm coined as Kernelized View Adaptive Subspace Learning (KVASL), which tries to learn different projection matrices for each camera view to compensate for the specific distortion brought by different views. The kernel trick is adopted to catch more information such that nonlinear transformation is possible. We present the motivating example in Figure 1, which not only demonstrates the benefit of adopting different projection matrices, but also illustrates the necessity of the kernel trick.

Denote  $|S|, |D|$  as the number of matched and mismatched image pairs in the training set. Our kernelized loss function is formulated as follows:

$$\begin{aligned} \ell_K(L_A^K, L_B^K) = & \frac{1}{|S|} \text{tr}(K_A^T K_A Q_A^T Q_A - 2K_B^T K_A Q_A^T Q_B + K_B^T K_B Q_B^T Q_B) \\ & + \mu_A \|L_A^K\|_F^2 + \mu_B \|L_B^K\|_F^2 - \frac{\lambda}{|D|} \text{tr}((|S| - 1)K_A^T K_A Q_A^T Q_A) \\ & + \text{tr}(2K_A^T X K_B Q_B^T Q_A - (|S| - 1)K_B^T K_B Q_B^T Q_B) \end{aligned} \quad (1)$$

where  $K_A = \phi(A_{train})^T \phi(A_{train}) \in R^{N \times N}$ ,  $K_B = \phi(B_{train})^T \phi(B_{train}) \in R^{N \times N}$  are symmetry matrices,  $\phi(A_{train}) = [\phi(x_{1A}^{train}), \dots, \phi(x_{NA}^{train})]$ ,  $\phi(B_{train}) = [\phi(x_{1B}^{train}), \dots, \phi(x_{NB}^{train})] \in R^{m \times N}$  are the matrices formed by feature vectors in the kernel space of the corresponding camera views,  $\text{tr}(\cdot)$  indicates the trace of the matrix and  $X$  corresponds to the matrix whose diagonal elements are all zeros and the other elements are all ones.

We adopt an alternately iterative gradient descent method to optimize our loss function. The gradients of  $Q_A, Q_B$  in Eq.(1) can be derived as follows:

$$\begin{aligned} \frac{\partial \ell}{\partial Q_A} = & \frac{2}{|S|} (Q_A K_A K_A - Q_B K_B K_B) \\ & - \frac{2\lambda}{|D|} [(|S| - 1)Q_A K_A K_A - Q_B K_B X K_A] + 2\mu_A Q_A K_A \\ \frac{\partial \ell}{\partial Q_B} = & \frac{2}{|S|} (Q_B K_B K_B - Q_A K_A K_B) \\ & - \frac{2\lambda}{|D|} [(|S| - 1)Q_B K_B K_B - Q_A K_A X K_B] + 2\mu_B Q_B K_B \end{aligned} \quad (2)$$

where  $X$  is the same matrix as in Eq.(1). In this case, the matrices to be learned can be updated in the following way:

$$Q_A(t+1) = Q_A(t) - \eta_{Q_A} \frac{\partial \ell}{\partial Q_A}, \quad Q_B(t+1) = Q_B(t) - \eta_{Q_B} \frac{\partial \ell}{\partial Q_B} \quad (3)$$

where  $\eta_{L_A}, \eta_{L_B}, \eta_{Q_A}, \eta_{Q_B}$  are the learning rate corresponding to each transform matrix. Once  $Q_A, Q_B$  are learned, the distance between two test images can be derived as follows:

$$\begin{aligned} d_K^2(x_{iA}^{test}, x_{jB}^{test}) = & e_i^T K_A^{test} Q_A^T Q_A (K_A^{test})^T e_i \\ & - 2e_j^T K_B^{test} Q_B^T Q_A (K_A^{test})^T e_i + e_j^T K_B^{test} Q_B^T Q_B (K_B^{test})^T e_j \end{aligned} \quad (4)$$

We implement our algorithm on four publicly available datasets, the comparison results with the baseline method demonstrate the superiority of view-adaptive projection matrices over using the same projection

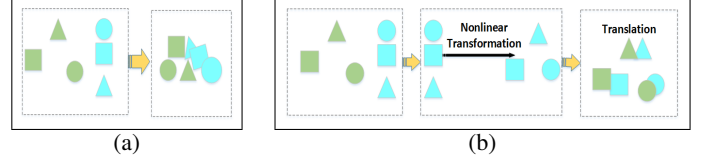


Figure 1: A conceptual illustration of how different projection functions can be more flexible than the same ones. (a) The same transformation; (b) Different transformations.

matrices. We also compare our algorithm with some state-of-the-art algorithms (SVMML [2], KISSME [1], kLFDA [3] et al.), which demonstrates the effectiveness of our algorithm. Detailed comparison results are shown in Figure 2 and Table 1.

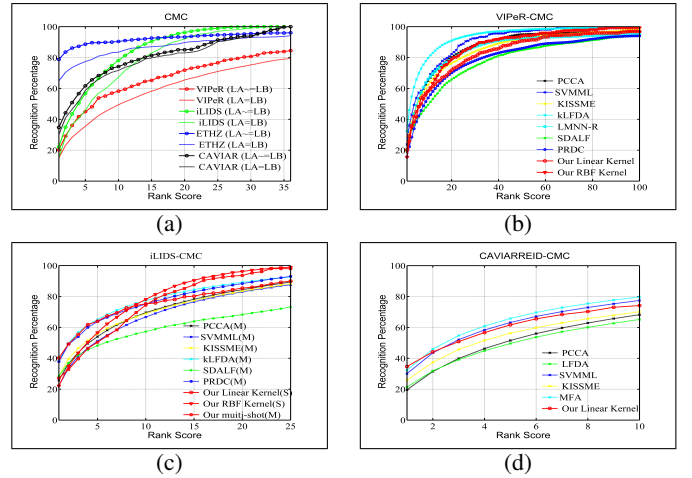


Figure 2: (a) Comparison results with the baseline method. (b) Results on VIPeR dataset. (c) Results on iLIDS dataset. (d) Results on CAVIAR4REID dataset. Red lines correspond to results of our algorithm.

Method	r=1	r=5	r=10	r=15
KISSME [1]	65.51	83.67	87.34	89.59
Our Linear Kernel	<b>78.77</b>	<b>88.63</b>	<b>90.41</b>	<b>92.60</b>
Our RBF Kernel	70.55	85.75	89.45	91.23

Table 1: Comparison of our algorithm with [1] on the ETHZ dataset.

- [1] Martin Köstinger, Martin Hirzer, Paul Wohlhart, Peter M. Roth, and Horst Bischof. Large scale metric learning from equivalence constraints. In *Proc. CVPR*, pages 2288–2295, 2012.
- [2] Zhen Li, Shiyu Chang, Feng Liang, Thomas S. Huang, Liangliang Cao, and John R. Smith. Learning locally-adaptive decision functions for person verification. In *Proc. CVPR*, pages 3610–3617, 2013.
- [3] Fei Xiong, Mengran Gou, Octavia I. Camps, and Mario Szaier. Person re-identification using kernel-based metric learning methods. In *Proc. ECCV*, pages 1–16, 2014.