# TLR: TRANSFER LATENT REPRESENTATION FOR UNSUPERVISED DOMAIN ADAPTATION

*Pan Xiao[1], Bo Du\*[1], Jia Wu[2], Lefei Zhang[1], Ruimin Hu[1] and Xuelong Li[3]*

[1]School of Computer, Wuhan University, Wuhan 430072, Hubei, China
[2]Department of Computing, Macquarie University, Sydney, NSW 2109, Australia
[3]Xi'an Institute of Optics and Precision Mechanics,
Chinese Academy of Sciences, Xi'an 710119, Shaanxi, China

## ABSTRACT

Domain adaptation refers to the process of learning prediction models in a target domain by making use of data from a source domain. Many classic methods solve the domain adaptation problem by establishing a common latent space, which may cause the loss of many important properties across both domains. In this manuscript, we develop a novel method, *transfer latent representation* (TLR), to learn a better latent space. Specifically, we design an objective function based on a simple linear autoencoder to derive the latent representations of both domains. The encoder in the autoencoder aims to project the data of both domains into a robust latent space. Besides, the decoder imposes an additional constraint to reconstruct the original data, which can preserve the common properties of both domains and reduce the noise that causes domain shift. Experiments on cross-domain tasks demonstrate the advantages of TLR over competing methods.

***Index Terms*—** Domain adaptation, linear autoencoder, object and action recognition

## 1. INTRODUCTION

Recently, online images and videos grow exponentially, which has created a strong demand for technologies to analyze the multimedia content. Unfortunately, labels for these new visual images are in short supply and it is nearly impossible to learn a good visual category model without enough labels. In real-world applications, there exist many labeled datasets in some old domains. Can we use these labeled datasets (i.e. the source domain) to handle unlabeled datasets (the target domain)? To answer this question, a technique named domain adaptation (DA) has been developed.

DA is very important when the labels for target domain data are lacking [1]. For example, we can obtain some labeled
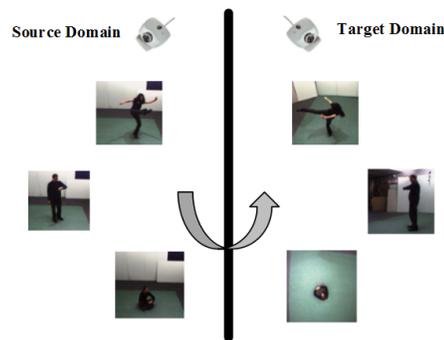
**Fig. 1**. Cross-camera action recognition.

images drawn from the Internet (i.e. the source domain) and some unlabeled images captured by cameras (the target domain), and that both domains contain the same objects. It is believed that a model trained in the source domain can significantly improve classification accuracy in the target domain after the common properties of both domains are extracted [2]. Taking action recognition by surveillance cameras in Fig.1 as an example, we have a series of action shots captured from two different angles. If we were to directly use a set of labeled images on the left to classify unlabeled pictures on the right, the classification accuracy might be unsatisfactory. However, considering that both sets of pictures contain the same set of actions, we believe a better prediction result can be obtained by recognizing and utilizing the commonalities between the image sets in classification.

Domain adaptation methods can be divided into two categories: semi-supervised DA and unsupervised DA according to the availability of labeled instances in the target domain. In this work, we focus on the unsupervised scenario, which is hard to solve since the labels for the target domain are totally non-existent. Many well-known methods have been proposed to solve the unsupervised domain adaptation problem. One straightforward solution is to project both domains into a common latent space. For example, Fernando *et al.* [3] proposed to learn a linear projection aligning the source

and target domains. Gong *et al.* [4] claimed that new latent representations could be obtained by regarding the subspaces of both domains as points in Grassmann manifolds. Pan *et al.* [5] and Yan *et al.* [6] projected the source and target domain data into a Reproducing Kernel Hilbert Space (RKHS) to obtain the latent representations of both domains. Although all these state-of-the-art methods have achieved promising results, there is still room for improvement, mainly because those methods may result in the loss of many important properties of both domains that are helpful for model building when projection is performed.

In this paper, we propose a new method called Transfer Latent Representation (TLR) to learn a better latent space. Specifically, we first follow the procedure outlined in [5] to obtain linearly separable source and target domain data by projecting both domains into an RKHS. To avoid the loss of useful properties, we then design an objective function based on a linear autoencoder to derive the latent representations of both domains. The encoder of the autoencoder is set up to project both domains into a latent space, in the same way as the existing domain adaptation methods. Besides, the decoder exerts an additional constraint, that is, the original data must be reconstructed by the projection. It is supposed that the use of this additional reconstruction constraint can assist in preserving the common properties of both domains and reducing the noise that causes domain shift. The Maximum Mean Discrepancy (MMD) [7] between the latent representations is also integrated into the objective function, so that the function is able to further narrow the distance between different domain distributions. Finally, we obtain the latent representations of both domains in a latent space.

## 2. PRELIMINARIES

In this work, we aim to solve the unsupervised domain adaptation problem: how to best label the unlabeled target domain data in an unsupervised manner by training a model on labeled data in a relevant source domain. Firstly, we denote $X_S = \{x_{S_1}, ..., x_{S_{n_1}}\} \in \mathbb{R}^{d \times n_1}$ as the source domain data and $X_T = \{x_{T_1}, ..., x_{T_{n_2}}\} \in \mathbb{R}^{d \times n_2}$ as the target domain data. Here, $d$ is the dimension of each instance, while $n_1$ and $n_2$ are the number of samples in the source and target domains respectively. The source domain data labels are denoted as $Y_S = \{y_{S_1}, ..., y_{S_{n_1}}\} \in \mathbb{R}^{n_1}$, where $y_{S_i}$ is the label of the corresponding source domain sample $x_{S_i}$. Similarly, the predicted labels of the target domain are denoted as $\hat{Y}_T = \{y_{T_1}, ..., y_{T_{n_2}}\} \in \mathbb{R}^{n_2}$. Our goal is to train a classifier based on $X_S$ and $\hat{Y}_S$, then predict the labels of the target domain data as accurately as possible.

### 2.1. Maximum Mean Discrepancy

*Maximum Mean Discrepancy* (MMD) has been successfully used to solve the domain adaptation problem [8, 9]. By com-

puting on $X_S$ and $X_T$, a non-parametric distance estimate between domain distributions can be directly obtained. Here, let

$$MMD(X_S, X_T) = \left\| \frac{1}{n_1} \sum_{i=1}^{n_1} f(x_{S_i}) - \frac{1}{n_2} \sum_{i=1}^{n_2} f(x_{T_i}) \right\|_{\mathcal{H}}^2 \tag{1}$$

where $\mathcal{H}$ is a universal Reproducing Kernel Hilbert Space (RKHS), and $f : \mathcal{X} \to \mathcal{H}$ denotes the non-linear transformation. By means of the kernel trick, (i.e., $k(x_i, x_j) = f(x_i)f(x_j)'$), we can rewrite (1) as

$$MMD(X_S, X_T) = tr(KL). \tag{2}$$

in which

$$K = \begin{bmatrix} K_{S,S} & K_{S,T} \\ K_{T,S} & K_{T,T} \end{bmatrix} = \begin{bmatrix} H_S \\ H_T \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times (n_1+n_2)} \tag{3}$$

is a symmetric kernel matrix. The elements of $K_{S,S}$, $K_{T,T}$, $K_{T,S}$ and $K_{S,T}$ are the values of $k(x_i, x_j)$ when $(x_i, x_j)$ belongs to the source domain, target domain, and two cross domains respectively. $H_S$ and $H_T$ denote the source and target domain samples mapped into the RKHS. $L$ is the MMD matrix and can be described as follows:

$$(L)_{i,j} = \begin{cases} \frac{1}{n_1 n_1}, & x_i, x_j \in X_S \\ \frac{1}{n_2 n_2}, & x_i, x_j \in X_T \\ \frac{-1}{n_1 n_2}, & otherwise \end{cases} \tag{4}$$

## 3. TRANSFER LATENT REPRESENTATION

The proposed model consists of two main stages. First, we map the data of both domains into the RKHS and obtain $H_S$ and $H_T$ according to (3). We then derive a matrix $W$ based on the simple linear autoencoder, which projects $H_S$ and $H_T$ into a latent space.

### 3.1. Simple Linear Autoencoder

The simplest form of an autoencoder is linear. Here, there is no activation function in the single hidden layer. The encoder is used to project the input data into the single hidden layer, while the decoder projects it back to the original feature space. Suppose that $X \in \mathbb{R}^{n \times (n_1+n_2)}$ is the input data matrix with $n$ samples. We want to obtain a projection matrix $W \in \mathbb{R}^{(n_1+n_2) \times k}$ in order to explore the $k$-dimensional latent representation $P \in \mathbb{R}^{n \times k}$. The obtained $P$ is projected back to the original feature space by means of transpose matrix of $W$ so that it becomes $\hat{X} \in \mathbb{R}^{n \times (n_1+n_2)}$. Note here that $k$ is smaller than $(n_1 + n_2)$. By minimizing the reconstruction error, we have

$$\min_{W, W^\top} \left\| PW^\top - X \right\|_F^2 \quad s.t. \quad P = XW. \tag{5}$$

The above objective makes $X$ and $\hat{X}$ as similar as possible.

## 3.2. Model Formulation

We apply the simple linear autoencoder to the source domain in RKHS ($H_S$), and the target domain in RKHS ($H_T$). One significant advantage of the autoencoder is that it can reconstruct the input features of the source and target domains, which forces the latent representations of both domains to maintain as many important properties as possible. Formally, we have

$$\min_W \left\| P_S W^\top - H_S \right\|_F^2 \quad s.t. \quad P_S = H_S W. \qquad (6)$$

$$\min_W \left\| P_T W^\top - H_T \right\|_F^2 \quad s.t. \quad P_T = H_T W. \qquad (7)$$

where $P_S$ and $P_T$ are the latent representations for the source and target domains respectively.

To further narrow the distance between the distributions of both domains, we minimize the MMD of the two latent representations ($P_S$ and $P_T$). More specifically, we have,

$$\min_W MMD(P_S, P_T). \qquad (8)$$

By combining (6) and (7) with (8), the proposed model can be summarized as follows:

$$\min_W F(W) = MMD(P_S, P_T) + \alpha \left\| P_S W^\top - H_S \right\|_F^2$$
$$+ \beta \left\| P_T W^\top - H_T \right\|_F^2. \qquad (9)$$

where $\alpha$ and $\beta$ are trade-off parameters.

At this point, our goal is to obtain the optimal $W$ that will minimize the objective function (9). The derived $W$ can project $H_S$ and $H_T$ into a common latent space. It is believed that, in the latent space, the noise that causes domain shift will be reduced and the common properties of different domains will be extracted.

## 3.3. Optimization

In this section, we introduce three propositions in turn. An efficient optimization algorithm is then designed.

**Proposition 1.** *The term (8) can be rewritten as*

$$\min_W tr(W^\top KLKW). \qquad (10)$$

*Proof.* According to (2), we have

$$MMD(P_S, P_T) = tr(K_{\mathcal{H}} L). \qquad (11)$$

where

$$K_{\mathcal{H}} = \begin{bmatrix} H_S WW^\top H_S^\top & H_S WW^\top H_T^\top \\ H_T WW^\top H_S^\top & H_T WW^\top H_T^\top \end{bmatrix}$$
$$= \begin{bmatrix} H_S \\ H_T \end{bmatrix} WW^\top \begin{bmatrix} H_S^\top & H_T^\top \end{bmatrix}$$
$$= KWW^\top K^\top. \qquad (12)$$

It is worth noting that we use $k(x_i, x_j) = x_i x_j'$ directly, since the source and target domain data have been mapped into the RKHS.

By substituting equation (12) into (11), the MMD of both domains in the latent space can be described as

$$MMD(P_S, P_T) = tr(KWW^\top K^\top L)$$
$$= tr(W^\top K^\top LKW). \qquad (13)$$

Since $K$ is a symmetric matrix, we have

$$MMD(P_S, P_T) = tr(W^\top KLKW). \qquad (14)$$

Finally, we obtain an equivalent problem (10). □

By combining (14) with (9), we can summarize our objective function as

$$\min_W F(W) = tr(W^\top KLKW) + \alpha \left\| H_S WW^\top - H_S \right\|_F^2$$
$$+ \beta \left\| H_T WW^\top - H_T \right\|_F^2. \qquad (15)$$

**Proposition 2.** *The objective function (15) can be rewritten more compactly as*

$$\min_W tr(WW^\top AWW^\top + W^\top BW - 2W^\top AW + A). \qquad (16)$$

*where*

$$A = KMK, B = KLK. \qquad (17)$$

*and*

$$M = \begin{bmatrix} \alpha I_{n_1 \times n_1} & 0_{n_1 \times n_2} \\ 0_{n_2 \times n_1} & \beta I_{n_2 \times n_2} \end{bmatrix}. \qquad (18)$$

The proof is given in the **Appendix**.

Here, it is evident that the solution may collapse to one point ($W = 0$). To avoid such an occurrence, we impose a constraint $W^\top AW = I$ into (16). Accordingly, the optimization problem with constraint can be summarized as follows:

$$\min_W \quad tr(W^\top W + W^\top BW)$$
$$s.t. \quad W^\top AW = I. \qquad (19)$$

We can solve problem (19) efficiently using the Lagrangian multiplier method, so that it eventually becomes the following optimization problem:

**Proposition 3.** *Problem (19) can be rewritten as*

$$\max_W tr((W^\top(I + B)W)^{-1} W^\top AW). \qquad (20)$$

*Proof.* The Lagrangian function of (19) is

$$L(W, Z) = tr(W^\top(I + B)W) - tr((W^\top AW - I)Z). \qquad (21)$$

where $Z$ is a matrix with the Lagrange multipliers in the diagonal. Setting $\frac{\partial L(W,Z)}{\partial W} = 0$, we have

$$(I + B)W = AWZ. \qquad (22)$$

To simplify the Lagrangian function $L(W, Z)$, we multiply both sides of equation (22) on the left by $W^\top$ and combine it with (21), so that we have

$$\min_W tr((W^\top AW)^{-1} W^\top (I + B)W). \qquad (23)$$

As the matrix $I + B$ is non-singular, an equivalent trace maximization problem (20) can be obtained. □

The solution of $W$ in (20) is the eigenvectors corresponding to the $k$ leading eigenvalues of $(I + B)^{-1}A$. The whole procedure of TLR is summarized in Algorithm 1.

---

**Algorithm 1** Transfer Latent Representation (TLR)

---

**Input:** Labeled source domain data $X_S$, unlabeled target domain data $X_T$, source labels $Y_S$, latent space dimension $k$, trade-off parameters $\alpha$ and $\beta$;
**Output:** Predicted target labels $\hat{Y}_T$;
  1: Compute matrices $K$, $H_S$, and $H_T$ according to (3);
  2: Compute matrices $L$ and $M$ using (4) and (18) respectively;
  3: Compute matrices $A$ and $B$ according to (17);
  4: Obtain projection matrix $W$ according to the $k$ leading eigenvalues of $(I + B)^{-1}A$;
  5: $P_S = H_S W$;
  6: $P_T = H_T W$;
  7: $\hat{Y}_T \leftarrow Classifier(P_S, P_T, Y_S )$;

---

## 4. EXPERIMENTS

To demonstrate the efficacy of our proposed TLR approach, we perform experiments on two cross-domain datasets: 1) the 4DA dataset, and 2) the IXMAS dataset.

### 4.1. Data Preparation

**4DA Dataset:** We adopt the public 4DA dataset [10], which contains four domains, namely **Amazon**, **Webcam**, **DSLR**, and **Caltech-256**. Fig.2 shows the MONITOR examples from the four domains. The differences between them are obvious. For example, the monitor screens from Amazon and Caltech-256 display colorful images while the monitor screens from Webcam and DSLR are black. For the experiments, we follow the procedure in [10] to extract SURF features. Each image corresponds to an 800-dimensional vector. The instances are then standardized by z-score. We randomly select two different domains from A (Amazon), W (Webcam), D (DSLR), and C (Caltech-256). Thus there are a total of $4 \times 3 = 12$ cross-domain pairs, e.g., $A \rightarrow W$, $A \rightarrow D$, $A \rightarrow C$,...,$C \rightarrow D$.
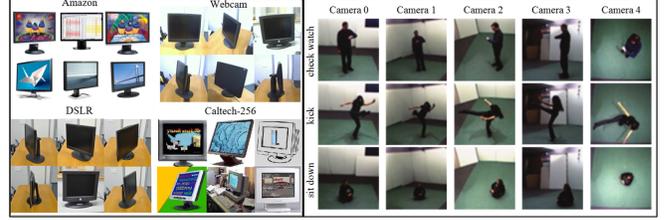


**Fig. 2**. Example images from the 4DA and IXMAS datasets.

**IXMAS Dataset:** The Inria Xmas Motion Acquisition Sequences (IXMAS)[1] is a multi-view action recognition dataset containing 11 actions. Each action is regarded as a category. As can be seen in Fig. 2, five cameras (cam0, cam1, ..., cam4) are used to capture the actions from different perspectives. Each perspective represents a domain. Thus, five domains are included in this dataset. Twelve actors are invited to perform each action three times, giving $12 \times 3 = 36$ instances per class. The feature extraction is based on the settings in [11]. We conduct experiments on 5 cross-domain pairs ($c0 \rightarrow c1$, $c1 \rightarrow c2$, ..., $c4 \rightarrow c0$).

### 4.2. Comparison methods

To evaluate the robustness of the proposed TLR approach, we compare TLR with six competitive methods: Principal Component Analysis (**PCA**), Information-Theoretical Learning (**ITL** [1]), Subspace Alignment (**SA** [3]), Transfer Component Analysis (**TCA** [5]), Geodesic Flow Kernel (**GFK** [4]), Maximum Independence Domain Adaptation (**MIDA** [6]). Following the settings in [4], 1-NN is chosen as the base classifier, where the source and target domains are regarded as the training set and test set respectively. This is so that we do not need to tune cross-validation parameters when training a model. We first compare our method with PCA, where both domains are mapped into their respective subspaces. In particular, ITL [1], SA [3], TCA [5], GFK [4], and MIDA [6] obtain a domain-invariant feature subspace in different ways. ITL optimizes an information-theoretic metric and learns the feature space discriminatively. SA learns a linear projection that aligns both domains using subspace alignment, while TCA maps data from the source and target domains into an RKHS in order to transfer components across domains. GFK extracts an infinite number of subspaces and constructs geodesic flows between them. Here, the subspaces of both domains are regarded as points in Grassmann manifolds. Finally, MIDA maximizes the independence of the derived and the instance features in order to reduce the difference between domains.

---

### 4.3. Implementation Details

In TLR, we need to tune two model parameters: the trade-off parameters $\alpha$ and $\beta$. Since the distributions of both domains are different, obtaining the optimal parameters by cross validation is impossible. We thus evaluate TLR by designing a search space according to [9]. The search range for $\alpha$ and $\beta$ is $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}\}$ separately. Similarly, the optimal dimension of latent space $k$ is obtained by searching $\{10, 20, ..., 200\}$. The best results are then reported. For the other comparison methods mentioned above, we tune the parameters according to the original paper and report their best performance.

**Table 1**. Classification accuracy (%) for all methods on the 4DA and IXMAS datasets.

| Methods | PCA | SA | ITL | TCA | GFK | MIDA | TLR |
|---|---|---|---|---|---|---|---|
| C→A | 32.5 | 31.4 | 35.4 | 37.6 | 35.8 | 37.3 | **38.7** |
| C→D | 26.3 | 33.0 | 33.1 | 35.0 | 35.7 | 35.3 | **38.1** |
| C→W | 25.1 | 26.9 | 28.0 | 31.9 | 31.0 | 31.8 | **34.4** |
| A→C | 31.9 | 32.2 | 34.9 | 34.9 | 33.8 | 34.6 | **35.2** |
| A→D | 25.7 | 28.4 | 30.0 | 30.1 | 33.2 | 29.7 | **34.8** |
| A→W | 28.7 | 28.8 | 29.6 | 32.6 | 33.0 | 32.6 | **35.1** |
| D→C | 27.7 | 30.8 | 31.6 | 31.1 | 27.8 | 30.7 | **32.2** |
| D→A | 31.0 | 31.8 | 34.1 | 34.2 | 31.5 | 33.9 | **35.1** |
| D→W | 59.5 | 78.7 | **78.8** | 75.3 | 69.1 | 75.5 | **78.8** |
| W→C | 25.5 | 25.1 | 27.8 | 29.7 | 28.9 | 29.7 | **30.6** |
| W→A | 31.2 | 30.0 | 32.4 | 30.3 | **33.7** | 31.4 | 30.0 |
| W→D | 68.1 | 81.1 | 82.4 | 80.1 | 78.2 | 81.5 | **86.3** |
| Average | 34.4 | 38.2 | 39.8 | 40.2 | 39.3 | 40.3 | **42.4** |
| c0→c1 | 8.6 | 14.7 | 15.4 | 30.0 | 14.2 | 25.1 | **35.2** |
| c1→c2 | 9.4 | 13.6 | 20.8 | 18.8 | 15.5 | 19.8 | **21.1** |
| c2→c3 | 10.5 | 9.0 | 13.6 | 12.3 | 8.9 | 12.4 | **19.9** |
| c3→c4 | 8.2 | 14.3 | 17.2 | 21.5 | 17.8 | 20.9 | **23.9** |
| c4→c0 | 13.5 | 14.2 | 16.5 | 24.6 | 16.3 | 24.3 | **27.5** |
| Average | 10.0 | 13.2 | 16.7 | 21.4 | 14.5 | 20.5 | **25.5** |

We follow the settings in [4] to select the training set and test set when conducting experiments on 4DA dataset. For the IXMAS dataset, we randomly select 30 labeled source domain instances per category as the training set and treat all target domain samples as the test set. We run experiments ten times at random for the 17 cross-domain image (object and action) pairs and the average classification accuracy is then reported in Table 1.

### 4.4. Experimental Results

The best result for each cross-domain pair is shown in bold. We observe that TLR outperforms all classic unsupervised domain adaptation methods. TLR's average classification accuracies on the 4DA and IXMAS datasets are 42.4% and 25.5% respectively, and the performance is improved by 2.1% and 4.1% relative to the best comparison method. This demonstrates that TLR can obtain more robust latent representations than its competitors when facing cross-domain recognition tasks.

Secondly, out of all methods studied, the classification results of PCA are the worst. This is because PCA is not designed to solve domain adaptation problem. SA's performance is slightly better than that of PCA, since SA adapts both domains in PCA subspaces, which further improves the classification accuracy.

Thirdly, TLR significantly outperforms ITL. The classification accuracy of ITL on IXMAS is not very good. A major limitation of ITL is that it assumes that data from the source domain and target domains are tightly clustered. This assumption may be invalid on many datasets.

Note that TCA, somewhat like TLR, also learns latent representations using MMD. However, the proposed method maintains more common properties between domains, as the input features can be reconstructed by means of the simple linear autoencoder.

The performance of GFK on the 4DA dataset is good, but poor on the IXMAS dataset. In GFK, the latent space dimension should be small enough to guarantee that subspaces transit smoothly along the geodesic flow. However, this may result in the loss of some important properties. TLR, on the other hand, can obtain a more accurate latent space.

Lastly, MIDA achieves better performance than the other compared algorithms. Theoretically, MIDA can learn features containing maximal independence with the domain features. However, one possible drawback of MIDA is that it retains fewer common properties than TLR does.
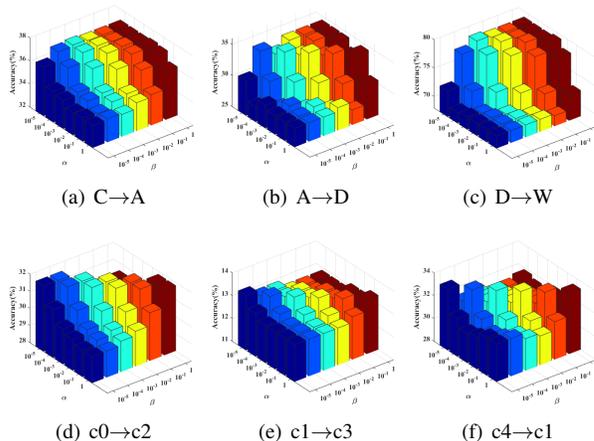
(a) C→A     (b) A→D     (c) D→W

(d) c0→c2     (e) c1→c3     (f) c4→c1

**Fig. 3**. Parameter sensitivity analysis for TLR.

### 4.5. Parameter Sensitivity Analysis

In the proposed TLR method, we need to tune two key parameters $\alpha$ and $\beta$. $\alpha$ represents how much we weight the source domain data, and $\beta$ denotes how much we weight the target domain data. To evaluate the effect of $\alpha$ and $\beta$ on the experimental results, we run the experiments on the 4DA dataset (three tasks: C→A, A→D, and D→W) and IXMAS dataset (another three tasks: c0→c2, c1→c3, and c4→c1)

with different parameter values. The two parameters $\alpha$ and $\beta$ are tuned from $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ separately. For the 4DA dataset, the results in Fig.3(a,b,c) show that the performance of our model using $\alpha$ with small values and $\beta$ with large values is often better than other settings. Moreover, the classification accuracy decreases greatly when $\alpha$ becomes larger and $\beta$ becomes smaller. For the IXMAS dataset, the results in Fig.3(d,e,f) show that a better classification accuracy can be obtained when the values of $\alpha$ and $\beta$ are the same, and the classification accuracy will decrease largely with the increasing of the difference between $\alpha$ and $\beta$.

## 5. CONCLUSION

This paper proposes an unsupervised domain adaptation method called Transfer Latent Representation (TLR). TLR aims to learn latent representations of the source and target domains. In latent space, the common properties of both domains are preserved and noise that causes domain shift is reduced. Experimental results on real-world cross-domain datasets demonstrate the effectiveness of our method.

In the future, we plan to extend TLR to solve the semi-supervised domain adaptation problem, in which there are only a few data labels in the target domain.

## 6. APPENDIX

### 6.1. Proof of Proposition 2

The objective function $F(W)$ can be rewritten as

$$
\begin{aligned}
F(W) = {}& tr(W^\top KLKW) \\
& + \alpha tr((WW^\top H_S^\top - H_S^\top)(H_S WW^\top - H_S)) \\
& + \beta tr((WW^\top H_T^\top - H_T^\top)(H_T WW^\top - H_T)) \\
= {}& tr(W^\top KLKW) \\
& + tr(WW^\top(\alpha H_S^\top H_S + \beta H_T^\top H_T)WW^\top) \\
& - 2tr(W^\top(\alpha H_S^\top H_S + \beta H_T^\top H_T)W) \\
& + tr(\alpha H_S^\top H_S + \beta H_T^\top H_T).
\end{aligned}
$$

We let $A = \alpha H_S^\top H_S + \beta H_T^\top H_T$, so that $A$ can be described as

$$
A = \begin{bmatrix} H_S^\top & H_T^\top \end{bmatrix} \begin{bmatrix} \alpha I_{n_1 \times n_1} & 0_{n_1 \times n_2} \\ 0_{n_2 \times n_1} & \beta I_{n_2 \times n_2} \end{bmatrix} \begin{bmatrix} H_S \\ H_T \end{bmatrix} = KMK.
$$
(24)

in which $M = \begin{bmatrix} \alpha I_{n_1 \times n_1} & 0_{n_1 \times n_2} \\ 0_{n_2 \times n_1} & \beta I_{n_2 \times n_2} \end{bmatrix}$. The objective function $F(W)$ then can be compactly rewritten as

$$
\begin{aligned}
F(W) = {}& tr(W^\top KLKW) + tr(WW^\top AWW^\top) \\
& - 2tr(W^\top AW) + tr(A).
\end{aligned}
$$

Futhermore, we let

$$
B = KLK. \tag{25}
$$

By substituting (25) into $F(W)$, the proposed model can be summarized as follows:

$$
\begin{aligned}
F(W) = {}& tr(W^\top BW) + tr(WW^\top AWW^\top) \\
& - 2tr(W^\top AW) + tr(A) \\
= {}& tr(WW^\top AWW^\top + W^\top BW - 2W^\top AW + A)
\end{aligned}
$$
(26)

Then, Proposition 2 is proven.

## 7. REFERENCES

[1] Shi Yuan and Sha Fei, "Information-theoretical learning of discriminative clusters for unsupervised domain adaptation," in *International Conference on Machine Learning (ICML)*, 2012, pp. 1275–1282.

[2] Sinno Jialin Pan and Qiang Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge & Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.

[3] Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars, "Unsupervised visual domain adaptation using subspace alignment," in *IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 2960–2967.

[4] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2066–2073.

[5] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang, "Domain adaptation via transfer component analysis," *IEEE Transactions on Neural Networks*, vol. 22, no. 2, pp. 199–210, 2011.

[6] K. Yan, L. Kou, and D. Zhang, "Learning domain-invariant subspace using domain features and independence maximization," *IEEE Transactions on Cybernetics*, vol. PP, no. 99, pp. 1–12, 2017.

[7] Arthur Gretton, Karsten M Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola, "A kernel method for the two-sample-problem," in *Advances in Neural Information Processing Systems*, 2007, pp. 513–520.

[8] Sinno Jialin Pan, James T Kwok, and Qiang Yang, "Transfer learning via dimensionality reduction," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2008, pp. 677–682.

[9] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu, "Transfer feature learning with joint distribution adaptation," in *IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 2200–2207.

[10] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell, "Adapting visual category models to new domains," in *European Conference on Computer Vision (ECCV)*, 2010, pp. 213–226.

[11] Jingen Liu, M Shah, B Kuipers, and S Savarese, "Cross-view action recognition via view knowledge transfer," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011, pp. 3209–3216.