

# Supplementary Material for Domain Adaptation by Joint Distribution Invariant Projections

Sentao Chen, Mehrtash Harandi, *Member, IEEE*, Xiaona Jin, and Xiaowei Yang

## I. PROOFS, OPTIMIZATION ALGORITHM, AND STATISTICAL TEST

In this supplementary material, we present the details of the proofs, the optimization algorithm, and the statistical test for the main paper.

1, In Part A, we provide the proof for Proposition 1 in Section III-C.

2, In Part B, we elaborate on the Riemannian Gradient Descent (RGD) technique in Section IV.

3, In Part C, we provide the proof for Theorem 1 in Section IV-A.

4, In Part D, we provide the proof for Lemma 1 in Section V.

5, In Part E, we provide the proof for Theorem 2 in Section V.

6, In Part F, we describe the procedure of the Wilcoxon signed-ranks test conducted in Section VI-C.

This work was supported by the National Natural Science Foundation of China (61906069), Guangdong Basic and Applied Basic Research Foundation (2019A1515011411, 2019A1515011700), Project Funded by China Postdoctoral Science Foundation (2019M662912), Science and Technology Program of Guangzhou (202002030355), and Fundamental Research Funds for the Central Universities (2019MS088). (*Corresponding author: Sentao Chen; Xiaowei Yang.*)

S. Chen, X. Jin and X. Yang are with the School of Software Engineering, South China University of Technology, Guangzhou 510006, China (e-mail: sentaochen@yahoo.com).

M. Harandi is with the Department of Electrical and Computer Systems Engineering, Monash University, Clayton VIC 3800, Australia.

A. *Proof of Proposition 1 in Section III-C of the Main Paper*  
*Proof.* By Eq. (3) in the main paper, we have

$$\begin{aligned} L^2(P^s(\mathbf{z}, y), P^t(\mathbf{z}, y)) &= \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y))^2 d\mathbf{z}dy \\ &= \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y)^2 - 2P^s(\mathbf{z}, y)P^t(\mathbf{z}, y) + P^t(\mathbf{z}, y)^2) d\mathbf{z}dy \end{aligned} \quad (1)$$

$$\begin{aligned} &= \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y)^2 - P^s(\mathbf{z}, y)P^t(\mathbf{z}, y)) d\mathbf{z}dy \\ &\quad - \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y)P^t(\mathbf{z}, y) - P^t(\mathbf{z}, y)^2) d\mathbf{z}dy \end{aligned} \quad (2)$$

$$\begin{aligned} &= \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y))P^s(\mathbf{z}, y) d\mathbf{z}dy \\ &\quad - \int_{(\mathbf{z}, y)} (P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y))P^t(\mathbf{z}, y) d\mathbf{z}dy. \end{aligned} \quad (3)$$

Note that the two terms in the last equation are expectations. According to the law of large numbers [1], given two sets of samples  $\mathcal{D}_{\varphi_s}^s$  and  $\mathcal{D}_{\varphi_t}^t$ , the empirical estimate of the  $L^2$ -distance can therefore be obtained as

$$\begin{aligned} \widehat{L}^2(\mathcal{D}_{\varphi_s}^s, \mathcal{D}_{\varphi_t}^t) &= \frac{1}{m_s} \sum_{i=1}^{m_s} (P^s(\mathbf{z}_i^s, y_i^s) - P^t(\mathbf{z}_i^s, y_i^s)) \\ &\quad - \frac{1}{m_t} \sum_{i=1}^{m_t} (P^s(\mathbf{z}_i^t, y_i^t) - P^t(\mathbf{z}_i^t, y_i^t)), \end{aligned} \quad (4)$$

which concludes the proof.  $\square$

B. *Optimization Framework for Section IV of the Main paper*

Recently, optimization on Riemannian manifolds have become increasingly prevalent in domain adaptation [2], [3], [4], [5]. In the main paper, we will face non-convex optimization problem defined in the product space of two Stiefel manifolds, which can be characterized as

$$\begin{aligned} &\min_{\mathbf{W}_s, \mathbf{W}_t} f(\mathbf{W}_s, \mathbf{W}_t) \\ &\text{s.t. } (\mathbf{W}_s, \mathbf{W}_t) \in \text{St}(d, D) \times \text{St}(d, D), \end{aligned} \quad (5)$$

where  $f(\mathbf{W}_s, \mathbf{W}_t)$  is an objective function and  $\text{St}(d, D) \times \text{St}(d, D)$  is the product set of two Stiefel manifolds. The Stiefel manifold  $\text{St}(d, D)$  is defined as the space of  $(D \times d)$ -dimensional matrices with orthonormal columns,  $d \leq D$  [6], [4]. That is,  $\text{St}(d, D) = \{\mathbf{W} \in \mathbb{R}^{D \times d} | \mathbf{W}^\top \mathbf{W} = \mathbf{I}\}$ , where  $\mathbf{I}$  is the identity matrix. Due to the concept of product topology [6], [4], the product set  $\text{St}(d, D) \times \text{St}(d, D)$  can be given

the structure of a Riemannian manifold. In the following, this product manifold is denoted as  $\mathcal{M}_{prod.} = \text{St}(d, D) \times \text{St}(d, D)$ . In the experiments of the main paper, we will make use of the first-order Riemannian Gradient Descent (RGD) method [6], [4], [7] to optimize over  $\mathcal{M}_{prod.}$ .

In general, the RGD method involves calculating the Euclidean gradient of the objective function, then projecting it onto the *tangent space*<sup>1</sup> to get the Riemannian gradient, and finally performing a *retraction* that maps the tangent vector back to the manifold. On the product manifold  $\mathcal{M}_{prod.}$ , the Riemannian gradient of  $f(\mathbf{W}_s, \mathbf{W}_t)$  is expressed as

$$\text{grad}f(\mathbf{W}_s, \mathbf{W}_t) = \left( \nabla_{\mathbf{W}_s}(f) - \mathbf{W}_s \text{sym}(\mathbf{W}_s^\top \nabla_{\mathbf{W}_s}(f)), \nabla_{\mathbf{W}_t}(f) - \mathbf{W}_t \text{sym}(\mathbf{W}_t^\top \nabla_{\mathbf{W}_t}(f)) \right), \quad (7)$$

where  $\nabla_{\mathbf{W}}(f)$  is the Euclidean gradient at  $\mathbf{W}$  and  $\text{sym}(\mathbf{W}^\top \nabla_{\mathbf{W}}(f)) = \frac{1}{2}(\mathbf{W}^\top \nabla_{\mathbf{W}}(f) + (\nabla_{\mathbf{W}}(f))^\top \mathbf{W})$ . Besides, the *retraction* is given by

$$r_{(\mathbf{W}_s, \mathbf{W}_t)}(\mathbf{Z}_s, \mathbf{Z}_t) = (\text{uf}(\mathbf{W}_s + \mathbf{Z}_s), \text{uf}(\mathbf{W}_t + \mathbf{Z}_t)), \quad (8)$$

where  $\text{uf}(\mathbf{W} + \mathbf{Z}) = (\mathbf{W} + \mathbf{Z})[(\mathbf{W} + \mathbf{Z})^\top (\mathbf{W} + \mathbf{Z})]^{-\frac{1}{2}}$ . Finally, RGD on the product manifold  $\mathcal{M}_{prod.}$  can be summarized by the following steps:

- (i) Compute  $\text{grad}f(\mathbf{W}_s, \mathbf{W}_t)$  at the current solution  $(\mathbf{W}_s, \mathbf{W}_t)$ .
- (ii) Update the current solution as

$$(\mathbf{W}_s, \mathbf{W}_t) \leftarrow r_{(\mathbf{W}_s, \mathbf{W}_t)}(-\alpha \text{grad}f(\mathbf{W}_s, \mathbf{W}_t)),$$

where  $\alpha > 0$  is the step size.

These steps are repeated until convergence to a stationary point, or until a maximum number of iterations is reached.

### C. Proof of Theorem 1 in Section IV-A of the Main paper

*Proof.* In Eq. (11) of the main paper, we denote  $\mathbf{A} = (\mathbf{H} + \lambda \mathbf{I})^{-1}$ , then  $\mathbf{A} = \mathbf{A}^\top$  and  $\hat{\boldsymbol{\theta}} = \mathbf{A}\hat{\mathbf{b}}$ . With these notations,

$$f(\mathbf{W}_s, \mathbf{W}_t) = \sum_{i=1}^{m_{st}} \hat{b}_i \hat{\theta}_i = \sum_{i,j=1}^{m_{st}} \hat{b}_i \hat{b}_j a_{ij}, \quad (9)$$

where  $a_{ij}$  is the  $(i, j)$ -th element of  $\mathbf{A}$ .

$$\begin{aligned} \nabla_{\mathbf{W}_s}(f) &= \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_s}(\hat{b}_i) + \sum_{i=1}^{m_{st}} \hat{b}_i \nabla_{\mathbf{W}_s}(\hat{\theta}_i) \\ &= \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_s}(\hat{b}_i) + \sum_{i,j=1}^{m_{st}} \hat{b}_i a_{ij} \nabla_{\mathbf{W}_s}(\hat{b}_j) \\ &\quad + \sum_{i,j=1}^{m_{st}} \hat{b}_i \hat{b}_j \nabla_{\mathbf{W}_s}(a_{ij}) \\ &= 2 \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_s}(\hat{b}_i) + \sum_{i,j=1}^{m_{st}} \hat{b}_i \hat{b}_j \nabla_{\mathbf{W}_s}(a_{ij}). \end{aligned} \quad (10)$$

<sup>1</sup> The tangent space at a point on a manifold is a vector space that consists of the tangent vectors of all possible curves passing through this point [3].

Similarly,

$$\nabla_{\mathbf{W}_t}(f) = 2 \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_t}(\hat{b}_i) + \sum_{i=1}^{m_{st}} \hat{b}_i \hat{b}_j \nabla_{\mathbf{W}_t}(a_{ij}). \quad (11)$$

Since  $\frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1}$  [8], it holds that  $\frac{\partial \mathbf{A}}{\partial w_{uv}} = -\mathbf{A} \frac{\partial \mathbf{H}}{\partial w_{uv}} \mathbf{A}$ , where  $w_{uv}$  is the  $(u, v)$ -th element of  $\mathbf{W}_s$  or  $\mathbf{W}_t$ . Hence, we have

$$\begin{aligned} \sum_{i=1}^{m_{st}} \hat{b}_i \hat{b}_j \frac{\partial a_{ij}}{\partial w_{uv}} &= \sum_{i=1}^{m_{st}} \hat{b}_i \hat{b}_j \left( \frac{\partial \mathbf{A}}{\partial w_{uv}} \right)_{ij} \\ &= -\hat{\mathbf{b}}^\top \mathbf{A}^\top \frac{\partial \mathbf{H}}{\partial w_{uv}} \mathbf{A} \hat{\mathbf{b}} \\ &= -\hat{\boldsymbol{\theta}}^\top \frac{\partial \mathbf{H}}{\partial w_{uv}} \hat{\boldsymbol{\theta}} \\ &= -\sum_{i=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \left( \frac{\partial \mathbf{H}}{\partial w_{uv}} \right)_{ij} \\ &= -\sum_{i=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \frac{\partial h_{ij}}{\partial w_{uv}}, \end{aligned} \quad (12)$$

where  $\left( \frac{\partial \mathbf{A}}{\partial w_{uv}} \right)_{ij}$  denotes the  $(i, j)$ -th element of the matrix  $\frac{\partial \mathbf{A}}{\partial w_{uv}}$ . This immediately leads to the following equations:

$$\sum_{i=1}^{m_{st}} \hat{b}_i \hat{b}_j \nabla_{\mathbf{W}_s}(a_{ij}) = -\sum_{i=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \nabla_{\mathbf{W}_s}(h_{ij}), \quad (13)$$

$$\sum_{i=1}^{m_{st}} \hat{b}_i \hat{b}_j \nabla_{\mathbf{W}_t}(a_{ij}) = -\sum_{i=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \nabla_{\mathbf{W}_t}(h_{ij}). \quad (14)$$

Therefore, we have

$$\nabla_{\mathbf{W}_s}(f) = 2 \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_s}(\hat{b}_i) - \sum_{i,j=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \nabla_{\mathbf{W}_s}(h_{ij}), \quad (15)$$

$$\nabla_{\mathbf{W}_t}(f) = 2 \sum_{i=1}^{m_{st}} \hat{\theta}_i \nabla_{\mathbf{W}_t}(\hat{b}_i) - \sum_{i,j=1}^{m_{st}} \hat{\theta}_i \hat{\theta}_j \nabla_{\mathbf{W}_t}(h_{ij}). \quad (16)$$

We now elaborate on deriving the gradient  $\nabla_{\mathbf{W}_s}(\hat{b}_i)$ . The remaining ones  $\nabla_{\mathbf{W}_t}(\hat{b}_i)$ ,  $\nabla_{\mathbf{W}_s}(h_{ij})$ ,  $\nabla_{\mathbf{W}_t}(h_{ij})$  can also be obtained in a similar way. In particular, we first derive the gradient for the Gaussian kernel function, which is the building block of  $\nabla_{\mathbf{W}_s}(\hat{b}_i)$ :

$$\begin{aligned} \nabla_{\mathbf{W}} \exp \left( -\frac{\|\mathbf{W}^\top \mathbf{x}_i - \mathbf{W}^\top \mathbf{x}_j\|^2}{2\sigma^2} \right) \\ = \frac{-1}{\sigma^2} \exp \left( -\frac{\|\mathbf{W}^\top \mathbf{x}_i - \mathbf{W}^\top \mathbf{x}_j\|^2}{2\sigma^2} \right) (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{W}, \end{aligned} \quad (17)$$

$$\begin{aligned} \nabla_{\mathbf{W}} \exp \left( -\frac{\|\mathbf{W}^\top \mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) \\ = \frac{-1}{\sigma^2} \exp \left( -\frac{\|\mathbf{W}^\top \mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) (\mathbf{x}_i \mathbf{x}_i^\top \mathbf{W} - \mathbf{x}_i \mathbf{x}_j^\top). \end{aligned} \quad (18)$$

Therefore, when  $1 \leq i \leq m_s$ ,

$$\begin{aligned}
& \nabla_{\mathbf{W}_s}(\hat{b}_i) \\
&= \nabla_{\mathbf{W}_s} \left( \frac{1}{m_s} \sum_{n=1}^{m_s} \exp \left( - \frac{\|\mathbf{W}_s^\top \mathbf{x}_n^s - \mathbf{W}_s^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \delta(y_n^s = y_i) \right. \\
&\quad \left. - \frac{1}{m_t} \sum_{n=1}^{m_t} \exp \left( - \frac{\|\mathbf{W}_t^\top \mathbf{x}_n^t - \mathbf{W}_s^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \delta(y_n^t = y_i) \right) \\
&= \frac{-1}{m_s \sigma^2} \sum_{n=1}^{m_s} \exp \left( - \frac{\|\mathbf{W}_s^\top \mathbf{x}_n^s - \mathbf{W}_s^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \\
&\quad \times \delta(y_n^s = y_i) (\mathbf{x}_n^s - \mathbf{x}_i) (\mathbf{x}_n^s - \mathbf{x}_i)^\top \mathbf{W}_s \\
&\quad + \frac{1}{m_t \sigma^2} \sum_{n=1}^{m_t} \exp \left( - \frac{\|\mathbf{W}_t^\top \mathbf{x}_n^t - \mathbf{W}_s^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \\
&\quad \times \delta(y_n^t = y_i) (\mathbf{x}_i \mathbf{x}_i^\top \mathbf{W}_s - \mathbf{x}_i (\mathbf{x}_n^t)^\top \mathbf{W}_t). \tag{19}
\end{aligned}$$

When  $m_s + 1 \leq i \leq m_{st}$ ,

$$\begin{aligned}
& \nabla_{\mathbf{W}_s}(\hat{b}_i) \\
&= \nabla_{\mathbf{W}_s} \left( \frac{1}{m_s} \sum_{n=1}^{m_s} \exp \left( - \frac{\|\mathbf{W}_s^\top \mathbf{x}_n^s - \mathbf{W}_t^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \delta(y_n^s = y_i) \right. \\
&\quad \left. - \frac{1}{m_t} \sum_{n=1}^{m_t} \exp \left( - \frac{\|\mathbf{W}_t^\top \mathbf{x}_n^t - \mathbf{W}_t^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \delta(y_n^t = y_i) \right) \\
&= \frac{1}{m_s \sigma^2} \sum_{n=1}^{m_s} \exp \left( - \frac{\|\mathbf{W}_s^\top \mathbf{x}_n^s - \mathbf{W}_t^\top \mathbf{x}_i\|^2}{2\sigma^2} \right) \\
&\quad \times \delta(y_n^s = y_i) (\mathbf{x}_n^s (\mathbf{x}_n^s)^\top \mathbf{W}_s - \mathbf{x}_n^s \mathbf{x}_i^\top \mathbf{W}_t). \tag{20}
\end{aligned}$$

#### D. Proof of Lemma 1 in Section V of the Main Paper

*Proof.* For any hypothesis  $h \in \mathcal{H}$ ,

$$\begin{aligned}
& |\mathbb{E}_{(\mathbf{z}, y) \sim P^t(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] - \mathbb{E}_{(\mathbf{z}, y) \sim P^s(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)]| \\
&= \left| \int_{(\mathbf{z}, y)} \ell(h(\mathbf{z}), y) (P^t(\mathbf{z}, y) - P^s(\mathbf{z}, y)) d\mathbf{z} dy \right| \tag{21}
\end{aligned}$$

$$\leq \int_{(\mathbf{z}, y)} |P^t(\mathbf{z}, y) - P^s(\mathbf{z}, y)| d\mathbf{z} dy \tag{22}$$

$$\leq \frac{1}{M} \int_{(\mathbf{z}, y)} (P^t(\mathbf{z}, y) - P^s(\mathbf{z}, y))^2 d\mathbf{z} dy. \tag{23}$$

The first inequality uses the assumption  $\ell(\cdot, \cdot) \leq 1$  and the property of integral. The second inequality exploits the assumption  $M = \inf_{(\mathbf{z}, y) \in \mathcal{Z} \times \mathcal{Y}} |P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y)| > 0$ , which means  $\frac{1}{M} |P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y)| \geq 1$  and thus  $|P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y)| \leq \frac{1}{M} (P^s(\mathbf{z}, y) - P^t(\mathbf{z}, y))^2, \forall (\mathbf{z}, y) \in \mathcal{Z} \times \mathcal{Y}$  follows. Therefore,

$$\begin{aligned}
& \mathbb{E}_{(\mathbf{z}, y) \sim P^t(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] \leq \mathbb{E}_{(\mathbf{z}, y) \sim P^s(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] \\
&\quad + \frac{1}{M} \int_{(\mathbf{z}, y)} (P^t(\mathbf{z}, y) - P^s(\mathbf{z}, y))^2 d\mathbf{z} dy, \tag{24}
\end{aligned}$$

which concludes the proof.  $\square$

#### E. Proof of Theorem 2 in Section V of the Main Paper

*Proof.* We first prove a fact, which will be used to obtain our final result.

**Fact 1.** Let  $\mathcal{A}$  and  $\mathcal{B}$  be two random events. For any  $\delta \in (0, 1)$ , if  $P(\mathcal{A}) \geq 1 - \delta/2$  and  $P(\mathcal{B}) \geq 1 - \delta/2$ , then the random event  $\mathcal{A} \cap \mathcal{B}$  holds with probability at least  $1 - \delta$ .

The proof is straightforward:  $P(\mathcal{A} \cap \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B}) - P(\mathcal{A} \cup \mathcal{B}) \geq 2 - \delta - P(\mathcal{A} \cup \mathcal{B}) \geq 1 - \delta$ .

We then bound the source domain generalization error and the  $L^2$ -distance respectively by their empirical estimates. For bounding the source error, we apply the Vapnik-Chervonenkis inequality [9]. That is, with probability exceeding  $1 - \delta/2$ , for any hypothesis  $h \in \mathcal{H}$ , the following inequality holds:

$$\begin{aligned}
& E_{(\mathbf{z}, y) \sim P^s(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] \\
&\leq \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{z}_i^s), y_i^s) + \sqrt{\frac{4}{m} \left( p \log \frac{2em}{p} + \log \frac{8}{\delta} \right)}. \tag{25}
\end{aligned}$$

Since the empirical estimate  $\widehat{L}^2(\mathcal{D}_{\varphi_s}^s, \mathcal{D}_{\varphi_t}^t)$  converges in probability to its true value  $L^2(P^s(\mathbf{z}, y), P^t(\mathbf{z}, y))$ , for a given  $\delta$ , there exists a positive integer  $N_\delta$ , such that when  $m > N_\delta$ , with probability at least  $1 - \delta/2$ , the following inequality holds:

$$L^2(P^s(\mathbf{z}, y), P^t(\mathbf{z}, y)) \leq \widehat{L}^2(\mathcal{D}_{\varphi_s}^s, \mathcal{D}_{\varphi_t}^t) + 1. \tag{26}$$

Now regard inequalities (25) and (26) respectively as random events  $\mathcal{A}$  and  $\mathcal{B}$ . Based on Lemma 1 in the main paper and Fact 1, we have that with probability at least  $1 - \delta$ , for  $m > N_\delta$  and any hypothesis  $h \in \mathcal{H}$ ,

$$\begin{aligned}
& E_{(\mathbf{z}, y) \sim P^t(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] \\
&\leq E_{(\mathbf{z}, y) \sim P^s(\mathbf{z}, y)}[\ell(h(\mathbf{z}), y)] + \frac{1}{M} L^2(P^s(\mathbf{z}, y), P^t(\mathbf{z}, y)) \\
&\leq \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{z}_i^s), y_i^s) + \sqrt{\frac{4}{m} \left( p \log \frac{2em}{p} + \log \frac{8}{\delta} \right)} \\
&\quad + \frac{1}{M} L^2(P^s(\mathbf{z}, y), P^t(\mathbf{z}, y)) \tag{27}
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{z}_i^s), y_i^s) + \sqrt{\frac{4}{m} \left( p \log \frac{2em}{p} + \log \frac{8}{\delta} \right)} \\
&\quad + \frac{1}{M} (\widehat{L}^2(\mathcal{D}_{\varphi_s}^s, \mathcal{D}_{\varphi_t}^t) + 1), \tag{28}
\end{aligned}$$

which concludes the proof.  $\square$

#### F. Statistical Test in Section VI-C of the Main Paper

We describe the procedure of the Wilcoxon signed-ranks test [10], [11] on the tasks from Table I in the main paper. The test compares the performance of two algorithms over multiple tasks. Specifically, in each task the mean classification accuracy is adopted as the performance measure of the algorithms. We fix JDIP as a control algorithm, and conduct 13 pairs of tests: SVM versus JDIP, ..., and JDIP-rbf versus JDIP. To run the test, we rank the differences in performance of two algorithms for each task out of  $N$  tasks. The differences are ranked according to their absolute values. The smallest absolute value gets the rank of 1, the second smallest gets the

rank of 2, and so on. In case of equality, average ranks are assigned. The statistic of the Wilcoxon signed-ranks test is:

$$z(a, b) = \frac{T(a, b) - N(N + 1)/4}{\sqrt{N(N + 1)(2N + 1)/24}}, \quad (29)$$

where  $T(a, b) = \min\{R^+(a, b), R^-(a, b)\}$ .  $R^+(a, b)$  is the sum of ranks for the tasks on which algorithm  $b$  outperforms algorithm  $a$  and  $R^-(a, b)$  is the sum of ranks for the opposite. They are defined as follows:

$$R^+(a, b) = \sum_{\text{diff}_i > 0} \text{rank}(\text{diff}_i) + \frac{1}{2} \sum_{\text{diff}_i = 0} \text{rank}(\text{diff}_i), \quad (30)$$

$$R^-(a, b) = \sum_{\text{diff}_i < 0} \text{rank}(\text{diff}_i) + \frac{1}{2} \sum_{\text{diff}_i = 0} \text{rank}(\text{diff}_i), \quad (31)$$

where  $\text{diff}_i$  is the difference between the accuracy of two algorithms on the  $i$ -th task out of  $N$  tasks, and  $\text{rank}(\text{diff}_i)$  is the rank of  $|\text{diff}_i|$ . We fix  $b$  as JDIP, and let  $a$  vary from SVM to JDIP-rbf in turn. Based on formulas (29)-(31), we can compute  $z(a, b)$  for the 13 pairs of tests.

## REFERENCES

- [1] L. Wasserman, *All of statistics: a concise course in statistical inference*. Springer Science & Business Media, 2013.
- [2] M. Gong, K. Zhang, T. Liu, D. Tao, C. Glymour, and B. Schölkopf, "Domain adaptation with conditional transferable components," in *International conference on machine learning*, 2016, pp. 2839–2848.
- [3] M. Baktashmotlagh, M. Harandi, and M. Salzmann, "Distribution-matching embedding for visual domain adaptation," *Journal of Machine Learning Research*, vol. 17, no. 108, pp. 1–30, 2016.
- [4] S. Herath, M. Harandi, and F. Porikli, "Learning an invariant hilbert space for domain adaptation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 3956–3965.
- [5] Y. Liu, W. Tu, B. Du, L. Zhang, and D. Tao, "Homologous component analysis for domain adaptation," *IEEE Transactions on Image Processing*, vol. 29, pp. 1074–1089, 2020.
- [6] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [7] T. Zhang and Y. Yang, "Robust pca by manifold optimization," *Journal of Machine Learning Research*, vol. 19, no. 80, pp. 1–39, 2018.
- [8] K. B. Petersen and M. S. Pedersen, *The Matrix Cookbook*, 2012.
- [9] V. N. Vapnik, *Statistical Learning Theory*. John Wiley & Sons, 1998.
- [10] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *Journal of Machine Learning Research*, vol. 7, pp. 1–30, 2006.
- [11] S. Chen, L. Han, X. Liu, Z. He, and X. Yang, "Subspace distribution adaptation frameworks for domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.