Aligning Infinite-Dimensional Covariance Matrices in Reproducing Kernel Hilbert Spaces for Domain Adaption: Supplementary Material

Zhen Zhang, Mianzhi Wang, Yan Huang, Arye Nehorai Washington University in St. Louis

{zhen.zhang, mianzhi.wang, yanhuang640, nehorai}@wustl.edu

Abstract

The supplementary material consists of three parts. In the first part, we provide insights on our framework by discussing the expressiveness of covariance descriptors in RKHS. In the second part, we provide more discussion on the experiments. In the third part, we prove all the mathematical results presented in the paper.

1. Discussion on RKHS covariance descriptors

Given RKHS data matrix $\Phi_X = [\phi(x_1), \phi(x_2), ..., \phi(x_N)]$, the maximum likelihood estimation of the RKHS covariance descriptor is

$$MC = \frac{1}{N} \sum_{i=1}^{N} \left[\phi(x_i) - \hat{\mu} \right] \left[\phi(x_i) - \hat{\mu} \right]^T = \Phi_{\boldsymbol{X}} \boldsymbol{J}_N \boldsymbol{J}_N^T \Phi_{\boldsymbol{X}}^T,$$
(1)

where $\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} \phi(x_i)$ is the empirical mean of samples. If we use the linear kernel, *i.e.*, $k(\vec{x}, \vec{y}) = \vec{x}^T \vec{y}$, then the corresponding feature map is just the identity function, *i.e.*, $\phi(\vec{x}) = \vec{x}$. As a result, the expression (1) degenerates to the MLE of covariance matrices in \mathbb{R}^n . For the computationally efficient estimation (See Section 3.2 in the paper), we have the same explanation. Therefore, RKHS covariance descriptors are the generalization of covariance matrices.

From the standpoint of kernel methods, with the nonlinear kernel, *e.g.*, the RBF kernel, RKHS covariance descriptors, which can capture nonlinear structure and higher-order correlations, are more informative than the covariance matrices. When aligning two infinite-dimensional covariance descriptors, we in fact match infinitely many orders of statistics.

Similar strategy of employing RKHS covariance descriptors to characterize a set of samples can be found in [1, 3], where the authors represent each image by a covariance descriptor in RKHS and quantify the discrepancies between covariance descriptors to classify images.

2. Discussion on the experiments

2.1. More about the datasets

Fig. 1(a) shows sample images in the COIL20 dataset. For each object in COIL20, we provide one example image of the 72 total. Fig. 1(b) shows sample images from the monitor category in the Office-Caltech dataset.

Table 1 lists the top categories and subcategories in the 20-Newsgroups dataset.

2.2. Visualization using kernel PCA

We use kernel principal components analysis [4] (kernel PCA) to visualize source and target samples in RKHS. We implement experiments on a cross-domain dataset generated from 20-Newsgroups. The source dataset consists of four subcategories of Comp and Rec: comp.graphics, comp.sys.mac.hardware, rec.sport.baseball, and rec.sport.hokey. The target dataset consists of the other four subcategories: comp.os.ms-windows.misc, comp.sys.ibm.pc.hardware, rec.autos, and rec.motorcycles. So there are 970+958+991+997 = 3916 samples in the source domain, and 963+979+987+993 = 3922



Figure 1: (a) Sample images from the COIL20 dataset. (b) Sample images from the Office-Caltech dataset.

Top Category	Subcategory	Number of Samples
Comp	comp.graphics	970
	comp.os.ms-windows.misc	963
	comp.sys.ibm.pc.hardware	979
	comp.sys.mac.hardware	958
Rec	rec.autos	987
	rec.motorcycles	993
	rec.sport.baseball	991
	rec.sport.hokey	997
Sci	sci.crypt	989
	sci.electronics	984
	sci.med	987
	sci.space	985
Talk	talk.politics.mideast	940
	talk.politics.misc	774
	talk.politics.guns	909
	talk.region.misc	627

Table 1: Top categories and subcategories in the 20-Newsgroups dataset. All features have a dimensional of 25,804.

samples in the target domain. We employ the RBF kernel, and visualize the coefficients of RKHS samples with respect to the first three principal components. Fig. 2(a) shows the results with the non-adapted kernel matrix K. We can see that the source and target distributions are very different. Fig. 2(b) and Fig. 2(c) show the results with the domain-invariant kernel matrices WC \tilde{K} and OT \tilde{K} (see Section 5 in the paper), respectively. From Fig. 2(b) and (c), we conclude that after "moving" the source data by the kernel whitening-coloring map or the kernel optimal transport map, the transformed source samples and target samples are closely distributed in RKHS. In addition, we note that the recognition accuracies with non-adapted kernel matrix K and domain-invariant kernel matrices WC \tilde{K} and OT \tilde{K} are 61.19%, 94.11% and 95.99%, respectively. The highly superior performances of our approaches demonstrate the effectiveness of aligning covariance descriptors.

2.3. Out-of-sample generalization on the Reuters-21578 dataset

In this subsection, we measure our approaches' ability to generalize out-of-sample patterns. We follow the experimental protocol in [2], and conduct experiments on the preprocessed Reuters-21578 datasets. To train the model, we randomly select 500 labeled samples from the source domain and 300 unlabeled samples from the target domain. We test the model on the remaining unlabeled samples in the target domain. We repeat the above procedures 500 times, and report the average



Figure 2: We implement experiments on a cross-domain dataset generated from 20-Newsgroups. We use kernel PCA [4] to visualize the data in RKHS. (a) Visualization of the source and target samples with the non-adapted kernel. (b) Visualization of the transformed source and target samples with the domain-invariant kernel WC \tilde{K} . (c) Visualization of the transformed source and target samples with the domain-invariant kernel OT \tilde{K} .



Figure 3: Accuracies and confidence intervals in recognizing out-of-sample data of the Reuters-21578 dataset. For all four methods, we use the linear kernel.

accuracies and standard errors. In these experiments, we compare our approaches with only the standard SVM and TCA, both of which possess generalizability. The parameters are set to be the same as those in the transductive setting. In Fig. 3 and Fig. 4, the experimental results with the both linear kernel and RBF kernel show that our approaches KWC and KOT outperform TCA and SVM with statistical significance.



Figure 4: Accuracies and confidence intervals in recognizing out-of-sample data of the Reuters-21578 dataset. For all four methods, we use the RBF kernel.

3. Proofs of the mathematical results in the paper

We first provide some useful lemmas (corollaries), which will be frequently used.

Lemma 1. Let \mathcal{H}_1 and \mathcal{H}_2 be two separable Hilbert spaces. Let $G : \mathcal{H}_1 \to \mathcal{H}_2$ be a linear operator with finite rank, and let $G^* : \mathcal{H}_2 \to \mathcal{H}_1$ be its adjoint operator. Then $\operatorname{Im}(G) = \operatorname{Im}(GG^*)$.

Proof. We need to show that $\text{Im}(G) \subseteq \text{Im}(GG^*)$, and $\text{Im}(GG^*) \subseteq \text{Im}(G)$.

- 1. Since $\operatorname{Im}(G) = \operatorname{Ker}^{\perp}(G^*)$ and $\operatorname{Im}(GG^*) = \operatorname{Ker}^{\perp}[(GG^*)^*] = \operatorname{Ker}^{\perp}(GG^*)^1$, to obtain $\operatorname{Im}(G) \subseteq \operatorname{Im}(GG^*)$, it is sufficient to show that $\operatorname{Ker}(GG^*) \subseteq \operatorname{Ker}(G^*)$. $\forall v \in \operatorname{Ker}(GG^*), GG^*v = 0 \implies \langle GG^*v, v \rangle_{\mathcal{H}_2} = \langle G^*v, G^*v \rangle_{\mathcal{H}_1} = 0 \implies G^*v = 0 \implies v \in \operatorname{Ker}(G^*)$. So $\operatorname{Ker}(GG^*) \subseteq \operatorname{Ker}(G^*)$.
- 2. $\forall v \in \text{Im}(GG^*)$, there exists $w \in \mathcal{H}_2$, such that $v = GG^*w$. Writing $v = G(G^*w)$, we have $v \in \text{Im}(G)$. So $\text{Im}(GG^*) \subseteq \text{Im}(G)$.

Lemma 2. Let \mathcal{H}_1 and \mathcal{H}_2 be two separable Hilbert spaces. Let $G : \mathcal{H}_1 \to \mathcal{H}_2$ be a linear operator with finite rank, and let $G^* : \mathcal{H}_2 \to \mathcal{H}_1$ be its adjoint operator. Let $\{\lambda_k(G^*G)\}_{k=1}^r$ and $\{\varphi_k(G^*G)\}_{k=1}^r$ be the positive eigenvalues and the corresponding orthonormal eigenvectors of the operator G^*G . Then,

$$GG^* = \sum_{k=1}^r \lambda_k(G^*G) \frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \otimes \frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}$$
(2)

¹Given two Hilbert spaces \mathcal{H}_1 and \mathcal{H}_2 , and a linear operator $G: \mathcal{H}_1 \to \mathcal{H}_2$, the kernel space of G is defined as $\operatorname{Ker}(G) = \{v \in \mathcal{H}_1, Gv = 0_{\mathcal{H}_2}\}$.

is the orthogonal eigen-decomposition of the operator GG^* .

Remark 1. The tensor product \otimes is defined such that $(u \otimes v)w = u \langle v, w \rangle_{\mathcal{H}}, \forall u, v, w \in \mathcal{H}$, which is the analogy of the outer product in \mathbb{R}^n , i.e., $\vec{u} \otimes \vec{v} = \vec{u}\vec{v}^T, \forall \vec{u}, \vec{v} \in \mathbb{R}^n$. In our paper, these two expressions $u \otimes v$ and $uv^T, \forall u, v \in \mathcal{H}$, have the same meaning.

Proof. $\{\lambda_k(G^*G)\}_{k=1}^r$ and $\{\varphi_k(G^*G)\}_{k=1}^r$ are the respective eigenvalues and eigenvectors of G^*G . Then $G^*G\varphi_k(G^*G) = \lambda_k(G^*G)\varphi_k(G^*G)$. After applying the operator G on both sides, we obtain $GG^*G\varphi_k(G^*G) = \lambda_k(G^*G)G\varphi_k(G^*G)$. So $\{\lambda_k(G^*G)\}_{k=1}^r$ are the positive eigenvalues of GG^* .

To show that $\{G\varphi_k(G^*G)\}_{k=1}^r$ are eigenvectors of GG^* , we need to show that $G\varphi_k(G^*G)$ is nonzero, $\forall k = 1, 2, ...r$. Suppose $G\varphi_k(G^*G) = 0$, then $G^*G\varphi_k(G^*G) = \lambda_k(G^*G)\varphi_k(G^*G) = 0$, which implies that $\varphi_k(G^*G) = 0$, contradicting the fact that $\varphi_k(G^*G)$ is an eigenvector. So $\{G\varphi_k(G^*G)\}_{k=1}^r$ are the eigenvectors of GG^* .

We also need to show $\{\lambda_k(G^*G)\}_{k=1}^r$ are the whole positive eigenvalues of GG^* . It is equivalent to showing that if λ is an positive eigenvalue of GG^* , then λ is also an eigenvalue of G^*G . To achieve this, we can just repeat the above procedure. Finally, we need to show that $\{\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}\}_{k=1}^r$ are orthonormal. For any k, l = 1, 2, ..., r,

$$\begin{split} \langle \frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}, \frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \rangle_{\mathcal{H}_2} &= \frac{1}{\lambda_k(G^*G)} \langle \varphi_k(G^*G), G^*G\varphi_k(G^*G) \rangle_{\mathcal{H}_1} \\ &= \frac{1}{\lambda_k(G^*G)} \langle \varphi_k(G^*G), \lambda_k(G^*G)\varphi_k(G^*G) \rangle_{\mathcal{H}_1} = 1, \\ \langle \frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}, \frac{G\varphi_l(G^*G)}{\sqrt{\lambda_l(G^*G)}} \rangle_{\mathcal{H}_2} &= \langle \frac{\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}, \frac{G^*G\varphi_l(G^*G)}{\sqrt{\lambda_l(G^*G)}} \rangle_{\mathcal{H}_1} = \lambda_l(G^*G) \langle \frac{\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}, \frac{\varphi_l(G^*G)}{\sqrt{\lambda_l(G^*G)}} \rangle_{\mathcal{H}_1} = 0. \end{split}$$

Corollary 1. The projection operator P_{GG^*} on the subspace $\text{Im}(GG^*)$ is

$$P_{GG^*} = \sum_{k=1}^{r} \left[\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \right] \otimes \left[\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \right]$$
$$= G\left[\sum_{k=1}^{r} \lambda_k(G^*G)^{-1} \varphi_k(G^*G) \otimes \varphi_k(G^*G) \right] G^*$$
$$= G(G^*G)^{\dagger} G^*.$$
(3)

Proof. From Lemma 2, we have that $\{\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}}\}_{k=1}^r$ are the orthonormal basis of $\text{Im}(GG^*)$. We can obtain P_{GG^*} immediately.

Corollary 2. The square root of the operator GG^* is

$$(GG^*)^{\frac{1}{2}} = \sum_{k=1}^r \left[\lambda_k(G^*G) \right]^{\frac{1}{2}} \left[\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \right] \otimes \left[\frac{G\varphi_k(G^*G)}{\sqrt{\lambda_k(G^*G)}} \right]$$
$$= G\left[\sum_{k=1}^r \lambda_k(G^*G)^{-\frac{1}{2}} \varphi_k(G^*G) \otimes \varphi_k(G^*G) \right] G^*$$
$$= G(G^*G)^{\frac{1}{2}} G^*.$$
(4)

Proof. Immediately by Lemma 2.

Corollary 3. The Moore-Penrose inverse of the square root of the operator GG^* is

$$(GG^{*})^{\dagger \frac{1}{2}} = \sum_{k=1}^{r} \left[\lambda_{k}(G^{*}G) \right]^{-\frac{1}{2}} \left[\frac{G\varphi_{k}(G^{*}G)}{\sqrt{\lambda_{k}(G^{*}G)}} \right] \otimes \left[\frac{G\varphi_{k}(G^{*}G)}{\sqrt{\lambda_{k}(G^{*}G)}} \right]$$
$$= G \left[\sum_{k=1}^{r} \lambda_{k}(G^{*}G)^{-\frac{3}{2}} \varphi_{k}(G^{*}G) \otimes \varphi_{k}(G^{*}G) \right] G^{*}$$
$$= G (G^{*}G)^{\frac{1}{2}} G^{*}.$$
 (5)

Table 2: The list of notations in the paper and the supplementary material

RKHS	$\mathcal{H}_{\mathcal{K}}$	centered kernel matrix C_{XX}	$oldsymbol{C}_{XX} = oldsymbol{J}_{N_s}^T oldsymbol{K}_{XX} oldsymbol{J}_{N_s}$
kernel function	k	centered kernel matrix C_{YY}	$oldsymbol{C}_{YY} = oldsymbol{J}_{N_t}^T oldsymbol{K}_{YY} oldsymbol{J}_{N_t}$
implicit feature map	ϕ	centered kernel matrix C_{YX}	$oldsymbol{C}_{YX} = oldsymbol{J}_{N_t}^Toldsymbol{K}_{YX}oldsymbol{J}_{N_s}$
identity operator in $\mathcal{H}_{\mathcal{K}}$	$\mathcal{I}_{\mathcal{H}_{\mathcal{K}}}$	top d eigenvalues and eigenvectors of C_{XX}	$(oldsymbol{\Lambda}_X,oldsymbol{V}_X)$
source samples number	N_s	top d eigenvalues and eigenvectors of C_{YY}	$(oldsymbol{\Lambda}_Y,oldsymbol{V}_Y)$
target samples number	N_t	matrix W_X	$oldsymbol{W}_X = oldsymbol{J}_{N_s}oldsymbol{V}_X(oldsymbol{I}_d - hooldsymbol{\Lambda}_X^{-1})^{rac{1}{2}}$
source samples in $\mathcal{H}_{\mathcal{K}}$	Φ_X	matrix W_Y	$\boldsymbol{W}_{Y} = \boldsymbol{J}_{N_{t}} \boldsymbol{V}_{Y} (\boldsymbol{I}_{d} - \rho \boldsymbol{\Lambda}_{Y}^{-1})^{\frac{1}{2}}$
target samples in $\mathcal{H}_{\mathcal{K}}$	Φ_Y	matrix C_{XY}^w	$oldsymbol{C}_{XY}^w = oldsymbol{W}_X^Toldsymbol{K}_{XY}oldsymbol{W}_Y$
$N_s \times N_s$ centering matrix	$oldsymbol{J}_{N_s} = rac{1}{\sqrt{N_s}} (oldsymbol{I} - rac{1}{N_s} ec{1} ec{1}^T)$	matrix $oldsymbol{C}_{YX}^w$	$oldsymbol{C}^w_{YX} = (oldsymbol{C}^w_{XY})^T$
$N_t \times N_t$ centering matrix	$\boldsymbol{J}_{N_t} = rac{1}{\sqrt{N_t}} (\boldsymbol{I} - rac{1}{N_t} \vec{1} \vec{1}^T)$	matrix B	$\boldsymbol{B} = \boldsymbol{C}_{YX} (\boldsymbol{C}_{XX} + \rho \boldsymbol{I}_{N_s})^{-\frac{1}{2}}$
kernel matrix K_{XX}	$oldsymbol{K}_{XX}=\Phi^T_X\Phi_X$		$\mathbf{D} = \begin{bmatrix} \mathbf{C}^w & \mathbf{C}^w & \bot c(\mathbf{A} & z\mathbf{I}) \end{bmatrix}^{\dagger \frac{1}{2}}$
kernel matrix K_{YY}	$oldsymbol{K}_{YY}=\Phi_Y^T\Phi_Y$	matrix D	$D = \begin{bmatrix} \mathbf{U}_{YX} \mathbf{U}_{XY} + \rho(\mathbf{\Lambda}_Y - \rho \mathbf{I}_d) \end{bmatrix}^{-2}$
kernel matrix K_{YX}	$oldsymbol{K}_{YX}=\Phi_Y^T\Phi_X$		$W_Y^{T}K_{YX}J_{N_s}$

Proof. Immediately by Lemma 2.

Lemma 3. Let \mathcal{H} be a separable Hilbert space. Let $\psi_1, \psi_2, ..., \psi_n$ be n elements in \mathcal{H} . We define the operator $\Psi : \mathbb{R}^n \to \mathcal{H}$ as $\Psi(\mathbf{x}) = \sum_{i=1}^n x_i \psi_i, \forall \mathbf{x} \in \mathbb{R}^n$. Then $\Psi^T : \mathcal{H} \to \mathbb{R}^n$ defined as $\Psi^T(u) = [\langle \psi_1, u \rangle_{\mathcal{H}}, \langle \psi_2, u \rangle_{\mathcal{H}}, ..., \langle \psi_n, u \rangle_{\mathcal{H}}]^T, \forall u \in \mathcal{H}$ is the adjoint operator of Ψ , i.e., $\Psi^* = \Psi^T$.

Proof.

$$\forall \boldsymbol{x} \in \mathbb{R}^{n}, \forall u \in \mathcal{H}, \quad \langle \Psi(\boldsymbol{x}), u \rangle_{\mathcal{H}} = \sum_{i=1}^{n} x_{i} \langle \psi_{i}, u \rangle_{\mathcal{H}} = \langle \boldsymbol{x}, \Psi^{T}(u) \rangle.$$
(6)

Note that if \mathcal{H}_1 and \mathcal{H}_2 are Euclidean spaces, say \mathbb{R}^{n_1} and \mathbb{R}^{n_2} , then the operator $G : \mathcal{H}_1 \to \mathcal{H}_2$ is just an $n_2 \times n_1$ matrix. All the above conclusions still hold. In the next section, we provide the proofs of all mathematical results in the paper. For convenience, we list the relevant notations in Table 2.

3.1. Proving Theorem 1

In \mathbb{R}^n , the whitening-coloring map is $\hat{T}_{WC} = \Sigma_t^{\frac{1}{2}} (\Sigma_s^{\frac{1}{2}})^{\dagger}$.

Theorem 1. If
$$\operatorname{Im}(\Sigma_t) \subseteq \operatorname{Im}(\Sigma_s)$$
, then $T_{\mathrm{WC}}\Sigma_s T_{\mathrm{WC}}^T = \Sigma_t$.

Proof. Substituting $\hat{T}_{\mathrm{WC}} = \Sigma_t^{\frac{1}{2}} (\Sigma_s^{\frac{1}{2}})^{\dagger}$ into the left part, we have

$$\hat{\boldsymbol{T}}_{\mathrm{WC}}\boldsymbol{\Sigma}_{s}\hat{\boldsymbol{T}}_{\mathrm{WC}}^{T} = \left[\boldsymbol{\Sigma}_{t}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\right]\boldsymbol{\Sigma}_{s}\left[\boldsymbol{\Sigma}_{t}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\right]^{T} = \boldsymbol{\Sigma}_{t}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\boldsymbol{\Sigma}_{s}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}} = \boldsymbol{\Sigma}_{t}^{\frac{1}{2}}\boldsymbol{P}_{s}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}} = \boldsymbol{\Sigma}_{t},\tag{7}$$

where P_s is the projection matrix onto the image space $\text{Im}(\Sigma_s)$, and the last equality holds because $\text{Im}(\Sigma_t) \subseteq \text{Im}(\Sigma_s)$. \Box

3.2. Proving the equivalence between two expressions of the optimal transport map

Given two positive definite covariance matrix Σ_t and Σ_s , we have

$$\Sigma_{t}^{\frac{1}{2}} (\Sigma_{t}^{\frac{1}{2}} \Sigma_{s} \Sigma_{t}^{\frac{1}{2}})^{-\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}} = \Sigma_{s}^{-\frac{1}{2}} (\Sigma_{s}^{\frac{1}{2}} \Sigma_{t} \Sigma_{s}^{\frac{1}{2}})^{\frac{1}{2}} \Sigma_{s}^{-\frac{1}{2}}.$$
(8)

Proof. We write $\Sigma_t = \Sigma_t^{\frac{1}{2}} (\Sigma_t^{\frac{1}{2}})^T$, and then substitute it into the right part of (8). We have

$$\boldsymbol{\Sigma}_{s}^{-\frac{1}{2}} (\boldsymbol{\Sigma}_{s}^{\frac{1}{2}} \boldsymbol{\Sigma}_{t} \boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\frac{1}{2}} \boldsymbol{\Sigma}_{s}^{-\frac{1}{2}} = \boldsymbol{\Sigma}_{s}^{-\frac{1}{2}} [(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}} \boldsymbol{\Sigma}_{t}^{\frac{1}{2}}) (\boldsymbol{\Sigma}_{s}^{\frac{1}{2}} \boldsymbol{\Sigma}_{t}^{\frac{1}{2}})^{T}]^{\frac{1}{2}} \boldsymbol{\Sigma}_{s}^{-\frac{1}{2}}$$
(9a)

$$= \Sigma_{s}^{-\frac{1}{2}} (\Sigma_{s}^{\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}}) [(\Sigma_{s}^{\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}})^{T} (\Sigma_{s}^{\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}})]^{-\frac{1}{2}} (\Sigma_{s}^{\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}})^{T} \Sigma_{s}^{-\frac{1}{2}}$$
(9b)

$$= \Sigma_{s}^{-\frac{1}{2}} \Sigma_{s}^{\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}} (\Sigma_{t}^{\frac{1}{2}} \Sigma_{s} \Sigma_{t}^{\frac{1}{2}})^{-\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}} \Sigma_{s}^{\frac{1}{2}} \Sigma_{s}^{-\frac{1}{2}} \Sigma_{s}^{-\frac{1}{2}}$$
(9c)

$$= \Sigma_{t}^{\frac{1}{2}} (\Sigma_{t}^{\frac{1}{2}} \Sigma_{s} \Sigma_{t}^{\frac{1}{2}})^{-\frac{1}{2}} \Sigma_{t}^{\frac{1}{2}}, \tag{9d}$$

where (9b) holds because of Corollary 2.

3.3. Proving Theorem 2

In \mathbb{R}^n , the optimal transport map is $\hat{T}_{\text{OT}} = \Sigma_t^{\frac{1}{2}} (\Sigma_t^{\frac{1}{2}} \Sigma_s \Sigma_t^{\frac{1}{2}})^{\dagger \frac{1}{2}} \Sigma_t^{\frac{1}{2}}$.

Theorem 2. If $\operatorname{Ker}(\Sigma_s) \cap \operatorname{Im}(\Sigma_t) = \{\vec{0}\}$, then we have $\hat{T}_{OT} \Sigma_s \hat{T}_{OT}^T = \Sigma_t$.

Proof. Let $\Sigma_t = V_{n \times r} \Lambda_{r \times r} V_{n \times r}^T$, where Λ is a diagonal matrix whose diagonal terms are the r positive eigenvalues and V consists of the corresponding eigenvectors.

Claim: rank $(\Sigma_s^{\frac{1}{2}}V) = r.$ **(I)**

Let $\vec{v}_1, \vec{v}_2, ..., \vec{v}_r$ be the columns of V. It is sufficient to show that $\sum_s^{\frac{1}{2}} \vec{v}_1, \sum_s^{\frac{1}{2}} \vec{v}_2, ..., \sum_s^{\frac{1}{2}} \vec{v}_r$ are linearly independent.

Suppose there exist $\lambda_1, \lambda_2, ..., \lambda_r$, such that $\lambda_1 \Sigma_s^{\frac{1}{2}} \vec{v}_1 + \lambda_2 \Sigma_s^{\frac{1}{2}} \vec{v}_2 + ... + \lambda_r \Sigma_s^{\frac{1}{2}} \vec{v}_r = \vec{0}$. Then we have that $\lambda_1 \vec{v}_1 + \lambda_2 \vec{v}_2 + ... + \lambda_r \Sigma_s^{\frac{1}{2}} \vec{v}_r = \vec{0}$. $\dots + \lambda_r \vec{v}_r \in \operatorname{Ker}(\boldsymbol{\Sigma}_s^{\frac{1}{2}}) = \operatorname{Ker}(\boldsymbol{\Sigma}_s). \text{ Since } \operatorname{Im}(\boldsymbol{\Sigma}_t) = \operatorname{Span}\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_r\}, \text{ immediately, } \lambda_1 \vec{v}_1 + \lambda_2 \vec{v}_2 + ... + \lambda_r \vec{v}_r \in \operatorname{Im}(\boldsymbol{\Sigma}_t).$ By the condition that $\operatorname{Ker}(\boldsymbol{\Sigma}_s) \cap \operatorname{Im}(\boldsymbol{\Sigma}_t) = \{\vec{0}\},$ we have $\lambda_1 \vec{v}_1 + \lambda_2 \vec{v}_2 + ... + \lambda_r \vec{v}_r = \vec{0}.$ And $\vec{v}_1, \vec{v}_2, ..., \vec{v}_r$ are linearly independent $\implies \lambda_1 = \lambda_2 = ... = \lambda_r = 0.$ So $\Sigma_s^{\frac{1}{2}} \vec{v}_1, \Sigma_s^{\frac{1}{2}} \vec{v}_2, ..., \Sigma_s^{\frac{1}{2}} \vec{v}_r$ are linearly independent. (II) Now we start to prove that $\hat{T}_{\text{OT}} \Sigma_s \hat{T}_{\text{OT}}^T = \Sigma_t$. Substituting $\hat{T}_{\text{OT}} = \Sigma_t^{\frac{1}{2}} (\Sigma_t^{\frac{1}{2}} \Sigma_s \Sigma_t^{\frac{1}{2}})^{\frac{1}{2}} \Sigma_t^{\frac{1}{2}}$ into the left part, we obtain

$$\hat{T}_{\rm OT} \Sigma_s \hat{T}_{\rm OT}^T = \Sigma_t^{\frac{1}{2}} (\Sigma_t^{\frac{1}{2}} \Sigma_s \Sigma_t^{\frac{1}{2}})^{\dagger \frac{1}{2}} \Sigma_t^{\frac{1}{2}} \Sigma_s \Sigma_t^{\frac{1}{2}} (\Sigma_t^{\frac{1}{2}} \Sigma_s \Sigma_t^{\frac{1}{2}})^{\dagger \frac{1}{2}} \Sigma_t^{\frac{1}{2}} = \Sigma_t^{\frac{1}{2}} P_{ts} \Sigma_t^{\frac{1}{2}},$$
(10)

where P_{ts} is the projection matrix onto $\operatorname{Im}(\Sigma_t^{\frac{1}{2}}\Sigma_s\Sigma_t^{\frac{1}{2}}) = \operatorname{Im}[(\Sigma_t^{\frac{1}{2}}\Sigma_s^{\frac{1}{2}})(\Sigma_t^{\frac{1}{2}}\Sigma_s^{\frac{1}{2}})^T]$. We set $G = \Sigma_t^{\frac{1}{2}}\Sigma_s^{\frac{1}{2}}$, then we obtain $P_{ts} = \Sigma_t^{\frac{1}{2}} \Sigma_s^{\frac{1}{2}} (\Sigma_s^{\frac{1}{2}} \Sigma_t \Sigma_s^{\frac{1}{2}})^{\dagger} \Sigma_s^{\frac{1}{2}} \Sigma_t^{\frac{1}{2}}$ by Corollary 1. So,

$$\hat{\boldsymbol{T}}_{\text{OT}}\boldsymbol{\Sigma}_{s}\hat{\boldsymbol{T}}_{\text{OT}}^{T} = \boldsymbol{\Sigma}_{t}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}^{\frac{1}{2}} = \boldsymbol{\Sigma}_{t}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{\Sigma}_{t}.$$
(11)

Substituting $\Sigma_t = V_{n \times r} \Lambda_{r \times r} V_{n \times r}^T$ into the above, we get

$$\hat{\boldsymbol{T}}_{\text{OT}}\boldsymbol{\Sigma}_{s}\hat{\boldsymbol{T}}_{\text{OT}}^{T} = \boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}(\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{\dagger}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}}\boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{T}$$

$$= \boldsymbol{V}\boldsymbol{\Lambda}^{\frac{1}{2}}(\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})[(\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{T}(\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})]^{\dagger}(\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T}\boldsymbol{\Sigma}_{s}^{\frac{1}{2}})^{T}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T}$$

$$= \boldsymbol{V}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{P}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{V}^{T},$$
(12)

where P is the projection matrix onto $\text{Im}(\Lambda^{\frac{1}{2}}V^T\Sigma_sV\Lambda^{\frac{1}{2}}) = \text{Im}(\Lambda^{\frac{1}{2}}V^T\Sigma_s^{\frac{1}{2}})$, and the last equality in (12) holds because of Corollary 1. Note that $\Lambda^{\frac{1}{2}} V^T \Sigma_s^{\frac{1}{2}}$ is an $r \times n$ matrix, so $\operatorname{Im}(\Lambda^{\frac{1}{2}} V^T \Sigma_s^{\frac{1}{2}}) \subseteq \mathbb{R}^r$. In addition, by the above claim, $\operatorname{rank}(\Lambda^{\frac{1}{2}}V^{T}\Sigma_{s}^{\frac{1}{2}}) = \operatorname{rank}(V^{T}\Sigma_{s}^{\frac{1}{2}}) = r, \text{ so we have that } \operatorname{Im}(\Lambda^{\frac{1}{2}}V^{T}\Sigma_{s}^{\frac{1}{2}}) = \mathbb{R}^{r}. \text{ Therefore, the projection matrix } P \text{ is just the identity matrix } I_{r \times r}. \text{ Finally, Eq (12) can be written as } \hat{T}_{\mathrm{OT}}\Sigma_{s}\hat{T}_{\mathrm{OT}}^{T} = V\Lambda^{\frac{1}{2}}I_{r \times r}\Lambda^{\frac{1}{2}}V^{T} = V\Lambda V^{T} = \Sigma_{t}, \text{ and we get } I_{r \times r}$ the desired conclusion.

3.4. Proving Proposition 1

The maximum likelihood estimations of source and target covariance descriptors are given by

$$\mathcal{M}C_s = (\Phi_X \boldsymbol{J}_{N_s})(\Phi_X \boldsymbol{J}_{N_s})^T + \rho I_{\mathcal{H}\mathcal{K}}$$
(13a)

$$\mathbf{M}C_t = (\Phi_Y \boldsymbol{J}_{N_t})(\Phi_Y \boldsymbol{J}_{N_t})^T.$$
(13b)

Proposition 1. With the maximum likelihood estimators (13), the kernel whitening-coloring map is given by

$$k\hat{T}_{WC} = (MC_t)^{\frac{1}{2}} (MC_s)^{\frac{1}{2}} = \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\frac{1}{2}} [\boldsymbol{C}_{YX} \boldsymbol{A} \boldsymbol{J}_{N_s} \Phi_X^T + \frac{1}{\sqrt{\rho}} \boldsymbol{J}_{N_t} \Phi_Y^T],$$
(14)

where $C_{YY} = J_{N_t}^T K_{YY} J_{N_t}$ and $C_{YX} = J_{N_t}^T K_{YX} J_{N_s}$ are centered kernel matrices, and $A = \sum_{k=1}^r \frac{1}{\lambda_k} (\frac{1}{\sqrt{\lambda_k + \rho}} - \frac{1}{\sqrt{\rho}}) \vec{v}_k \vec{v}_k^T$, and $\{\lambda_k, \vec{v}_k\}_{k=1}^r$ are positive eigenpairs of C_{XX} .

Proof. We substitute maximum likelihood estimators (13) into the expressions of $k\hat{T}_{WC}$, and then we have,

$$k\hat{T}_{WC} = \left[(\Phi_Y \boldsymbol{J}_{N_t}) (\Phi_Y \boldsymbol{J}_{N_t})^T \right]^{\frac{1}{2}} \left[(\Phi_X \boldsymbol{J}_{N_s}) (\Phi_X \boldsymbol{J}_{N_s})^T + \rho I_{\mathcal{H}_{\mathcal{K}}} \right]^{-\frac{1}{2}}.$$
(15)

For simplicity, we set $G = \Phi_X J_{N_s}$, then by Lemma 3, $G^* = (\Phi_X J_{N_s})^T$. Let $\{\lambda_k\}_{k=1}^r$ and $\{\vec{v}_k\}_{k=1}^r$ be the positive eigenvalues and the corresponding orthonormal eigenvectors of $G^*G = (\Phi_X J_{N_s})^T (\Phi_X J_{N_s}) = C_{XX}$. Then by Lemma 2, $\{\lambda_k, \frac{G\vec{v}_k}{\sqrt{\lambda_k}}\}_{k=1}^r$ are the positive eigenpairs of GG^* . Let $\{\psi_k\}_{k=r+1}^{\dim(\mathcal{H}_{\mathcal{K}})}$ be a set of orthonormal vectors, such that

$$\left\{\frac{G\vec{v}_1}{\sqrt{\lambda_1}}, \frac{G\vec{v}_2}{\sqrt{\lambda_2}}, \dots, \frac{G\vec{v}_r}{\sqrt{\lambda_r}}, \psi_{r+1}, \psi_{r+2}, \dots\right\}$$
(16)

is a complete orthonormal basis of $\mathcal{H}_{\mathcal{K}}$. Then the identity operator $I_{\mathcal{H}_{\mathcal{K}}}$ can be written as $I_{\mathcal{H}_{\mathcal{K}}} = \sum_{k=1}^{r} \frac{G\vec{v}_{k}}{\sqrt{\lambda_{k}}} \otimes \frac{G\vec{v}_{k}}{\sqrt{\lambda_{k}}} + \sum_{k=r+1}^{\dim(\mathcal{H}_{\mathcal{K}})} \psi_{k} \otimes \psi_{k}$. Therefore, we have

$$\left[(\Phi_X \boldsymbol{J}_{N_s}) (\Phi_X \boldsymbol{J}_{N_s})^T + \rho I_{\mathcal{H}_{\mathcal{K}}} \right]^{-\frac{1}{2}}$$
(17a)

$$= (GG^* + \rho I_{\mathcal{H}_{\mathcal{K}}})^{-\frac{1}{2}} \tag{17b}$$

$$= \left[\sum_{k=1}^{r} \lambda_{k} \frac{G\vec{v}_{k}}{\sqrt{\lambda_{k}}} \otimes \frac{G\vec{v}_{k}}{\sqrt{\lambda_{k}}} + \rho I_{\mathcal{H}_{\mathcal{K}}}\right]^{-\frac{1}{2}}$$
(17c)

$$= \left[\sum_{k=1}^{r} (\lambda_k + \rho) \frac{G\vec{\boldsymbol{v}}_k}{\sqrt{\lambda_k}} \otimes \frac{G\vec{\boldsymbol{v}}_k}{\sqrt{\lambda_k}} + \sum_{k=r+1}^{\dim(\mathcal{H}_{\mathcal{K}})} \rho \psi_k \otimes \psi_k\right]^{-\frac{1}{2}}$$
(17d)

$$=\sum_{k=1}^{r} \frac{1}{\sqrt{\lambda_k + \rho}} \frac{G\vec{v}_k}{\sqrt{\lambda_k}} \otimes \frac{G\vec{v}_k}{\sqrt{\lambda_k}} + \sum_{k=r+1}^{\dim(\mathcal{H}_{\mathcal{K}})} \frac{1}{\sqrt{\rho}} \psi_k \otimes \psi_k$$
(17e)

$$=\sum_{k=1}^{r}\frac{1}{\sqrt{\lambda_{k}+\rho}}\frac{G\vec{\boldsymbol{v}}_{k}}{\sqrt{\lambda_{k}}}\otimes\frac{G\vec{\boldsymbol{v}}_{k}}{\sqrt{\lambda_{k}}}+\frac{1}{\sqrt{\rho}}(I_{\mathcal{H}_{\mathcal{K}}}-\sum_{k=1}^{r}\frac{G\vec{\boldsymbol{v}}_{k}}{\sqrt{\lambda_{k}}}\otimes\frac{G\vec{\boldsymbol{v}}_{k}}{\sqrt{\lambda_{k}}})$$
(17f)

$$=G\left(\sum_{k=1}^{r}\frac{1}{\lambda_{k}\sqrt{\lambda_{k}+\rho}}\vec{v}_{k}\vec{v}_{k}^{T}\right)G^{*}+\frac{1}{\sqrt{\rho}}\left[I_{\mathcal{H}_{\mathcal{K}}}-G\left(\sum_{k=1}^{r}\frac{1}{\lambda_{k}}\vec{v}_{k}\vec{v}_{k}^{T}\right)G^{*}\right]$$
(17g)

$$=GAG^* + \frac{1}{\sqrt{\rho}}I_{\mathcal{H}_{\mathcal{K}}}$$
(17h)

$$= \Phi_X \boldsymbol{J}_{N_s} \boldsymbol{A} \boldsymbol{J}_{N_s}^T \Phi_X^T + \frac{1}{\sqrt{\rho}} I_{\mathcal{H}_{\mathcal{K}}}.$$
(17i)

Substituting (17i) into (15), we get

$$k\hat{T}_{WC} = \left[(\Phi_Y \boldsymbol{J}_{N_t}) (\Phi_Y \boldsymbol{J}_{N_t})^T \right]^{\frac{1}{2}} \left[(\Phi_X \boldsymbol{J}_{N_s}) (\Phi_X \boldsymbol{J}_{N_s})^T + \rho I_{\mathcal{H}_{\mathcal{K}}} \right]^{-\frac{1}{2}}$$
(18a)

$$= \left[(\Phi_Y \boldsymbol{J}_{N_t}) (\Phi_Y \boldsymbol{J}_{N_t})^T \right]^{\frac{1}{2}} (\Phi_X \boldsymbol{J}_{N_s} \boldsymbol{A} \boldsymbol{J}_{N_s}^T \Phi_X^T + \frac{1}{\sqrt{\rho}} I_{\mathcal{H}_{\mathcal{K}}})$$
(18b)

$$= \Phi_Y \boldsymbol{J}_{N_t} \left[\boldsymbol{J}_{N_t}^T \Phi_Y^T \Phi_Y \boldsymbol{J}_{N_t} \right]^{\dagger \frac{1}{2}} \boldsymbol{J}_{N_t}^T \Phi_Y^T (\Phi_X \boldsymbol{J}_{N_s} \boldsymbol{A} \boldsymbol{J}_{N_s}^T \Phi_X^T + \frac{1}{\sqrt{\rho}} I_{\mathcal{H}_{\mathcal{K}}})$$
(18c)

$$= \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \big[\boldsymbol{C}_{YX} \boldsymbol{A} \boldsymbol{J}_{N_s} \Phi_X^T + \frac{1}{\sqrt{\rho}} \boldsymbol{J}_{N_t} \Phi_Y^T \big],$$
(18d)

where (18c) holds because of Corollary 2.

3.5. Proving Proposition 2

The computationally efficient estimations of source and target covariance descriptors are given by

$$EC_s = (\Phi_X \boldsymbol{W}_X)(\Phi_X \boldsymbol{W}_X)^T + \rho I_{\mathcal{H}_{\mathcal{K}}}$$
(19a)

$$EC_t = (\Phi_Y \boldsymbol{W}_Y)(\Phi_Y \boldsymbol{W}_Y)^T.$$
(19b)

Proposition 2. With the computationally efficient estimators (19), the kernel optimal transport map is given by

$$k\hat{T}_{\text{OT}} = (EC_t)^{\frac{1}{2}} \left[(EC_t)^{\frac{1}{2}} (EC_s) (EC_t)^{\frac{1}{2}} \right]^{\frac{1}{2}} (EC_t)^{\frac{1}{2}} = \Phi_Y W_Y \left[C_{YX}^w C_{XY}^w + \rho (\Lambda_Y - \rho I_d) \right]^{\frac{1}{2}} W_Y^T \Phi_Y^T,$$
(20)

where $C_{YX}^w = W_Y^T K_{YX} W_X$ and $C_{XY}^w = (C_{YX}^w)^T$, and Λ_Y is the diagonal matrix storing the top d eigenvalues of C_{YY} .

Proof. We substitute the computationally efficient estimators (19) into the expression of $k\hat{T}_{OT}$, and then we have,

$$k\hat{T}_{OT} = (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}} \left[(\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}} (\Phi_X W_X W_X^T \Phi_X^T + \rho I_{\mathcal{H}_{\mathcal{K}}}) (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}} \right]^{\frac{1}{2}} (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}}$$
(21)

We observe that $EC_s = \Phi_X W_X W_X^T \Phi_X^T + \rho I_{\mathcal{H}_{\mathcal{K}}}$ is strictly positive definite, which implies that there exists an operator \mathcal{A} (*e.g.*, $\mathcal{A} = (EC_s)^{\frac{1}{2}}$), such that $\mathcal{A}\mathcal{A}^T = EC_s$. Then $k\hat{T}_{OT}$ can be written as

$$\hat{kT}_{OT} = (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}} [(\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}} (\mathcal{A}\mathcal{A}^T) (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}}]^{\frac{1}{2}} (\Phi_Y W_Y W_Y^T \Phi_Y^T)^{\frac{1}{2}}$$
(22a)

$$= (\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \left(\left[(\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \mathcal{A} \right] \left[(\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \mathcal{A} \right]^{T} \right)^{\frac{1}{2}} (\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}}$$
(22b)

$$= (\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} [(\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \mathcal{A}] \left([(\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \mathcal{A}]^{T} [(\Phi_{Y} W_{Y} W_{Y}^{T} \Phi_{Y}^{T})^{\frac{1}{2}} \mathcal{A}] \right)^{\top \frac{1}{2}}$$
(22c)

$$\left[\left(\Phi_Y \boldsymbol{W}_Y \boldsymbol{W}_Y^T \Phi_Y^T\right)^{\frac{1}{2}} \mathcal{A}\right]^T \left(\Phi_Y \boldsymbol{W}_Y \boldsymbol{W}_Y^T \Phi_Y^T\right)^{\frac{1}{2}} \tag{22d}$$

$$= (\Phi_Y \boldsymbol{W}_Y \boldsymbol{W}_Y^T \Phi_Y^T) \mathcal{A} \left[\mathcal{A}^T \Phi_Y \boldsymbol{W}_Y \boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A} \right]^{\dagger \frac{3}{2}} \mathcal{A}^T (\Phi_Y \boldsymbol{W}_Y \boldsymbol{W}_Y^T \Phi_Y^T)$$
(22e)

$$= \Phi_Y \boldsymbol{W}_Y (\boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A}) \left[(\boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A})^T (\boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A}) \right]^{\dagger \frac{3}{2}} (\boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A})^T \boldsymbol{W}_Y^T \Phi_Y^T$$
(22f)

$$= \Phi_Y W_Y \left[(W_Y^T \Phi_Y^T \mathcal{A}) (W_Y^T \Phi_Y^T \mathcal{A})^T \right]^{\dagger \frac{1}{2}} W_Y^T \Phi_Y^T$$
(22g)

$$=\Phi_Y \boldsymbol{W}_Y (\boldsymbol{W}_Y^T \Phi_Y^T \mathcal{A} \mathcal{A}^T \Phi_Y \boldsymbol{W}_Y)^{\dagger \frac{1}{2}} \boldsymbol{W}_Y^T \Phi_Y^T$$
(22h)

$$=\Phi_Y W_Y \left[W_Y^T \Phi_Y^T (\Phi_X W_X W_X^T \Phi_X^T + \rho I_{\mathcal{H}_{\mathcal{K}}}) \Phi_Y W_Y \right]^{\dagger \frac{1}{2}} W_Y^T \Phi_Y^T$$
(22i)

$$=\Phi_Y W_Y (C_{YX}^w C_{XY}^w + \rho W_Y^T K_{YY} W_Y)^{\dagger \frac{1}{2}} W_Y^T \Phi_Y^T,$$
(22j)

where (22c) and (22g) hold because of Corollary 3. Next we prove that $W_Y^T K_{YY} W_Y = \Lambda_Y - \rho I_d$:

$$W_Y^T K_{YY} W_Y = \left[J_{N_t} V_Y (I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}} \right]^T K_{YY} \left[J_{N_t} V_Y (I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}} \right]$$

= $(I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}} V_Y^T C_{YY} V_Y^T (I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}}$
= $(I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}} \Lambda_Y (I_d - \rho \Lambda_Y^{-1})^{\frac{1}{2}}$
= $\Lambda_Y - \rho I_d.$ (23)

After substituting (23) into (22j), we obtain the desired result.

The domain-invariant kernel riangle ilde K is given by

$$\Delta \tilde{\boldsymbol{K}} = \begin{bmatrix} \Delta \tilde{\boldsymbol{K}}_{ss} & \Delta \tilde{\boldsymbol{K}}_{ts}^T \\ \Delta \tilde{\boldsymbol{K}}_{ts} & \Delta \tilde{\boldsymbol{K}}_{tt} \end{bmatrix} = \begin{bmatrix} \Psi_{s \to t}^T \Psi_{s \to t} & \Psi_{s \to t}^T \Psi_t \\ \Psi_t^T \Psi_{s \to t} & \Psi_t^T \Psi_t \end{bmatrix},$$
(24)

where the symbol \triangle represents the way of "moving" the source samples, *i.e.*, $\triangle = WC$ or OT, and $\Psi_{s \to t}$ denotes the transported source samples, *i.e.*, $\Psi_{s \to t} = kT_{\triangle}(\Psi_s)$, and Ψ_s and Ψ_t denote the centered source samples and target samples, respectively.

Using the kernel whitening-coloring map (14), we get

$$\Psi_{s \to t} = \mathbf{k} \hat{T}_{WC}(\Psi_s) = \sqrt{N_s} \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \boldsymbol{B}$$

$$WC \tilde{\boldsymbol{K}}_{ss} = N_s \boldsymbol{B}^T \boldsymbol{B}$$

$$WC \tilde{\boldsymbol{K}}_{ts} = \sqrt{N_s N_t} \boldsymbol{C}_{YY}^{\frac{1}{2}} \boldsymbol{B} = \sqrt{N_s N_t} \boldsymbol{U}_Y \boldsymbol{\Lambda}_Y^{\frac{1}{2}} \boldsymbol{U}_Y^T \boldsymbol{B},$$
(25)

where $\boldsymbol{B} = \boldsymbol{C}_{YX}(\boldsymbol{C}_{XX} + \rho \boldsymbol{I}_{N_s})^{-\frac{1}{2}}$, and $(\boldsymbol{U}_Y, \Lambda_Y^{\frac{1}{2}})$ stores the top *d* eigenpairs of \boldsymbol{C}_{YY} . Note that, in practice, in order to exploit the principal components and reduce the computational complexity, we artificially select *d* to be a small integer, *i.e.*, $d \ll N_t$.

Proof.

(I) The transported source samples are

$$\Psi_{s \to t} = \mathbf{k} T_{\mathrm{WC}}(\Psi_s) \tag{26a}$$

$$= \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \big[\boldsymbol{C}_{YX} \boldsymbol{A} \boldsymbol{J}_{N_s} \Phi_X^T + \frac{1}{\sqrt{\rho}} \boldsymbol{J}_{N_t} \Phi_Y^T \big] \big(\sqrt{N_s} \Phi_X \boldsymbol{J}_{N_s} \big)$$
(26b)

$$=\sqrt{N_s}\Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \boldsymbol{C}_{YX} (\boldsymbol{A}\boldsymbol{C}_{XX} + \frac{1}{\sqrt{\rho}} \boldsymbol{I}_{N_s}).$$
(26c)

Now we consider the term $AC_{XX} + \frac{1}{\sqrt{\rho}} I_{N_s}$. Recall that $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_r\}$ are C_{XX} 's eigenvectors, the corresponding eigenvalues of which are positive. Let $\{\vec{u}_{r+1}, \vec{u}_{r+2}, ..., \vec{u}_{N_s}\}$ be a set of orthonormal vectors, such that $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_r, \vec{u}_{r+1}, \vec{u}_{r+2}, ..., \vec{v}_{N_s}\}$

 $..., \vec{u}_{N_s}$ } is an orthonormal system for \mathbb{R}^{N_s} . Then, we have

$$\boldsymbol{A}\boldsymbol{C}_{XX} + \frac{1}{\sqrt{\rho}}\boldsymbol{I}_{N_s} = \sum_{k=1}^r \frac{1}{\lambda_k} (\frac{1}{\sqrt{\lambda_k + \rho}} - \frac{1}{\sqrt{\rho}}) \boldsymbol{\vec{v}}_k \boldsymbol{\vec{v}}_k^T \boldsymbol{C}_{XX} + \frac{1}{\sqrt{\rho}} \boldsymbol{I}_{N_s}$$
(27a)

$$=\sum_{k=1}^{r}\frac{1}{\lambda_{k}}\left(\frac{1}{\sqrt{\lambda_{k}+\rho}}-\frac{1}{\sqrt{\rho}}\right)\vec{v}_{k}(\boldsymbol{C}_{XX}\vec{v}_{k})^{T}+\frac{1}{\sqrt{\rho}}\boldsymbol{I}_{N_{s}}$$
(27b)

$$=\sum_{k=1}^{r} (\frac{1}{\sqrt{\lambda_{k}+\rho}} - \frac{1}{\sqrt{\rho}})\vec{v}_{k}\vec{v}_{k}^{T} + \frac{1}{\sqrt{\rho}}\sum_{k=1}^{r}\vec{v}_{k}\vec{v}_{k}^{T} + \frac{1}{\sqrt{\rho}}\sum_{k=r+1}^{N_{s}}\vec{u}_{k}\vec{u}_{k}^{T}$$
(27c)

$$=\sum_{k=1}^{r} \frac{1}{\sqrt{\lambda_k + \rho}} \vec{v}_k \vec{v}_k^T + \frac{1}{\sqrt{\rho}} \sum_{k=r+1}^{N_s} \vec{u}_k \vec{u}_k^T$$
(27d)

$$= \left[\sum_{k=1}^{r} (\lambda_k + \rho) \vec{v}_k \vec{v}_k^T + \sum_{k=r+1}^{N_s} \rho \vec{u}_k \vec{u}_k^T\right]^{-\frac{1}{2}}$$
(27e)

$$= (C_{XX} + \rho I_{N_s})^{-\frac{1}{2}}.$$
(27f)

Write $\boldsymbol{B} = \boldsymbol{C}_{YX} (\boldsymbol{C}_{XX} + \rho \boldsymbol{I}_{N_s})^{-\frac{1}{2}}$, then $\Psi_{s \to t} = k \hat{T}_{WC} (\Psi_s) = \sqrt{N_s} \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \boldsymbol{B}$. (II)

WC
$$\tilde{\mathbf{K}}_{ss} = \Psi_{s \to t}^T \Psi_{s \to t} = (\sqrt{N_s} \Phi_Y \mathbf{J}_{N_t} \mathbf{C}_{YY}^{\dagger \frac{1}{2}} \mathbf{B})^T (\sqrt{N_s} \Phi_Y \mathbf{J}_{N_t} \mathbf{C}_{YY}^{\dagger \frac{1}{2}} \mathbf{B})$$

= $N_s \mathbf{B}^T \mathbf{C}_{YY}^{\dagger \frac{1}{2}} \mathbf{J}_{N_t}^T \Phi_Y^T \Phi_Y \mathbf{J}_{N_t} \mathbf{C}_{YY}^{\dagger \frac{1}{2}} \mathbf{B} = N_s \mathbf{B}^T \mathbf{P}_{YY} \mathbf{B} = N_s \mathbf{B}^T \mathbf{B},$ (28)

where P_{YY} is the projection matrix onto $\text{Im}(C_{YY}) = \text{Im}(J_{N_t}^T \Phi_Y^T)$, and the last equality holds because $\text{Im}(B) \subseteq \text{Im}(C_{YX}) \subseteq \text{Im}(J_{N_t}^T \Phi_Y^T)$. (III)

WC
$$\tilde{\boldsymbol{K}}_{ts} = \Psi_t^T \Psi_{s \to t} = (\sqrt{N_t} \Phi_Y \boldsymbol{J}_{N_t})^T (\sqrt{N_s} \Phi_Y \boldsymbol{J}_{N_t} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \boldsymbol{B}) = \sqrt{N_t N_s} \boldsymbol{C}_{YY}^{\dagger \frac{1}{2}} \boldsymbol{B}.$$
 (29)

Using the kernel optimal transport map (20), we get

$$\Psi_{s \to t} = \mathbf{k} \hat{T}_{OT}(\Psi_s) = \sqrt{N_s} \Phi_Y W_Y D$$

$$OT \tilde{K}_{ss} = N_s D^T (\Lambda_Y - \rho I_d) D$$

$$OT \tilde{K}_{ts} = \sqrt{N_s N_t} J_{N_t} K_{YY} W_Y D,$$
(30)

where $\boldsymbol{D} = \left[\boldsymbol{C}_{YX}^{w}\boldsymbol{C}_{XY}^{w} + \rho(\boldsymbol{\Lambda}_{Y} - \rho\boldsymbol{I}_{d})\right]^{\dagger \frac{1}{2}} \boldsymbol{W}_{Y}^{T}\boldsymbol{K}_{YX}\boldsymbol{J}_{N_{s}}.$

Proof.

(I) The transported samples are

$$\Psi_{s \to t} = \mathbf{k} \hat{T}_{\mathrm{OT}}(\Psi_s) = \Phi_Y \mathbf{W}_Y \left[\mathbf{C}_{YX}^w \mathbf{C}_{XY}^w + \rho(\mathbf{\Lambda}_Y - \rho \mathbf{I}_d) \right]^{\dagger \frac{1}{2}} \mathbf{W}_Y^T \Phi_Y^T (\sqrt{N_s} \Phi_X \mathbf{J}_{N_s}) = \sqrt{N_s} \Phi_Y \mathbf{W}_Y \left[\mathbf{C}_{YX}^w \mathbf{C}_{XY}^w + \rho(\mathbf{\Lambda}_Y - \rho \mathbf{I}_d) \right]^{\dagger \frac{1}{2}} \mathbf{W}_Y^T \mathbf{K}_{YX} \mathbf{J}_{N_s} = \sqrt{N_s} \Phi_Y \mathbf{W}_Y \mathbf{D}.$$
(31)

(II)

$$OT\tilde{\boldsymbol{K}}_{ss} = \Psi_{s \to t}^{T} \Psi_{s \to t} = (\sqrt{N_s} \Phi_Y \boldsymbol{W}_Y \boldsymbol{D})^T (\sqrt{N_s} \Phi_Y \boldsymbol{W}_Y \boldsymbol{D}) = N_s \boldsymbol{D}^T \boldsymbol{W}_Y^T \boldsymbol{K}_{YY} \boldsymbol{W}_Y \boldsymbol{D} = N_s \boldsymbol{D}^T (\boldsymbol{\Lambda}_Y - \rho \boldsymbol{I}_d) \boldsymbol{D}, \quad (32)$$

where the last equality holds because of (23).

(III)

$$OT\tilde{\boldsymbol{K}}_{ts} = \Psi_t^T \Psi_{s \to t} = (\sqrt{N_t} \Phi_Y \boldsymbol{J}_{N_t})^T (\sqrt{N_s} \Phi_Y \boldsymbol{W}_Y \boldsymbol{D}) = \sqrt{N_s N_t} \boldsymbol{J}_{N_t} \boldsymbol{K}_{YY} \boldsymbol{W}_Y \boldsymbol{D}.$$
(33)

References

- [1] M. Harandi, M. Salzmann, and F. Porikli. Bregman divergences for infinite dimensional covariance matrices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1003–1010, 2014. 1
- [2] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2011. 2
- [3] M. H. Quang, M. San Biagio, and V. Murino. Log-hilbert-schmidt metric between positive definite operators on hilbert spaces. In Advances in Neural Information Processing Systems, pages 388–396, 2014.
- [4] B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998. 1, 3