



# **Optimal Transport in Reproducing Kernel Hilbert Space: Theory and Applications**

**IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE  
INTELLIGENCE (JULY 2020)**

**Zhen Zhang, Mianzhi Wang, and Arye Nehorai**  
**Reporter: Fengjiao Gong**

**June 9, 2022**

# Outline

1. Optimal Transport in  $\mathbb{R}^n$ 
  - 1.1 Monge's Formulation
  - 1.2 Kantorovich's Formulation
  - 1.3 OT between Gaussian Measures
2. Optimal Transport in RKHS
  - 2.1 Reproducing Kernel Hilbert Space(RKHS)
  - 2.2 Kantorovich's OT in RKHS
  - 2.3 OT between Gaussian Measures in RKHS
  - 2.4 Applications
    - 2.4.1 Image Classification
    - 2.4.2 Domain Adaptation
  - 2.5 Experiments

## Monge's Formulation

Find a **transport map**  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$  that pushes  $\mu$  to  $\nu$  (denoted as  $T_{\#}\mu = \nu$ ) to minimize the total transport cost

$$\inf_{T_{\#}\mu=\nu} \int_{\mathbb{R}^n} \|\vec{x} - T(\vec{x})\|_2^2 d\mu(\vec{x}) \quad (1)$$

where  $\mu, \nu \in \text{Pr}(\mathbb{R}^n)$  are two probability measures, and  $\text{Pr}(\mathbb{R}^n)$  is the set of Borel probability measures on  $\mathbb{R}^n$ .

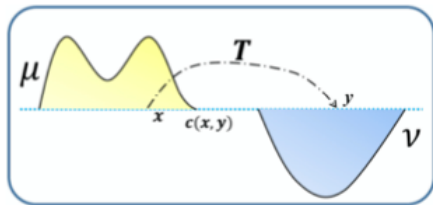


Figure 1: Illustration of the optimal transport problem.

# Monge's Formulation

Find a **transport map**  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$  that pushes  $\mu$  to  $\nu$  (denoted as  $T_{\#}\mu = \nu$ ) to minimize the total transport cost

$$\inf_{T_{\#}\mu=\nu} \int_{\mathbb{R}^n} \|\vec{x} - T(\vec{x})\|_2^2 d\mu(\vec{x}) \quad (1)$$

where

- ▶  $\text{Pr}(\mathbb{R}^n)$  — set of Borel probability measures on  $\mathbb{R}^n$
- ▶  $\mu, \nu \in \text{Pr}(\mathbb{R}^n)$  — two probability measures

**ISSUE:** Existence of  $T$  cannot be guaranteed!

**Kantorovich relaxation**

# Kantorovich's Formulation

Minimized over **all transport plans** instead of **transport map**

$$\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\vec{x} - \vec{y}\|_2^2 d\pi(\vec{x}, \vec{y}), \quad (2)$$

- ▶ Transport plan  $\pi(\vec{x}, \vec{y})$ : joint probability measure describing the amount of mass transported from location  $\vec{x}$  to location  $\vec{y}$
- ▶  $\Pi(\mu, \nu)$ : set of joint probability measures on  $\mathbb{R}^n \times \mathbb{R}^n$ , with marginals  $\mu, \nu$
- ▶ Splitting — mass at one location can be divided and transported to multiple destinations

# Kantorovich's Formulation

Minimized over **all transport plans** instead of **transport map**

$$\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\vec{x} - \vec{y}\|_2^2 d\pi(\vec{x}, \vec{y}), \quad (2)$$

**Wasserstein distance** — metric on  $\text{Pr}(\mathbb{R}^n)$

$$d_W(\mu, \nu) \triangleq \inf_{\pi \in \Pi(\mu, \nu)} \left[ \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\vec{x} - \vec{y}\|_2^2 d\pi(\vec{x}, \vec{y}) \right]^{\frac{1}{2}} \quad (3)$$

[reference] Cedric Villani, *Topics in Optimal Transportation*. Providence, RI, USA: American Mathematical Society, 2003, vol. 58.

# OT between Gaussian Measures

## Wasserstein distance

**Remark 1.**  $(\cdot)^{\frac{1}{2}}$  denotes the principle matrix square root, i.e., for any positive semi-definite (PSD) matrix  $\Sigma$ , then  $\Sigma^{\frac{1}{2}} = (U\Lambda U^T)^{\frac{1}{2}} = U\Lambda^{\frac{1}{2}}U^T$ .

[reference] D. Dowson and B. Landau, “The Fréchet distance between multivariate normal distributions,” *J. Multivariate Anal.*, vol. 12, no. 3, pp. 450–455, 1982.

# OT between Gaussian Measures

## Wasserstein distance

**Theorem 1 .** Let  $\mu$  and  $\nu$  be two probability measures on  $\mathbb{R}^n$  with finite first and second order moments. Let  $\vec{m}_\mu$  and  $\vec{m}_\nu$  , and  $\Sigma_\mu$  and  $\Sigma_\nu$  be the corresponding expectations and covariance matrices, respectively. Write

$$d_{\text{GaW}}(\mu, \nu) = \left[ \|\vec{m}_\mu - \vec{m}_\nu\|_2^2 + \text{tr}(\Sigma_\mu + \Sigma_\nu - 2\Sigma_{\mu\nu}) \right]^{\frac{1}{2}} \quad (4)$$

where  $\Sigma_{\mu\nu} = \left( \Sigma_\mu^{\frac{1}{2}} \Sigma_\nu \Sigma_\mu^{\frac{1}{2}} \right)^{\frac{1}{2}}$  . Then,

- 1)  $d_{\text{GaW}}(\mu, \nu) \leq d_{\text{W}}(\mu, \nu)$  , and
- 2) The equality will be valid if both  $\mu$  and  $\nu$  are Gaussian.



# OT between Gaussian Measures

## Wasserstein distance

**Theorem 1 .** Let  $\mu$  and  $\nu$  be two probability measures on  $\mathbb{R}^n$  with finite first and second order moments. Let  $\vec{m}_\mu$  and  $\vec{m}_\nu$ , and  $\Sigma_\mu$  and  $\Sigma_\nu$  be the corresponding expectations and covariance matrices, respectively. Write

$$d_{\text{GaW}}(\mu, \nu) = \left[ \|\vec{m}_\mu - \vec{m}_\nu\|_2^2 + \text{tr}(\Sigma_\mu + \Sigma_\nu - 2\Sigma_{\mu\nu}) \right]^{\frac{1}{2}} \quad (5)$$

- ▶  $d_{\text{GaW}}(\mu, \nu)$ : metric on set of all Gaussian measures — both Gaussian
- ▶  $d_{\text{GaW}}(\mu, \mu)$ : metric on covariance matrices — same expectations

## Bures metric

# OT between Gaussian Measures

## Wasserstein distance

**Corollary 1.** Let  $\text{Sym}^+(n)$  be the set of all positive semi-definite matrices of size  $n \times n$ . Then,

$$d_B(\Sigma_1, \Sigma_2) = [\text{tr}(\Sigma_1 + \Sigma_2 - 2\Sigma_{12})]^{1/2} \quad (6)$$

defines a metric on  $\text{Sym}^+(n)$ .

*[reference] R. Bhatia, T. Jain, and Y. Lim, “On the Bures Wasserstein distance between positive definite matrices,” Expositiones Mathematicae, 2018.*

# OT between Gaussian Measures

## Optimal transport map

**Theorem 2 .** Let  $\mu$  and  $\nu$  be two Gaussian measures on  $\mathbb{R}^n$  whose covariance matrices are of full rank. Let  $\vec{m}_\mu$  and  $\vec{m}_\nu$ , and  $\Sigma_\mu$  and  $\Sigma_\nu$ , denote the respective expectations and covariance matrices. Then the optimal transport map  $T_G$  between  $\mu$  and  $\nu$  exists, and can be written as

$$T_G(\vec{x}) = \Sigma_\mu^{-\frac{1}{2}} \Sigma_{\mu\nu} \Sigma_\mu^{-\frac{1}{2}} (\vec{x} - \vec{m}_\mu) + \vec{m}_\nu \quad (7)$$

Covariance matrices are *full-rank*.

*[reference] R. Bhatia, T. Jain, and Y. Lim, "On the Bures Wasserstein distance between positive definite matrices," Expositiones Mathematicae, 2018.*

# OT between Gaussian Measures

## Optimal transport map

**Remark 2.** “ $\dagger$ ” denotes the Moore-Penrose inverse.  $\text{Im}(\Sigma)$  denotes the image of the linear transform  $\Sigma$ , i.e.,  $\text{Im}(\Sigma) = \{\Sigma\vec{x}, \vec{x} \in \mathbb{R}^n\}$ .

*[reference]*

*[https://en.wikipedia.org/wiki/Moore%E2%80%93Penrose\\_inverse](https://en.wikipedia.org/wiki/Moore%E2%80%93Penrose_inverse)*

*M. Gelbrich, “On a formula for the L2 Wasserstein metric between measures on Euclidean and hilbert spaces,” Mathematische Nachrichten, vol. 147, no. 1, pp. 185–203, 1990.*

# OT between Gaussian Measures

## Optimal transport map

**Theorem 3.** Let  $\mu$  and  $\nu$  be two Gaussian measures defined on  $\mathbb{R}^n$ . Let  $\bar{\mu}$  and  $\bar{\nu}$  be the corresponding centered Gaussian measures which are derived from  $\mu$  and  $\nu$ , respectively, by translation. Let  $\mathbf{P}_\mu$  be the projection matrix onto  $\text{Im}(\Sigma_\mu)$ . Then the optimal transport map  $T_G$  from  $\bar{\mu}$  to  $P_{\mu\#}\bar{\nu}$  is linear and self-adjoint, and can be written as

$$T_G(\vec{\mathbf{x}}) = \left(\Sigma_\mu^{\frac{1}{2}}\right)^\dagger \Sigma_{\mu\nu} \left(\Sigma_\mu^{\frac{1}{2}}\right)^\dagger \vec{\mathbf{x}} \quad (8)$$

Covariance matrix is *rank-deficient* ( $\bar{\mu} \Rightarrow P_{\mu\#}\bar{\nu}$ ).

# Optimal Transport in Reproducing Kernel Hilbert Space

1. Reproducing Kernel Hilbert Space(RKHS)
2. Kantorovich'S OT in RKHS
3. OT between Gaussian Measures in RKHS
4. Applications
  - 4.1 Image Classification
  - 4.2 Domain Adaptation
5. Experiments

# Reproducing Kernel Hilbert Space(RKHS)

## ► RKHS

Function  $k : \mathcal{X} \times \mathcal{X}$  is called a **reproducing kernel** of  $\mathcal{H}$ , and  $\mathcal{H}$  is a **reproducing kernel Hilbert space**, if  $k$  satisfies:

- 1)  $\forall x \in \mathcal{X}, k(\cdot, x) \in \mathcal{H}$ ,
- 2)  $\forall x \in \mathcal{X}, f \in \mathcal{H}, \langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$ .

where  $\mathcal{H}$  is a Hilbert space of  $\mathbb{R}$ -valued functions defined on nonempty set  $\mathcal{X}$ .

## ► Feature map

Implicit feature map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$

$$\phi(x) = k(\cdot, x)$$

Then

$$\langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = k(x, y), \forall x, y \in \mathcal{X}$$

# Kantorovich'S OT in RKHS

## OT Formulation

Given  $\mu, \nu \in \Pr(\mathcal{X})$ , the Kantorovich optimal transport between **pushforward measures**  $\phi_{\#}\mu$  and  $\phi_{\#}\nu$  on  $\mathcal{H}_{\mathcal{K}}$  is written as

$$d_W(\phi_{\#}\mu, \phi_{\#}\nu) = \left[ \inf_{\pi_{\mathcal{K}} \in \Pi(\phi_{\#}\mu, \phi_{\#}\nu)} \int_{\mathcal{H}_{\mathcal{K}} \times \mathcal{H}_{\mathcal{K}}} \|u - v\|_{\mathcal{H}_{\mathcal{K}}}^2 d\pi_{\mathcal{K}}(u, v) \right]^{\frac{1}{2}} \quad (9)$$

where

- ▶  $\Pi(\phi_{\#}\mu, \phi_{\#}\nu)$ : set of joint probability measures on  $\mathcal{H}_{\mathcal{K}} \times \mathcal{H}_{\mathcal{K}}$  with marginals  $\phi_{\#}\mu, \phi_{\#}\nu$
- ▶  $k$ : positive definite kernel on  $\mathcal{X} \times \mathcal{X}$
- ▶  $(\mathcal{H}_{\mathcal{K}}, \mathcal{B}_{\mathcal{H}_{\mathcal{K}}})$ : reproducing kernel Hilbert space generated by  $k$



# Kantorovich'S OT in RKHS

## OT Formulation

**Theorem 4.** Let  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  be a Borel space, and let the reproducing kernel  $k$  be measurable. Given  $\mu, \nu \in \text{Pr}(\mathcal{X})$ , we write

$$d_{\mathcal{W}}^{\mathcal{H}}(\mu, \nu) = \left[ \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} d^2(x, y) d\pi(x, y) \right]^{\frac{1}{2}} \quad (10)$$

where  $d^2(x, y) = \|\phi(x) - \phi(y)\|_{\mathcal{H}_{\mathcal{K}}}^2 = k(x, x) + k(y, y) - 2k(x, y)$ . Then,

1)  $d_{\mathcal{W}}^{\mathcal{H}}(\mu, \nu) = d_{\mathcal{W}}(\phi_{\#}\mu, \phi_{\#}\nu)$ , and

2) If  $\pi^*$  is a minimizer of (10), then  $(\phi, \phi)_{\#}\pi^*$  is a minimizer of (9), where

$$(\phi, \phi)(x, y) = (\phi(x), \phi(y))$$

and  $(\phi, \phi) : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{H}_{\mathcal{K}} \times \mathcal{H}_{\mathcal{K}}$

# Kantorovich'S OT in RKHS

## Discrete OT

Discrete version of (10) can be written as

$$\min_{P \in U_{nm}} \text{tr} \left( P^T \mathbf{D} \right) \quad (11)$$

where

- ▶  $U_{nm}$  — set of  $n \times m$  non-negative matrices representing probabilistic couplings with marginals  $\hat{\mu}, \hat{\nu}$ ,

$$U_{nm} = \left\{ \mathbf{P} \in \mathbb{R}_+^{n \times m} \mid \mathbf{P} \mathbf{1}_m = \hat{\mu}, \mathbf{P}^T \mathbf{1}_n = \hat{\nu} \right\}$$

- ▶  $\mathbf{D}$  cost matrix

$$\mathbf{D}_{i,j} = k(x_i, x_i) + k(y_j, y_j) - 2k(x_i, y_j)$$

- ▶ Empirical histograms:

$$\hat{\mu} = \sum_{i=1}^n \hat{\mu}_i \delta_{x_i}, \quad \hat{\nu} = \sum_{j=1}^m \hat{\nu}_j \delta_{y_j}$$

# OT between Gaussian Measures in RKHS

Let  $\mu$  be a Borel probability measure on  $\mathcal{X}$ , and

- ▶ Mean:  $m_\mu = E_{X \sim \mu}(\phi(X))$
- ▶ Covariance operator:  $R_\mu = E_{X \sim \mu} ((\phi(X) - m_\mu) \otimes (\phi(X) - m_\mu))^2$

exist and be bounded with respect to the Hilbert norm and HilbertSchmidt norm, respectively.

Tensor product of  $\mathcal{H}$ :

$$(u \otimes v)w = \langle v, w \rangle_{\mathcal{H}} u, \quad \forall u, v, w \in \mathcal{H}$$

*[reference] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, “Measuring statistical dependence with hilbert-schmidt norms,” in Proc. 16th Int. Conf. Algorithmic Learn. Theory, 2005, pp. 63–77.*

# OT between Gaussian Measures in RKHS

## Remark 3 .

(1) The square root of a nonnegative, self-adjoint, and compact operator  $G$  is defined as  $G^{\frac{1}{2}} = \sum_{i=1} \sqrt{\lambda_i(G)} \varphi_i(G) \otimes \varphi_i(G)$ , where  $\lambda_i(G)$  and  $\varphi_i(G)$  are eigenvalues and eigenfunctions of  $G$ .

(2) The trace of a trace-class operator  $G$  on a separable Hilbert space  $\mathcal{H}$  is defined as  $\text{tr}(G) = \sum_{i=1}^{\dim(\mathcal{H})} \langle G e_i, e_i \rangle$ , where  $\{e_i\}_{i=1}^{\dim(\mathcal{H})}$  is an orthonormal system of  $\mathcal{H}$ .

If data distributions in RKHS (the corresponding pushforward measures) are **Gaussian**, the conclusions in RKHS are **similar** to the ones in Euclidean spaces.

# OT between Gaussian Measures in RKHS

## KGW distance

**Proposition 1.** Assume that the hypotheses in Theorem 4 hold. Let  $\mu, \nu \in \text{Pr}(\mathcal{X})$ . Let  $m_\mu$  and  $m_\nu$ , and  $R_\mu$  and  $R_\nu$ , be the corresponding means and covariance operators, respectively. Write

$$d_{\text{GW}}^{\mathcal{H}}(\mu, \nu) = \left[ \|m_\mu - m_\nu\|_{\mathcal{H}_{\mathcal{K}}}^2 + \text{tr}(R_\mu + R_\nu - 2R_{\mu\nu}) \right]^{\frac{1}{2}} \quad (12)$$

where  $R_{\mu\nu} = \left( R_\mu^{\frac{1}{2}} R_\nu R_\mu^{\frac{1}{2}} \right)^{\frac{1}{2}}$ . Then,

(1)  $d_{\text{GW}}^{\mathcal{H}}(\mu, \nu) \leq d_{\text{W}}^{\mathcal{H}}(\mu, \nu)$ , and

(2) The equality will be valid if both  $\phi_{\#}\mu$  and  $\phi_{\#}\nu$  are Gaussian.

# OT between Gaussian Measures in RKHS

**Corollary 2.** Let  $\text{Sym}^+(\mathcal{H}_{\mathcal{K}}) \subseteq \mathcal{H}_{\mathcal{K}} \otimes \mathcal{H}_{\mathcal{K}}$  be the set of nonnegative, self-adjoint, and trace-class operators in  $\mathcal{H}_{\mathcal{K}}$ . Then

$$d_{\text{B}}^{\mathcal{H}}(R_1, R_2) = [\text{tr}(R_1 + R_2 - 2R_{12})]^{\frac{1}{2}} \quad (13)$$

defines a metric on  $\text{Sym}^+(\mathcal{H}_{\mathcal{K}})$ .

*$d_{\text{B}}^{\mathcal{H}}$  quantifies the difference between the dispersions of data in RKHS.*

*$d_{\text{B}}^{\mathcal{H}}$  quantifies the discrepancy between distributions*

# OT between Gaussian Measures in RKHS

**Definition 1.** Let  $\mu \in \text{Pr}(\mathcal{X})$ . If there exist disjoint subsets  $\Omega_1$ ,  $\Omega_2$ , and  $\Omega_3$ , satisfying  $\mathcal{X} = \Omega_1 \cup \Omega_2 \cup \Omega_3$ , and  $\mu(\Omega_1), \mu(\Omega_2), \mu(\Omega_3) > 0$ , then we say  $\mu$  satisfies the 3-splitting property.

Here  $\text{Pr}^s(\mathcal{X})$  is the set of Borel measures satisfying the 3-splitting property.

# OT between Gaussian Measures in RKHS

**Theorem 5.** Let the measurable space  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  be locally compact and Hausdorff. Let  $k$  be a  $c_0$ -universal reproducing kernel. Then, the embedding  $\mu \rightarrow R_{\mu}, \forall \mu \in \text{Pr}^s(\mathcal{X})$  is injective.

$d_B^{\mathcal{H}}$  (KB) induces a metric on  $\text{Pr}^s(\mathcal{X})$ .

[reference]

[https://en.wikipedia.org/wiki/Compact\\_space](https://en.wikipedia.org/wiki/Compact_space)

[https://en.wikipedia.org/wiki/Hausdorff\\_space](https://en.wikipedia.org/wiki/Hausdorff_space)

B. K. Sriperumbudur, K. Fukumizu, and G. R. Lanckriet, “Universality, characteristic kernels and RKHS embedding of measures,” *J. Mach. Learn. Res.*, vol. 12, pp. 2389–2410, Jul. 2011.



# OT between Gaussian Measures in RKHS

## KGOP map

**Proposition 2.** Given  $\mu, \nu \in \text{Pr}(\mathcal{X})$ , assume the pushforward measures  $\phi_{\#}\mu$  and  $\phi_{\#}\nu$  on RKHS are Gaussian. Let  $\bar{\mu}_{\phi}$  and  $\bar{\nu}_{\phi}$  be the respective centered measures of  $\phi_{\#}\mu$  and  $\phi_{\#}\nu$ . Let  $P_{\mu}$  be the projection operator on  $\text{Im}(R_{\mu})$ . Then the kernel Gauss-optimal transport map  $T_G^{\mathcal{H}}$  between  $\bar{\mu}_{\phi}$  and  $P_{\mu\#}(\bar{\nu}_{\phi})$  is a linear and self-adjoint operator, and can be written as

$$T_G^{\mathcal{H}}(u) = \left(R_{\mu}^{\frac{1}{2}}\right)^{\dagger} R_{\mu\nu} \left(R_{\mu}^{\frac{1}{2}}\right)^{\dagger} u, \forall u \in \mathcal{H}_{\mathcal{K}} \quad (14)$$

# OT between Gaussian Measures in RKHS

Sample matrices

$$\mathbf{X} = [x_1, x_2, \dots, x_n], \quad \mathbf{Y} = [y_1, y_2, \dots, y_m]$$

Mapped data matrices

$$\Phi_X = [\phi(x_1), \phi(x_2), \dots, \phi(x_n)], \quad \Phi_Y = [\phi(y_1), \phi(y_2), \dots, \phi(y_m)]$$

Kernel matrices

$$(K_{XX})_{ij} = k(x_i, x_j), \quad (K_{XY})_{ij} = k(x_i, y_j), \quad (K_{YY})_{ij} = k(y_i, y_j)$$

Centering matrices

$$\mathbf{H}_n = \mathbf{I}_{n \times n} - \frac{1}{n} \vec{\mathbf{1}}_n \vec{\mathbf{1}}_n^T, \quad \mathbf{H}_m = \mathbf{I}_{m \times m} - \frac{1}{m} \vec{\mathbf{1}}_m \vec{\mathbf{1}}_m^T$$

# OT between Gaussian Measures in RKHS

Estimated mean:

$$\hat{m}_\mu = \frac{1}{n} \Phi_X \vec{1}_n, \quad \hat{m}_\nu = \frac{1}{m} \Phi_Y \vec{1}_m$$

Estimated covariance operators:

$$\hat{R}_\mu = \frac{1}{n} \Phi_X H_n \Phi_X^T, \quad \hat{R}_\nu = \frac{1}{m} \Phi_Y H_m \Phi_Y^T$$

**Remark 4.**  $\|\cdot\|_*$  denotes the nuclear norm, i.e.,  $\|A\|_* = \sum_{i=1}^r \sigma_i(A)$ , where  $\sigma_i(A)$  are the singular values of matrix  $A$

# OT between Gaussian Measures in RKHS

## KGW Distance Empirical Estimation

**Proposition 3.** The empirical kernel Gauss-Wasserstein distance is

$$\hat{d}_{\text{GW}}^{\mathcal{H}}(\mu, \nu) = \left[ \frac{1}{n} \text{tr}(\mathbf{K}_{XX}) + \frac{1}{m} \text{tr}(\mathbf{K}_{YY}) - \frac{2}{mn} \vec{\mathbf{1}}_n^T \mathbf{K}_{XY} \vec{\mathbf{1}}_m - \frac{2}{\sqrt{mn}} \|\mathbf{H}_n \mathbf{K}_{XY} \mathbf{H}_m\|_* \right]^{\frac{1}{2}} \quad (15)$$

The kernel Bures distance between  $\hat{R}_\mu$  and  $\hat{R}_\nu$  is

$$d_{\text{B}}^{\mathcal{H}}(\hat{R}_\mu, \hat{R}_\nu) = \left[ \frac{1}{n} \text{tr}(\mathbf{K}_{XX} \mathbf{H}_n) + \frac{1}{m} \text{tr}(\mathbf{K}_{YY} \mathbf{H}_m) - \frac{2}{\sqrt{mn}} \|\mathbf{H}_n \mathbf{K}_{XY} \mathbf{H}_m\|_* \right]^{\frac{1}{2}} \quad (16)$$

# OT between Gaussian Measures in RKHS

## Empirical Estimation of the KGOT Map

**Proposition 4.** Let  $X$  and  $Y$  be data matrices sampled from  $\mu$  and  $\nu$ , respectively. Then the empirical projection operator on  $\text{Im}(\hat{R}_\mu)$  is

$$\hat{P}_\mu = \Phi_X \mathbf{H}_n \mathbf{C}_{XX}^\dagger \mathbf{H}_n \Phi_X^T \quad (17)$$

and the empirical Gauss-optimal transport map from  $\bar{\mu}_\phi$  and  $P_{\mu\#}(\bar{\nu}_\phi)$  is

$$\hat{T}_G^{\mathcal{H}} = \sqrt{\frac{n}{m}} \Phi_X \mathbf{H}_n \mathbf{C}_{XX}^\dagger \mathbf{C}_{XY}^{\frac{1}{2}} \mathbf{C}_{YX}^\dagger \mathbf{H}_m \Phi_Y^T \quad (18)$$

where

$$\begin{aligned} \mathbf{C}_{XX} &= \mathbf{H}_n \mathbf{K}_{XX} \mathbf{H}_n \\ \mathbf{C}_{XY} &= \mathbf{H}_n \mathbf{K}_{XY} \mathbf{H}_m \\ \mathbf{C}_{YX} &= \mathbf{H}_m \mathbf{K}_{YX} \mathbf{H}_n \end{aligned} \quad (19)$$

# OT between Gaussian Measures in RKHS

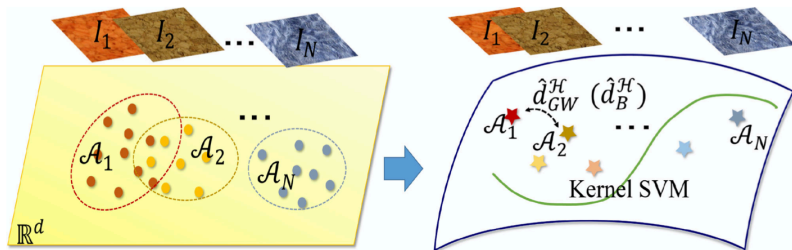
## KGOT Map

**Proposition 5 .**

$$\hat{T}_G^{\mathcal{H}} \hat{R}_\mu \hat{T}_G^{\mathcal{H}} = \hat{P}_\mu \hat{R}_\nu \hat{P}_\mu \quad (20)$$

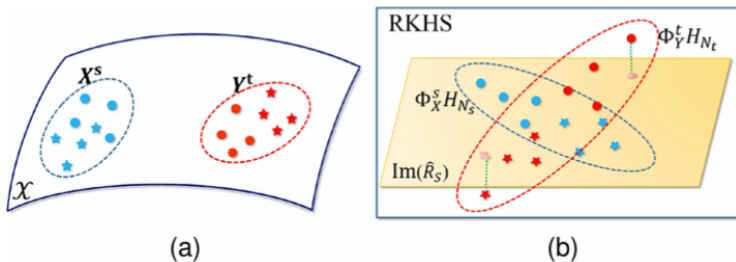
Align covariance operators in RKHS.

# Application — Image Classification



- Represent each image  $I_i$  by a collection of feature samples  $\mathcal{A}_i$ .
- Compute the KGW (or the KB) distances between any pair of images.
- Apply kernel SVM to conduct classification.

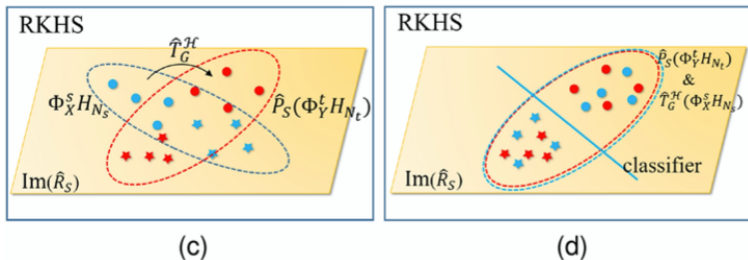
## Application – Domain Adaptation



- (a) The labeled dataset,  $X^s$ , in the source domain and the unlabeled dataset,  $Y^t$ , in the target domain. Dots and stars represent different classes;
- (b) Map  $X^s$  and  $Y^t$  to the RKHS  $\mathcal{H}_K$ , and centralize the mapped data. (The centered source dataset  $\Phi_X^s H_{N_s}$  lies in  $\text{Im}(\hat{R}_s)$ );

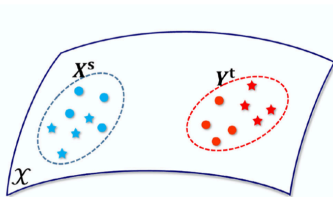


# Application – Domain Adaptation

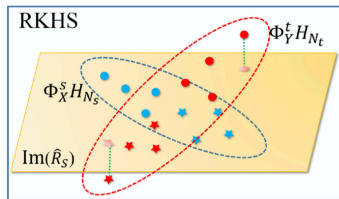


- (c) Project the target dataset  $\Phi_Y^t H_{N_t}$  onto  $\text{Im}(\hat{R}_S)$ . The projection is  $\hat{P}_S(\Phi_Y^t H_{N_t})$ ;
- (d) Apply the KGOT map to transport the source data to the target domain. The transported data is  $\hat{T}_G^{\mathcal{H}}(\Phi_X^s H_{N_s})$ .

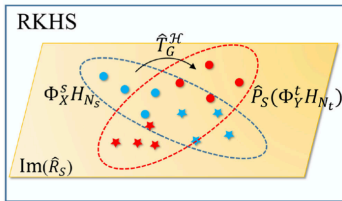
# Domain Adaptation



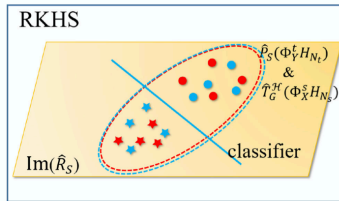
(a)



(b)



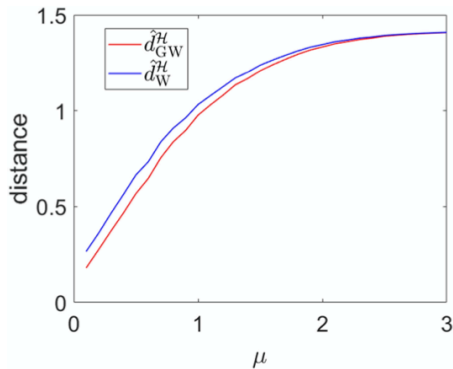
(c)



(d)

Finally, train a classifier using  $\hat{T}_G^{\mathcal{H}}(\Phi_X^s H_{N_s})$ , then apply to  $\hat{P}_S(\Phi_Y^t H_{N_t})$ .

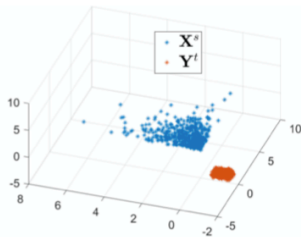
# Experiments



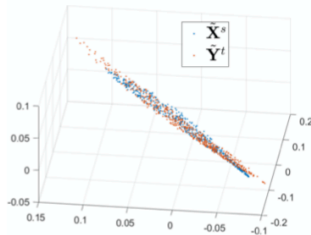
**Figure 2:** The estimated KGW and KW distances between Gaussian distributions  $N(m \vec{1}, I)$  and  $N(-m \vec{1}, I)$ .

Clearly, KGW is less than KW.

# Experiments



(a)



(b)

(a) Source dataset  $X^s$ , and target dataset  $Y^t$  ;  
(b) Representations of datasets  $\hat{T}_G^{\mathcal{H}} (\Phi_X^s H_{N_s})$   
and  $\hat{P}_s (\Phi_Y^t H_{N_t})$  under coordinate system  
 $(l_i)_{i=1}^3$ .

Distributions of  $\tilde{X}^s$  and  $\tilde{Y}^t$  are close to each other.

*Thanks!*