

# Information geometry of Wasserstein statistics on shapes and affine deformations

Shun-ichi Amari, Takeru Matsuda

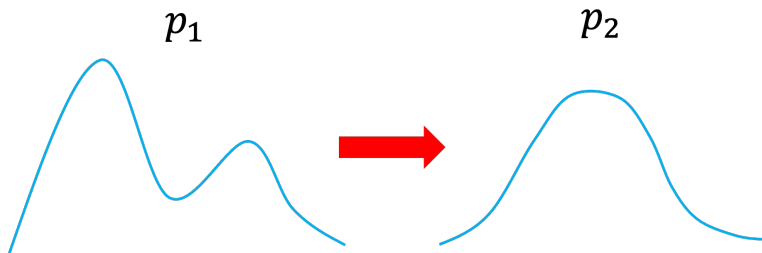
Information Geometry, 2024

# Wasserstein distance

- $L^2$  Wasserstein distance (= optimal transportation cost) between  $p_1$  and  $p_2$  on  $\mathbb{R}^d$

$$W_2(p_1, p_2) = \inf_{X_1, X_2} \mathbb{E}[\|X_1 - X_2\|^2]^{1/2}$$

- ▶ infimum over all joint distributions of  $(X_1, X_2)$  with  $X_1 \sim p_1$  and  $X_2 \sim p_2$  marginally (coupling)



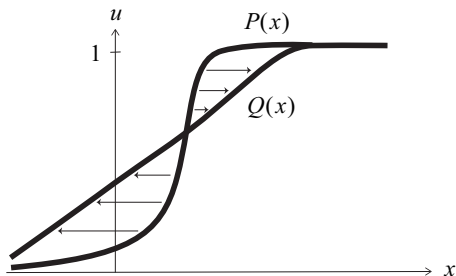
## One-dimensional case

- When  $d = 1$ ,  $W_2$  is explicitly given by the cdfs  $P_1$  and  $P_2$ :

$$W_2(p_1, p_2) = \left( \int_0^1 (P_1^{-1}(u) - P_2^{-1}(u))^2 du \right)^{1/2}$$

- optimal coupling = monotone map

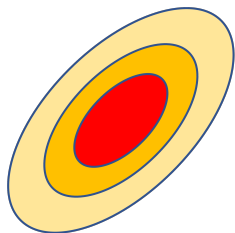
$$X_2 = P_2^{-1}(P_1(X_1))$$



## Elliptically contoured family

- When  $d \geq 2$ ,  $W_2$  is intractable in general..
- elliptically contoured family (e.g. Gaussian)
  - $\mu$ : mean,  $\Sigma$ : covariance,  $f$ : shape

$$p(x \mid \mu, \Sigma) = (\det \Sigma)^{-1/2} f(\|\Sigma^{-1/2}(x - \mu)\|)$$



### Proposition (Gelbrich, 1990)

$$\begin{aligned} & W_2(p(x \mid \mu_1, \Sigma_1), p(x \mid \mu_2, \Sigma_2)) \\ &= \left( \|\mu_1 - \mu_2\|^2 + \text{tr} \left( \Sigma_1 + \Sigma_2 - 2(\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \right) \right)^{1/2} \end{aligned}$$

- note:  $W_2$  does not depend on the shape  $f$

# Wasserstein v.s. Kullback–Leibler

- bijective variable transformation

$$y = g(x) \quad \rightarrow \quad \tilde{p}(y) = \left| \frac{dx}{dy} \right| p(x)$$

- Kullback–Leibler divergence: invariant

$$D_{\text{KL}}(p, q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

$$D_{\text{KL}}(\tilde{p}, \tilde{q}) = D_{\text{KL}}(p, q)$$

- Wasserstein distance: **not** invariant

$$W_2(\tilde{p}, \tilde{q}) \neq W_2(p, q)$$

## Li–Zhao framework

- Recently, Li and Zhao (2023) developed Wasserstein counterparts of information geometric concepts

Kullback–Leibler divergence	Wasserstein distance
Fisher score	Wasserstein score
Fisher information matrix	Wasserstein information matrix
covariance	Wasserstein covariance
Cramer–Rao	Wasserstein–Cramer–Rao
Fisher efficiency	Wasserstein efficiency

- We investigate their statistical meaning

# Continuity equation

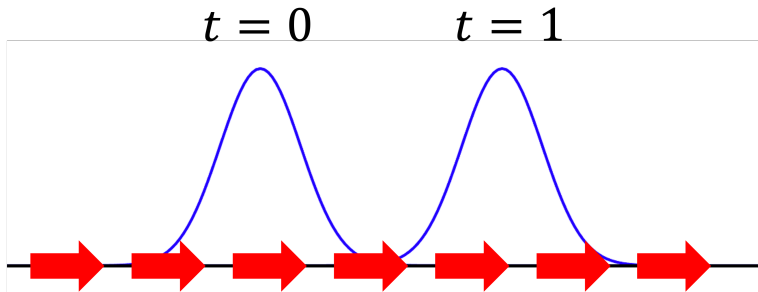
$$\frac{\partial}{\partial t} p(x, t) = -\nabla_x \cdot (p(x, t) \nabla_x \Phi(x))$$

- This PDE describes **dynamics of measure transport**
- intuition: Many particles are distributed with  $p(x, t)$  and they move with velocity  $\nabla_x \Phi(x)$

## Example: 1d linear potential

$$\frac{\partial}{\partial t} p(x, t) = -\nabla_x \cdot (p(x, t) \nabla_x \Phi(x))$$

- $\Phi(x) = x \rightarrow \nabla_x \Phi(x) \equiv 1$  (const.)
- $p(x, 0) = \mathcal{N}(0, 1) \rightarrow p(x, t) = \mathcal{N}(t, 1)$  (shift)

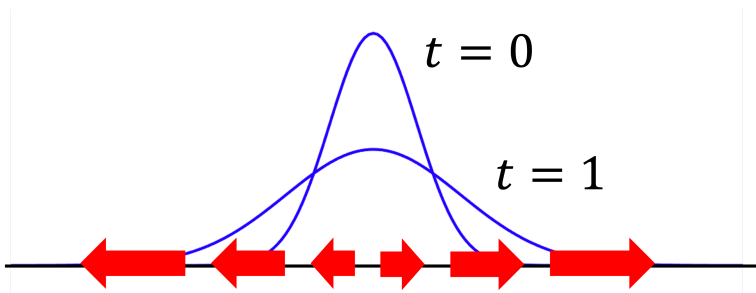




## Example: 1d quadratic potential

$$\frac{\partial}{\partial t} p(x, t) = -\nabla_x \cdot (p(x, t) \nabla_x \Phi(x))$$

- $\Phi(x) = x^2 \rightarrow \nabla_x \Phi(x) = 2x$
- $p(x, 0) = \mathcal{N}(0, 1) \rightarrow p(x, t) = \mathcal{N}(0, t + 1)$  (expansion)



# Wasserstein score function

## Definition (Li and Zhao, 2023)

For  $i = 1, \dots, p$ , the **Wasserstein score function**  $\Phi_i^W(x | \theta)$  is the solution of

$$-\nabla_x \cdot (p(x | \theta) \nabla_x \Phi_i^W(x | \theta)) = \frac{\partial}{\partial \theta_i} p(x | \theta), \quad \mathbb{E}_\theta[\Phi_i^W(x | \theta)] = 0.$$

- For infinitesimal  $\delta$ , the map  $x \mapsto x + \delta \nabla_x \Phi_i^W(x | \theta)$  is the optimal transport map from  $p(x | \theta)$  to  $p(x | \theta + \delta e_i)$  with transportation cost

$$W_2(p(x | \theta), p(x | \theta + \delta e_i)) = \left( \int \|\delta \nabla_x \Phi_i^W(x | \theta)\|^2 p(x | \theta) dx \right)^{1/2}$$

- $e_i$ :  $i$ -th standard unit vector

# Wasserstein information matrix (WIM)

## Definition (Li and Zhao, 2023)

The **Wasserstein information matrix**  $G_W(\theta)$  is the  $p \times p$  matrix given by

$$G_W(\theta) = \left( \int \frac{\partial}{\partial \theta_i} p(x | \theta) \cdot \Phi_j^W(x | \theta) dx \right)_{ij}$$

- cf. Fisher information matrix

$$G_F(\theta) = \left( \int \frac{\partial}{\partial \theta_i} p(x | \theta) \cdot \Phi_j^F(x | \theta) dx \right)_{ij}$$

$$\Phi_j^F(x | \theta) = \frac{\partial}{\partial \theta_j} \log p(x | \theta)$$

- inner product = pairing of **tangent vector** and **cotangent vector**
  - information geometry: **m-representation** and **e-representation**

# Wasserstein information matrix (WIM)

## Proposition (Li and Zhao, 2023)

$$G_W(\theta)_{ij} = \mathbb{E}_\theta[(\nabla_x \Phi_i^W(x | \theta))^\top (\nabla_x \Phi_j^W(x | \theta))]$$

## Proposition (Li and Zhao, 2023)

$$W_2(p(x | \theta), p(x | \theta + \delta))^2 = \delta^\top G_W(\theta) \delta + o(\|\delta\|^2)$$

- WIM = Hessian of Wasserstein distance
  - cf. Fisher information matrix = Hessian of Kullback–Leibler divergence
- WIM appears in Otto calculus and Wasserstein gradient flow

## Example: 1d Gaussian

$$p(x | \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right), \quad \theta = (\mu, \sigma)$$

- Wasserstein distance

$$W_2(p(x | \theta_1), p(x | \theta_2))^2 = (\mu_1 - \mu_2)^2 + (\sigma_1 - \sigma_2)^2$$

- Wasserstein score function

$$\Phi_\mu^W(x | \theta) = x - \mu, \quad \Phi_\sigma^W(x | \theta) = \frac{(x - \mu)^2 - \sigma^2}{2\sigma}$$

- Wasserstein information matrix

$$G_W(\theta) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

- More generally, 1d location-scale model is Euclidean (totally geodesic) in  $L^2$ -Wasserstein geometry

# Wasserstein estimator

## Definition (Li and Zhao, 2023)

The **Wasserstein estimator**  $\hat{\theta}_W(x)$  is the zero of the Wasserstein score function:

$$\Phi_i^W(x | \hat{\theta}_W(x)) = 0, \quad i = 1, \dots, p$$

- cf. MLE = zero of the Fisher score function = projection w.r.t. Kullback–Leibler divergence
- What does it mean??
  - ▶ It is different from the projection w.r.t. Wasserstein distance studied in Amari and M. (2022)

## Example: elliptically contoured family

$$p(x | \mu, \Sigma) = (\det \Sigma)^{-1/2} f(\|\Sigma^{-1/2}(x - \mu)\|)$$

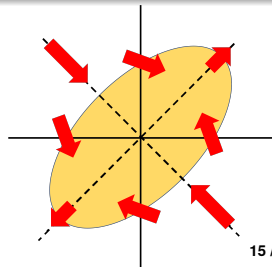
### Theorem (Amari and M., 2024)

- Wasserstein score functions are quadratic
- Wasserstein estimator = 2nd-order moment estimator

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top$$

e.g. 2d Gaussian  $N_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \theta \\ \theta & 1 \end{pmatrix} \right)$

$$\Phi^W(x | \theta) = \frac{1}{4} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}^\top \begin{pmatrix} -\theta & 1 \\ 1 & -\theta \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$



# Wasserstein covariance & Wasserstein–Cramer–Rao

## Definition (Li and Zhao, 2023)

The **Wasserstein covariance**  $\text{Var}_\theta^{\text{W}}[\hat{\theta}]$  of an estimator  $\hat{\theta}$  is the  $p \times p$  positive semidefinite matrix given by

$$\text{Var}_\theta^{\text{W}}[\hat{\theta}] = (\mathbf{E}_\theta[(\nabla_x \hat{\theta}_i)^\top (\nabla_x \hat{\theta}_j)])_{ij}$$

## Theorem (Li and Zhao, 2023)

When  $\hat{\theta}$  is unbiased ( $\mathbf{E}_\theta[\hat{\theta}] = \theta$ ),

$$\text{Var}_\theta^{\text{W}}(\hat{\theta}) \succeq G_{\text{W}}(\theta)^{-1}$$

- What does it mean??
  - cf. usual Cramer–Rao = lower bound of mean squared error



# Wasserstein covariance and robustness

$$X \sim p(x \mid \theta), \quad Z \sim q(z)$$

- We consider estimation of  $\theta$  from noisy observation  $X + Z$ 
  - $E[Z] = 0$ ,  $\text{Var}[Z] = \sigma^2 I$

## Theorem (Amari and M., 2024)

$$\begin{aligned} \text{Var}_\theta^W[\hat{\theta}] &= \lim_{\sigma^2 \rightarrow 0} \frac{\text{Var}_\theta[\hat{\theta}(X + Z)] - \text{Var}_\theta[\hat{\theta}(X)]}{\sigma^2} \\ &\quad - \frac{1}{2} \left( \text{Cov}_\theta[\hat{\theta}_a(X), \Delta\hat{\theta}_b(X)] + \text{Cov}_\theta[\hat{\theta}_b(X), \Delta\hat{\theta}_a(X)] \right) \end{aligned}$$

# Wasserstein covariance and robustness

## Corollary (Amari and M., 2024)

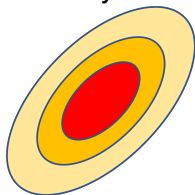
If  $\hat{\theta}$  is quadratic,

$$\text{Var}_{\theta}^{\text{W}}[\hat{\theta}] = \lim_{\sigma^2 \rightarrow 0} \frac{\text{Var}_{\theta}[\hat{\theta}(X + Z)] - \text{Var}_{\theta}[\hat{\theta}(X)]}{\sigma^2}$$

- Thus, Wasserstein covariance quantifies the **robustness against additive noise** of quadratic estimators.
- e.g. Wasserstein estimator for elliptically contoured family

$$p(x \mid \mu, \Sigma) = (\det \Sigma)^{-1/2} f(\|\Sigma^{-1/2}(x - \mu)\|)$$

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^{\top}$$



- “additive noise”: not invariant w.r.t. transformation of  $x$ 
  - noise contamination  $\approx$  (random) transportation