

# The Monge Gap: A Regularizer to Learn All Transport Maps

Theo Uscidda   Marco Cuturi  
**presenter:** Shen Yuan



中國人民大學  
RENMIN UNIVERSITY OF CHINA

高瓴人工智能学院  
Gaoling School of Artificial Intelligence

# Outline

Background

The Monge Gap

Learning with the Monge Gap

Experiments

Conclusion

# Outline

## Background

- Monge and Kantorovich formulation
- Entropic regularization

## The Monge Gap

## Learning with the Monge Gap

## Experiments

## Conclusion

# Monge formulation

Given a compact subset  $\Omega \subset \mathbb{R}^d$ , a continuous cost function  $c : \Omega \times \Omega \rightarrow \mathbb{R}$  and two probability distributions  $\mu, \nu \in \mathcal{P}(\Omega)$ , the Monge problem is to find  $T : \Omega \rightarrow \Omega$  that push-forward  $\mu$  onto  $\nu$ , which minimizes the averaged displacement cost:

$$W_c(\mu, \nu) := \inf_{T \# \mu = \nu} \int_{\Omega} c(\mathbf{x}, T(\mathbf{x})) d\mu(\mathbf{x}) \quad (1)$$

$c$ -OT means any solution to 1 between  $\mu$  and  $\nu$ .

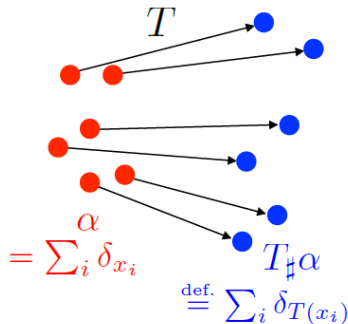


Figure 1: Push-forward of measures.

# Kantorovich formulation

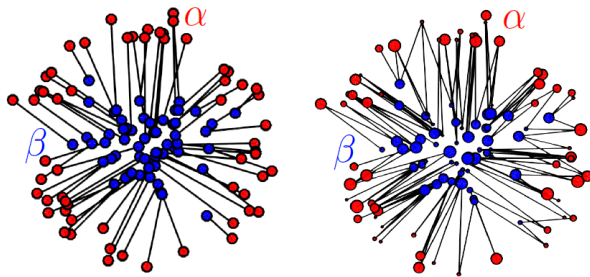


Figure 2: Comparison of transport maps and generic couplings.

# Kantorovich formulation

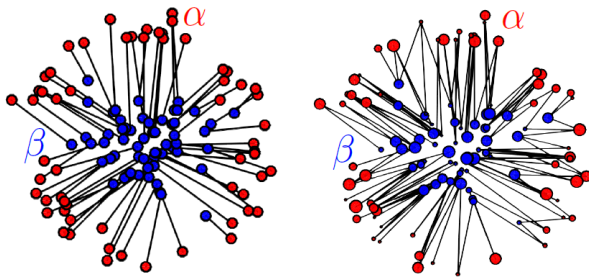


Figure 2: Comparison of transport maps and generic couplings.

$$W_c(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \int \int_{\Omega \times \Omega} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}) \quad (2)$$

An optimal coupling  $\pi^* = (\text{Id}, T^*)\# \mu$  always exists.

# Entropic regularization

For empirical measures  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$ ,  $\hat{\nu}_n = \frac{1}{n} \sum_{j=1}^n \delta_{\mathbf{y}_j}$  and  $\varepsilon > 0$ ,  $\mathbf{C} = [c(\mathbf{x}_i, \mathbf{y}_j)]_{ij}$ ,

$$W_{c,\varepsilon}(\hat{\mu}_n, \hat{\nu}_n) := \min_{\mathbf{P} \in U_n} \langle \mathbf{P}, \mathbf{C} \rangle - \varepsilon H(\mathbf{P}) \quad (3)$$

where  $U_n = \{\mathbf{P} \in \mathbb{R}_+^{n \times n}, \mathbf{P}\mathbf{1}_n = \frac{1}{n}\mathbf{1}_n, \mathbf{P}^T\mathbf{1}_n = \frac{1}{n}\mathbf{1}_n\}$  and  $H(\mathbf{P}) = -\sum_{i,j=1}^n \mathbf{P}_{ij} \log(\mathbf{P}_{ij})$ .

# Outline

## Background

Monge and Kantorovich formulation

Entropic regularization

## The Monge Gap

Learning with the Monge Gap

Experiments

Conclusion



# The Monge Gap

Given a cost  $c$  and a reference measure  $\rho \in \mathcal{P}$ , the Monge gap of a vector field  $T : \Omega \rightarrow \Omega$  is defined as:

$$\mathcal{M}_\rho^c(T) := \int_\Omega c(\mathbf{x}, T(\mathbf{x})) d\rho(\mathbf{x}) - W_c(\rho, T\# \rho) \quad (4)$$

# The Monge Gap

Given a cost  $c$  and a reference measure  $\rho \in \mathcal{P}$ , the Monge gap of a vector field  $T : \Omega \rightarrow \Omega$  is defined as:

$$\mathcal{M}_\rho^c(T) := \int_{\Omega} c(\mathbf{x}, T(\mathbf{x})) d\rho(\mathbf{x}) - W_c(\rho, T\# \rho) \quad (4)$$

- ▶ For any vector field  $T$ ,  $\mathcal{M}_\rho^c(T) \geq 0$ .
- ▶  $T$  is a  $c$ -OT map between  $\rho$  and  $T\# \rho \Leftrightarrow \mathcal{M}_\rho^c(T) = 0$

# The Monge Gap

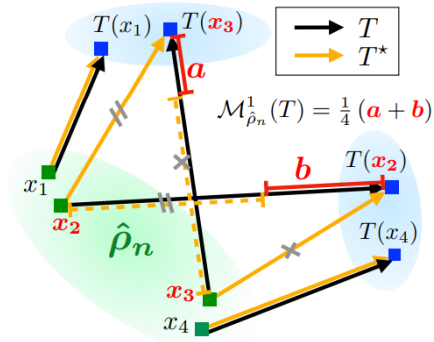


Figure 1. Sketch of the Monge Gap  $\mathcal{M}_{\hat{\rho}_n}^1(T)$  instantiated with the euclidean cost  $c(\cdot, \cdot) = \|\cdot - \cdot\|_2$ , where  $\hat{\rho}_n$  is a discrete measure supported on four points. Because the OT map  $T^*$  between  $\hat{\rho}_n$  and  $T\# \hat{\rho}_n$  does not coincide with  $T$  (notably on on points  $x_2, x_3$ ), the Monge gap here is positive, and equal to differences in lengths that amount to  $(a+b)/4$  in the plot.

# Consistency of the Monge Gap

**Lemma 3.2**(Consistency). Given empirical measures  $\hat{\rho}_n := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$ , provided that  $T$  is continuous, it almost surely holds

$$\lim_{n \rightarrow +\infty} \mathcal{M}_{\hat{\rho}_n}^c(T) = \mathcal{M}_{\rho}^c(T) \quad (5)$$

# Consistency of the Monge Gap

**Lemma 3.2**(Consistency). Given empirical measures  $\hat{\rho}_n := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$ , provided that  $T$  is continuous, it almost surely holds

$$\lim_{n \rightarrow +\infty} \mathcal{M}_{\hat{\rho}_n}^c(T) = \mathcal{M}_{\rho}^c(T) \quad (5)$$

$$\begin{aligned} \mathcal{M}_{\rho}^c(T) &:= \int_{\Omega} c(\mathbf{x}, T(\mathbf{x})) d\rho(\mathbf{x}) - W_c(\rho, T\# \rho) \\ \mathcal{M}_{\hat{\rho}_n, \varepsilon}^c(T) &:= \frac{1}{n} \sum_{i=1}^n c(\mathbf{x}_i, T(\mathbf{x}_i)) - W_{c, \varepsilon}(\hat{\rho}_n, T\# \hat{\rho}_n) \end{aligned} \quad (6)$$

## Relation to Cyclical Monotonicity

For any  $n \in \mathbb{N}$ , any set  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \times \{\mathbf{y}_1, \dots, \mathbf{y}_n\} \subset \Gamma$  and permutation  $\sigma \in \mathcal{S}_n$ , a set  $\Gamma \subset \Omega \times \Omega$  is  $c$ -CM(Cyclical Monotonicity) if

$$\sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_i) \leq \sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_{\sigma(i)}) \quad (7)$$

## Relation to Cyclical Monotonicity

For any  $n \in \mathbb{N}$ , any set  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \times \{\mathbf{y}_1, \dots, \mathbf{y}_n\} \subset \Gamma$  and permutation  $\sigma \in \mathcal{S}_n$ , a set  $\Gamma \subset \Omega \times \Omega$  is  $c$ -CM(Cyclical Monotonicity) if

$$\sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_i) \leq \sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_{\sigma(i)}) \quad (7)$$

With  $\mathbf{y}_i := T(\mathbf{x}_i)$ , the Monge gap estimator using permutations is:

$$\mathcal{M}_{\hat{\rho}_n}^c(T) = \frac{1}{n} \sum_{i=1}^n c(\mathbf{x}_i, T(\mathbf{x}_i)) - \min_{\sigma \in \mathcal{S}_n} \frac{1}{n} \sum_{i=1}^n c(\mathbf{x}_i, T(\mathbf{x}_{\sigma(i)})) \quad (8)$$

The cyclical monotonicity of that set is equivalent to the optimality of  $T$ .

# Properties of the Monge Gap

**Proposition 3.3.** Let  $\mu, \nu \in \mathcal{P}(\Omega)$  such that  $\text{Spt}(\mu) \subset \text{Spt}(\rho)$ , and a map  $T$  s.t.  $T\# \mu = \nu$ . Then  $\mathcal{M}_\rho^c(T) = 0$  implies that  $T$  is a c-OT map between  $\mu$  and  $\nu$ .



# Outline

## Background

- Monge and Kantorovich formulation
- Entropic regularization

## The Monge Gap

## Learning with the Monge Gap

## Experiments

## Conclusion

## Using directly the Monge gap as a regularizer

Given a loss function defined through a divergence  $\Delta$ , a regularization weight  $\lambda_{\text{MG}} > 0$  is introduced:

$$\min_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta) := \Delta(T_\theta \# \mu, \nu) + \lambda_{\text{MG}} \mathcal{M}_\rho^c(T_\theta) \quad (9)$$

## Gradient of Monge Gap

According to the Danskin (1967) Theorems,  $\mathcal{M}_{\hat{\rho}_n, \varepsilon}^c$  is differentiable and its gradient reads:

$$\nabla_{\theta} \mathcal{M}_{\hat{\rho}_n, \varepsilon}^c(T_{\theta}) = \sum_{i,j=1}^n \left( \frac{1}{n} \delta_{ij} - \mathbf{P}_{ij}^{\varepsilon} \right) \nabla_{\theta} c(\mathbf{x}_i, T_{\theta}(\mathbf{x}_j)) \quad (10)$$

# Gradient of Monge Gap

According to the Danskin (1967) Theorems,  $\mathcal{M}_{\hat{\rho}_n, \varepsilon}^c$  is differentiable and its gradient reads:

$$\nabla_{\theta} \mathcal{M}_{\hat{\rho}_n, \varepsilon}^c(T_{\theta}) = \sum_{i,j=1}^n \left( \frac{1}{n} \delta_{ij} - \mathbf{P}_{ij}^{\varepsilon} \right) \nabla_{\theta} c(\mathbf{x}_i, T_{\theta}(\mathbf{x}_j)) \quad (10)$$

Since  $\mathbf{P}^{\varepsilon} \in U_n$ ,  $\forall i, j$ ,  $0 \leq \mathbf{P}_{ij}^{\varepsilon} \leq 1/n$ , so:

$$\begin{cases} (1/n) \delta_{ij} - \mathbf{P}_{ij}^{\varepsilon} \geq 0 & \text{if } i = j \\ (1/n) \delta_{ij} - \mathbf{P}_{ij}^{\varepsilon} \leq 0 & \text{if } i \neq j \end{cases} \quad (11)$$

# Handling Costs with Structure

For cost  $c(\mathbf{x}, \mathbf{y}) = h(\mathbf{x} - \mathbf{y})$  with  $h$  strictly convex, the map has structure, as a known functional depending on  $h^*$  applied to the gradient a dual potential. A parametrized vector field  $F_\theta$  is introduced to model directly the dual potential gradient  $\nabla\psi^*$ :

$$T_\theta : \mathbf{x} \mapsto \mathbf{x} - \nabla h^* \circ F_\theta(\mathbf{x}) \quad (12)$$

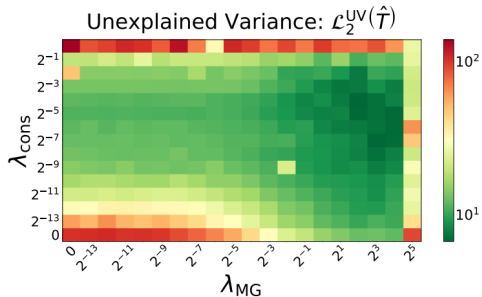
# The final loss function

The regularizer penalizes the asymmetry of  $\text{Jax}_{\mathbf{x}}F$  for  $\mathbf{x} \sim \rho$ :

$$\mathcal{C}_{\rho}(F) = \mathbb{E}_{X \sim \rho} \left[ \|\text{Jac}_X F - \text{Jac}_X^T F\|_2^2 \right] \quad (13)$$

$$\min_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta) := \Delta((I_d - \nabla h^* \circ F_{\theta}) \# \mu, \nu) + \lambda_{\text{MG}} \mathcal{M}_{\rho}^c(I_d - \nabla h^* \circ F_{\theta}) + \lambda_{\text{cons}} \mathcal{C}_{\rho}(F_{\theta}) \quad (14)$$

# The selection of the regularization weights ( $\lambda_{\text{MG}}$ , $\lambda_{\text{cons}}$ )



$$\mathcal{L}_2^{\text{UV}}(\hat{T}) := 100 \frac{\mathbb{E}_{\mu}[\|\hat{T}(X) - T^*(X)\|^2]}{\text{Var}_{\nu}(X)} \quad (15)$$

Figure 5. Heatmap showing the influence of the Monge gap  $\mathcal{M}_{\mu}^2$  and the conservative regularizer  $\mathcal{C}_{\mu}^2$ , when learning the Monge map for the  $\ell_2^2$  cost between Korotin et al. (2021) benchmark pair of dimension  $d = 32$ . For each pair of regularization weights  $(\lambda_{\text{MG}}, \lambda_{\text{cons}})$  on a regular grid, we report the unexplained variance  $\mathcal{L}_2^{\text{UV}}(\hat{T})$  provided by the the estimated map  $\hat{T}$ .

# Outline

## Background

- Monge and Kantorovich formulation
- Entropic regularization

## The Monge Gap

## Learning with the Monge Gap

## Experiments

## Conclusion



# Experiments

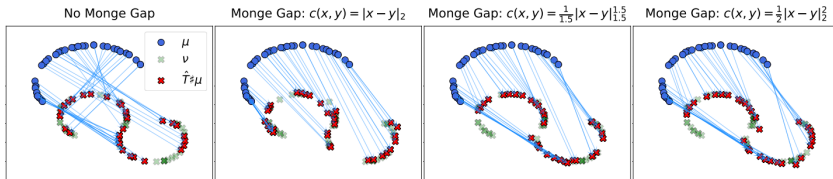


Figure 2. Fitting of transport maps between synthetic measures  $\mu, \nu$  in dimension  $d = 2$ , with the same fitting loss  $\Delta = W_{2,\varepsilon}$  but Monge gap  $\mathcal{M}_\mu^c$  instantiated with various costs  $c$ . We also fit an MLP without Monge gap, minimizing only the fitting loss. For  $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$ , we use the method for generic costs §4.1, directly parameterizing  $T_\theta$  as an MLP and using  $\lambda_{\text{MG}} = 5$ . For  $c(\mathbf{x}, \mathbf{y}) = \frac{1}{1.5}\|\mathbf{x} - \mathbf{y}\|_{1.5}^{1.5}$  and  $c(\mathbf{x}, \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$ , since they have the form  $c(\mathbf{x}, \mathbf{y}) = h(\mathbf{x} - \mathbf{y})$  with  $h$  strictly convex and known Legendre transform  $h^*$ , we use the method for costs with structure §4.2. Accordingly, we parameterize  $T_\theta = \text{Id} - \nabla h^* \circ F_\theta$  with an MLP  $F_\theta$  and penalize lack of conservativity with  $C_\mu$ . Moreover, we use  $\lambda_{\text{MG}} = 1$  and  $\lambda_{\text{cons}} = 0.01$ .

# Experiments

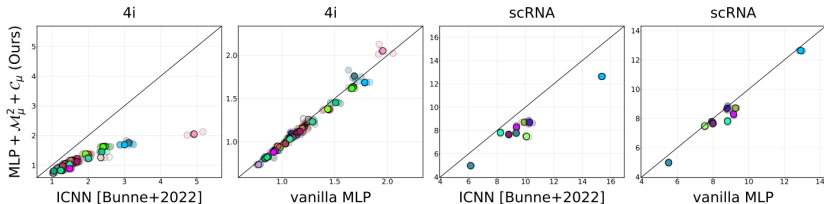


Figure 3. Fitting of a transport map  $\hat{T}$  to predict the responses of cells populations to cancer treatments, on 4i and scRNA datasets, providing respectively 34 and 9 treatment responses. For each profiling technology and each treatment, we compare the predictions of a MLP trained with Monge gap  $\mathcal{M}_\mu^2(F)$  + conservative regularizer  $C_\mu$  to those provided by a vanilla MLP (trained without regularization), and a gradient-ICNN learned via the neural dual formulation (Makkuva et al., 2020). We measure predictive performance using the Sinkhorn divergence between a batch of unseen (test) treated cells and a batch of unseen control cells mapped with  $\hat{T}$ , see § 6.4 and Appendix B.5 for details. Each scatter plot displays points  $z_i = (x_i, y_i)$  where  $y_i$  is the divergence obtained by our method and  $x_i$  that of the other baseline, on all treatments. A point below the diagonal  $y = x$  refers to an experiment in which our methods outperforms the baseline. To each treatment, we assign a color and plot 5 runs, along with their mean (the brighter point).

# Experiments

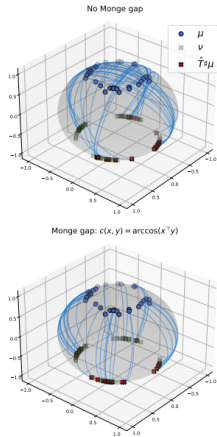


Figure 4. Fitting of transport maps between synthetic measures on the 2-sphere. In both cases, we parameterize the map as  $T_{\theta} = F_{\theta} / \|F_{\theta}\|_2$  where  $F_{\theta}$  is an MLP, and we use  $\Delta = W_{\ell_2, \ell_2}$  as fitting loss. On the upper plot, we do not use any regularizer while on the lower plot we regularize with the Monge gap instantiated for the geodesic cost  $c(\mathbf{x}, \mathbf{y}) = \arccos(\mathbf{x}^T \mathbf{y})$  and use  $\lambda_{MG} = 1$ .

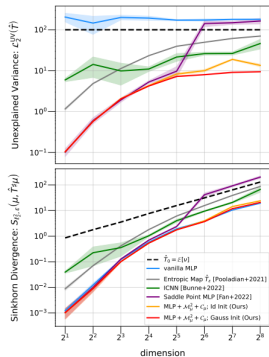


Figure 6. Performances of Monge gap-based learning and baselines on estimating the ground-truth maps between each pair of Gaussian mixtures  $\mu, \nu$  in dimension  $d \in \{2, 4, 8, \dots, 256\}$  of the Korotin et al. (2021) benchmark. We report both Sinkhorn divergence  $S_{\ell_2, \ell_2}(\tilde{T}_{\theta}^{\mu}, \nu)$  and the unexplained variance  $L_2^{UV}(\tilde{T})$  averaged over 5 fittings.

# Outline

## Background

- Monge and Kantorovich formulation
- Entropic regularization

## The Monge Gap

## Learning with the Monge Gap

## Experiments

## Conclusion

# Conclusion

- ▶ This paper provides a new strategy to train optimal transport maps.
- ▶ The regularizer adapts to any cost  $c$ , but requires defining a reference measure  $\rho$ .