

# Self-study: Variational Wasserstein Problems

Tianpei Xie

Aug. 19th., 2022

## Contents

<b>1</b>	<b>Previous lessons</b>	<b>2</b>
<b>2</b>	<b>Differentiating the Wasserstein Loss</b>	<b>5</b>
2.1	Eulerian and Lagrangian discretization . . . . .	6
2.2	Automatic differentiation . . . . .	7
<b>3</b>	<b>Wasserstein Barycenters and Clustering</b>	<b>8</b>
3.1	Discrete measures . . . . .	8
3.2	Arbitrary measures . . . . .	8
3.3	K-means as a Wasserstein variational problem . . . . .	9
3.4	Infinite mixtures and non-parametric Bayesian prior . . . . .	9
3.5	Special cases . . . . .	9
3.5.1	Barycenter of Gaussian measures is Gaussian . . . . .	9
3.5.2	1-dimensional distributions on $\mathbb{R}$ . . . . .	10
3.5.3	Affine push-forward measures . . . . .	11
3.6	Entropic approximation . . . . .	11
3.6.1	Primal solution and generalized Sinkhorn algorithm . . . . .	12
3.6.2	Dual problem . . . . .	13
3.7	Wasserstein propagation . . . . .	13
<b>4</b>	<b>Minimum Kantorovich Estimators</b>	<b>13</b>
4.1	Challenges for Maximum Likelihood Estimation . . . . .	13
4.2	Learning generative models with Minimum Kantorovich Estimators . . . . .	14

# 1 Previous lessons

- Recall the definition of Wasserstein distance between two measures  $\alpha, \beta \in \mathcal{M}(\mathcal{X})$ , where  $\mathcal{M}(\mathcal{X})$  is the space of Radon measure on  $\mathcal{X}$ :

**Definition** We suppose  $\mathcal{X} = \mathcal{Y}$  and that for some  $p \geq 1$ ,  $c(x, y) = d(x, y)^p$ , where  $d$  is a distance on  $\mathcal{X}$ , i.e.

- $d(x, y) = d(y, x)$  for all  $x, y \in \mathcal{X}$ ;
- $d(x, y) = 0$  iff  $x = y$ ;
- $d(x, y) \leq d(x, z) + d(z, y)$ , for all  $x, y, z \in \mathcal{X}$

Then the  **$p$ -Wasserstein distance** between  $\alpha, \beta \in \mathcal{M}_+^1(\mathcal{X})$  on  $\mathcal{X}$  is defined by  $\mathcal{W}_p(\alpha, \beta) := \mathcal{L}_{d^p}(\alpha, \beta)^{\frac{1}{p}}$ , where  $\mathcal{X} = \mathcal{Y}$  and

$$\begin{aligned} \mathcal{L}_{d^p}(\alpha, \beta) &:= \min_{\pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})} \int_{\mathcal{X} \times \mathcal{Y}} d^p(x, y) d\pi(x, y) \\ &\text{s.t. } P_{\mathcal{X}\#}\pi = \alpha, \\ &\quad P_{\mathcal{Y}\#}\pi = \beta. \end{aligned}$$

Here  $P_{\mathcal{X}\#}$  and  $P_{\mathcal{Y}\#}$  are the **push-forwards** of the **projections**  $P_{\mathcal{X}}(x, y) = x$  and  $P_{\mathcal{Y}}(x, y) = y$ .

- For two random variables  $X, Y$ , the  **$p$ -Wasserstein distance between their distributions**  $\mathcal{W}_p(P_X, P_Y) := (\min_{\pi} \mathbb{E}_{(X, Y) \sim \pi} [d^p(X, Y)])^{\frac{1}{p}}$ .
- For discrete measures,  $\alpha = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i}$  and  $\beta := \sum_{i=1}^m b_i \delta_{\mathbf{y}_i}$ , the primal problem for optimal transport is

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times m}} \langle \mathbf{P}, \mathbf{C} \rangle = \sum_{i,j} C_{i,j} P_{i,j} \quad (1)$$

$$\text{s.t. } \mathbf{P} \mathbf{1}_m = \mathbf{a} \quad (2)$$

$$\mathbf{P}^T \mathbf{1}_n = \mathbf{b} \quad (3)$$

$$P_{i,j} \geq 0$$

where  $\mathbf{C}_{n,m} := [C_{i,j}]_{i \in [1:n], j \in [1:m]}$ ,  $C_{i,j} := c(\mathbf{x}_i, \mathbf{y}_j) \geq 0$ . The feasible set is defined as

$$U(\mathbf{a}, \mathbf{b}) := \{\mathbf{P} \in \mathbb{R}_+^{n \times m} : \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^T \mathbf{1}_n = \mathbf{b}\} \quad (4)$$

- the corresponding dual problem with respect to primal problem is

$$\max_{\boldsymbol{\lambda} \in \mathbb{R}^n, \boldsymbol{\mu} \in \mathbb{R}^m} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle \quad (5)$$

$$\text{s.t. } \lambda_i + \mu_j \leq C_{i,j} \quad \forall i \in [1:n], j \in [1:m] \quad (6)$$

where  $\boldsymbol{\lambda} = [\lambda_i]_n$ ,  $\boldsymbol{\mu} = [\mu_j]_m$  are **dual variables** (slack variables) for marginal distribution constrain  $\mathbf{a}$  and  $\mathbf{b}$ . We denote  $\boldsymbol{\lambda} \oplus \boldsymbol{\mu} := \boldsymbol{\lambda} \mathbf{1}_m^T + \mathbf{1}_n \boldsymbol{\mu}^T \in \mathbb{R}^{n \times m}$  so that the linear constraints is  $\boldsymbol{\lambda} \oplus \boldsymbol{\mu} \leq \mathbf{C}$ . Such dual variables  $\boldsymbol{\lambda}, \boldsymbol{\mu}$  are often referred to as "**Kantorovich potentials**." The feasible set of the dual problem is defined as

$$R(\mathbf{C}) := \{\boldsymbol{\lambda} \in \mathbb{R}^n, \boldsymbol{\mu} \in \mathbb{R}^m : \boldsymbol{\lambda} \oplus \boldsymbol{\mu} \leq \mathbf{C}\} \quad (7)$$

where  $\lambda \oplus \mu = \lambda \mathbf{1}_m + \mathbf{1}_n \mu^T$ .

- By strong duality:

$$\begin{aligned} \mathcal{W}_p(\alpha, \beta)^p &= \mathcal{L}_{d^p}(\alpha, \beta) := L_C(\mathbf{a}, \mathbf{b}) \\ &= \min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \\ &= \max_{\lambda \oplus \mu \leq \mathbf{C}} \langle \lambda, \mathbf{a} \rangle + \langle \mu, \mathbf{b} \rangle \end{aligned} \quad (8)$$

or at optimal solution

$$\begin{aligned} \mathcal{W}_p(\alpha, \beta)^p &= \langle \mathbf{P}^*, \mathbf{C} \rangle \\ &= \langle \lambda^*, \mathbf{a} \rangle + \langle \mu^*, \mathbf{b} \rangle \end{aligned} \quad (9)$$

where  $\mathbf{P}^*$  is the optimal coupling matrix for primal problem and  $(\lambda^*, \mu^*)$  are optimal Kantorovich potentials or optimal dual solutions that satisfies

$$\lambda^* \oplus \mu^* = \mathbf{C} = [D_{i,j}^p] \quad (10)$$

That is both primal objective function and dual objective function can be used to estimate the Wasserstein distance.

- We also have the **maximum entropy optimal transport problem**

$$L_C^\epsilon(\mathbf{a}, \mathbf{b}) = \min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle - \epsilon H(\mathbf{P}) \quad (11)$$

where the second term is entropy

$$H(\mathbf{P}) := - \sum_{i,j} P_{i,j} (\log(P_{i,j}) - 1) \quad (12)$$

This problem has a unique optimal solution

$$\mathbf{P}^* = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \quad (13)$$

where  $\mathbf{u} = [\exp(\lambda_i/\epsilon)] = \exp(\lambda/\epsilon)$  and  $\mathbf{v} = [\exp(\mu_j/\epsilon)] = \exp(\mu/\epsilon)$ .

- The dual problem of the **maximum entropy optimal transport problem**:

$$L_C^\epsilon(\mathbf{a}, \mathbf{b}) = \max_{\lambda \in \mathbb{R}^n, \mu \in \mathbb{R}^m} \langle \lambda, \mathbf{a} \rangle + \langle \mu, \mathbf{b} \rangle - \epsilon \langle \exp(\lambda/\epsilon), \mathbf{K} \exp(\mu/\epsilon) \rangle \quad (14)$$

where  $\mathbf{K} = \exp(-\mathbf{C}/\epsilon)$  is the Gibbs distribution.

- The **probability interpretation** of original primal and dual Kantorovich optimal transport problem:

$$\begin{aligned} (P) \quad \mathcal{L}_c(\alpha, \beta) &= \min_{(X,Y) \sim \pi} \mathbb{E}_{(X,Y)} [c(X, Y)] \\ &\text{s.t. } X \sim \alpha, \\ &\quad Y \sim \beta \end{aligned} \quad (15)$$

$$\begin{aligned} (D) \quad \mathcal{L}_c(\alpha, \beta) &= \max_{(\lambda, \mu) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y})} \mathbb{E}_{X \sim \alpha} [\lambda(X)] + \mathbb{E}_{Y \sim \beta} [\mu(Y)] \\ &\text{s.t. } \lambda(x) + \mu(y) \leq c(x, y), \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}, \end{aligned} \quad (16)$$

- The **probability interpretation** of primal and dual maximum entropy optimal transport problem:

$$(P) \quad \mathcal{L}^\epsilon(\alpha, \beta) := \min_{(X,Y) \sim \pi} \mathbb{E}_{(X,Y)} [c(X, Y)] + \epsilon I(X; Y) \quad (17)$$

$$\text{s.t. } X \sim \alpha$$

$$Y \sim \beta$$

where  $I(X; Y) := \text{KL}(\pi \parallel \alpha \otimes \beta)$  is the mutual information between  $X$  and  $Y$ .

$$(D) \quad \mathcal{L}^\epsilon(\alpha, \beta) := \sup_{(\lambda, \mu) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y})} \mathbb{E}_{X \sim \alpha} [\lambda(X)] + \mathbb{E}_{Y \sim \beta} [\mu(Y)]$$

$$- \epsilon \mathbb{E}_{X \sim \alpha, Y \sim \beta} \left[ \exp \left( \frac{-c(X, Y) + \lambda(X) + \mu(Y)}{\epsilon} \right) \right]. \quad (18)$$

- Given a cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , the ***c-transform*** of  $f$  is defined as

$$f^c(y) := \inf_{x \in \mathcal{X}} c(x, y) - f(x) \quad (19)$$

The function  $f^c : \mathcal{Y} \rightarrow \mathbb{R}$  is also called the ***c-conjugate function*** of  $f$ . Similarly,  $g : \mathcal{Y} \rightarrow \mathbb{R}$ , the  ***$\bar{c}$ -transform*** of  $g$  is defined as

$$g^{\bar{c}}(x) := \inf_{y \in \mathcal{Y}} c(x, y) - g(y) \quad (20)$$

A function  $\psi : \mathcal{X} \rightarrow \mathbb{R}$  is ***c-concave*** if there exists some function  $\phi : \mathcal{Y} \rightarrow \mathbb{R}$  and cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  so that  $\psi$  is the  $\bar{c}$ -transform of  $\phi$ , i.e.  $\psi = \phi^{\bar{c}}$ . Denote  $\psi$  as  $c\text{-concave}(\mathcal{X})$ .

A function  $\phi : \mathcal{Y} \rightarrow \mathbb{R}$  is  ***$\bar{c}$ -concave*** if there exists some function  $\psi : \mathcal{X} \rightarrow \mathbb{R}$  and cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  so that  $\phi$  is the  $c$ -transform of  $\psi$ , i.e.  $\phi = \psi^c$ . Denote  $\phi$  as  $\bar{c}\text{-concave}(\mathcal{Y})$ .

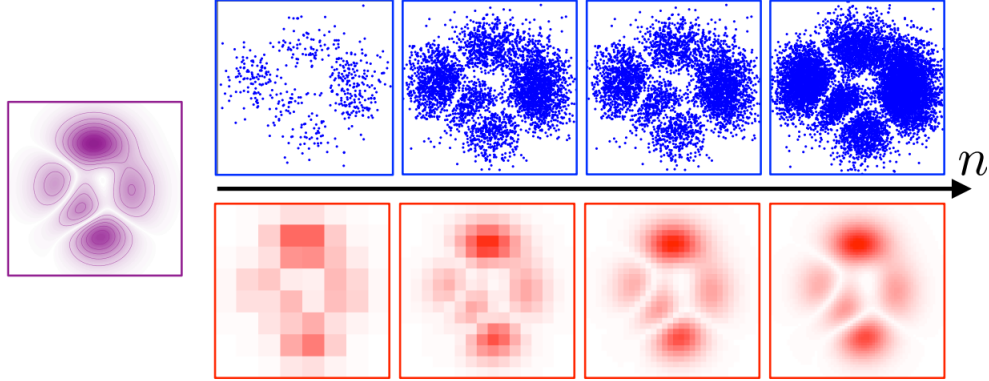
For distance  $c = d$ ,  $f^c = f^{\bar{c}}$ , thus we drop their distinctions.

**Proposition 1.1** *If  $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a distance, then the function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is ***c-concave*** if and only if  $f$  is ***Lipschitz continuous*** with Lipschitz constant less than 1 w.r.t. the distance  $c$ . We will denote by  $Lip_1$  the set of these functions. Moreover, for every  $f \in Lip_1$ , i.e.  $\|f\|_L \leq 1$ , we have the  $c$ -transform of  $f$ ,  $f^c = -f$ . [Santambrogio, 2015]*

- Thus the dual problem (5) is equivalent to an **unconstrained optimization problem**

$$\mathcal{L}_c(\alpha, \beta) := \max_{\lambda \in \mathcal{C}(\mathcal{X})} \int_{\mathcal{X}} \lambda d\alpha + \int_{\mathcal{Y}} \lambda^c d\beta \quad (21)$$

$$= \max_{\mu \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} \mu^{\bar{c}} d\alpha + \int_{\mathcal{Y}} \mu d\beta \quad (22)$$



**Figure 9.1:** Increasing fine discretization of a continuous distribution having a density (violet, left) using a Lagrangian representation  $\frac{1}{n} \sum_i \delta_{x_i}$  (blue, top) and an Eulerian representation  $\sum_i \mathbf{a}_i \delta_{x_i}$  with  $x_i$  representing cells on a grid of increasing size (red, bottom). The Eulerian perspective starts from a pixelated image down to one with such fine resolution that it almost matches the original density. Weights  $\mathbf{a}_i$  are directly proportional to each pixel-cell’s intensity.

**Figure 1:** The comparison between Eulerian Discretization (down) and Lagrangian Discretization (top) [Peyr and Cuturi, 2019]

## 2 Differentiating the Wasserstein Loss

The optimal transport geometry has a unique ability, not shared with other information divergences, to leverage physical ideas (mass displacement) and geometry (a ground cost between observations or bins) to compare measures. These two facts combined make it thus very tempting to use the Wasserstein distance as a loss function.

However, the main technical challenge associated with that idea lies in approximating and **differentiating** efficiently the Wasserstein distance. We start this section by presenting methods to approximate such gradients, then follow with three important applications that can be cast as variational Wasserstein problems.

We consider a **statistical estimation problem**: given a probability distribution  $\beta$  arising from measurements, we need to find a model from a parameterized family of distributions  $\{\alpha_\theta, \theta \in \Theta\}$  that **minimizes** the *Wasserstein distance*, where  $\Theta$  is a subset of a Euclidean space.

$$\min_{\theta \in \Theta} \mathcal{L}_c(\alpha_\theta, \beta) \quad (23)$$

where  $\mathcal{L}_c(\alpha_\theta, \beta)$  is defined as the optimal value of (15). This problem in general is **non-convex**. For cases such as Gaussian measures or elliptically contoured distributions, the **Wasserstein distance based objective** function has closed form, which can be solved directly.

In most cases, however, one has to resort to a careful **discretization** of  $\alpha_\theta$  to compute a **local minimizer** for Problem (23). Two approaches can be envisioned: Eulerian or Lagrangian. Note that the discrete measure is  $\alpha = \sum_i^n a_i \delta_{x_i}$ , which depends on two elements: the **position** of mass  $\mathbf{X} := \{x_1, \dots, x_n\}$  and the **probability measure** at each position  $\mathbf{a} := \{a_1, \dots, a_n\}$ . Figure 1 compares the Eulerian discretization and Lagrangian discretization.

## 2.1 Eulerian and Lagrangian discretization

- **Eulerian Discretization** (when the measure of mass changes but the position of mass does not)

A first way to discretize the problem is to suppose that both distributions  $\alpha = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i}$  and  $\beta = \sum_{i=1}^m b_i \delta_{\mathbf{y}_i}$  are discrete distributions defined on fixed locations  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$ . The parameterized measure  $\alpha_\theta$  is in that case entirely represented through the weight vector  $\mathbf{a} : \theta \rightarrow \mathbf{a}(\theta) \in \Delta_n$ , which, in practice, might be very sparse if the grid is large. This setting corresponds to the so-called class of **Eulerian discretization methods**.

In order to obtain differentiable objective function, we use the maximum entropy OT (11). The statistical estimation problem becomes

$$\min_{\theta \in \Theta} L_C^\epsilon(\mathbf{a}(\theta), \mathbf{b}) := \mathcal{E}_E(\theta) \quad (24)$$

where  $C_{i,j} = d(\mathbf{x}_i, \mathbf{y}_j)$ .

Recall that the maximum entropy OT objective is **differentiable** and convex w.r.t. input histograms, with gradient

$$\nabla_{(\mathbf{a}, \mathbf{b})} L_C^\epsilon(\mathbf{a}, \mathbf{b}) = \begin{bmatrix} \boldsymbol{\lambda}_* \\ \boldsymbol{\mu}_* \end{bmatrix} \quad (25)$$

where  $(\boldsymbol{\lambda}_*, \boldsymbol{\mu}_*)$  are the **dual potentials** i.e. the optimal solutions of regularized dual problem (14) chosen so that  $\sum_i \lambda_i = \sum_j \mu_j = 0$ . The zero mean condition on  $(\boldsymbol{\lambda}_*, \boldsymbol{\mu}_*)$  is important when using gradient descent to guarantee conservation of mass.

By chain rule, we have

$$\nabla_\theta \mathcal{E}_E(\theta) = \nabla_\theta \mathbf{a}(\theta)^T \boldsymbol{\lambda}_* \quad (26)$$

where  $\nabla_\theta \mathbf{a}(\theta) \in \mathbb{R}^{n \times d_\theta}$  is the Jacobian (differential) of the map  $\mathbf{a}(\theta)$  and  $d_\theta = \dim(\theta)$ .

The problem (24) is a convex optimization with differentiable gradient. It can be solved directly.

- **Lagrangian Discretization** (when the position of mass changes)

A different approach consists in using instead fixed (typically uniform) weights and approximating an input measure  $\alpha$  as an **empirical measure**  $\alpha_\theta = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i(\theta)}$  for a point-cloud parameterization map  $\mathbf{X} : \theta \rightarrow \mathbf{X}(\theta) = (\mathbf{x}(\theta)_i)$ , where  $\mathbf{x}(\theta)_i \in \mathcal{X}$  for  $i = 1, \dots, n$ . We assume here that  $\mathcal{X}$  is Euclidean.

Note that the measure changes when the **position** of mass changes, which will affect the cost function  $C_{i,j}$  as well. The optimization becomes

$$\min_{\theta \in \Theta} L_{C(\theta)}^\epsilon \left( \frac{1}{n} \mathbf{1}_n, \mathbf{b} \right) := \mathcal{E}_L(\theta) \quad (27)$$

where  $C(\theta) := [c(\mathbf{x}_i(\theta), \mathbf{y}_j)]_{i,j}$ . Note that here the cost matrix  $C(\theta)$  now depends on  $\theta$  since the support of  $\alpha_\theta$  changes with  $\theta$ .

The objective function  $L_C^\epsilon(\mathbf{a}, \mathbf{b})$  can be represented as  $L_C^\epsilon(\mathbf{a}, \mathbf{b}) = \langle \mathbf{P}^*, \mathbf{C} \rangle - \epsilon H(\mathbf{P}^*)$ . We can show that the function  $g : \mathbf{C} \rightarrow g(\mathbf{C}) := L_C^\epsilon(\mathbf{a}, \mathbf{b})$  is a **concave** and **smooth** function. The gradient of  $g$  is

$$\nabla_C L_C^\epsilon(\mathbf{a}, \mathbf{b}) = \mathbf{P}^* \quad (28)$$

where  $\mathbf{P}^* = \text{diag}(\mathbf{u}) \exp(-\mathbf{C}/\epsilon) \text{diag}(\mathbf{v})$  is the optimal primal solution (11) and it is also a function of  $\mathbf{C}$ .

By chain rule, we have

$$\nabla_\theta \mathcal{E}_L(\theta) = \nabla_\theta \mathbf{X}(\theta)^T (\nabla F(\mathbf{X}(\theta))) \quad (29)$$

where  $\mathbf{X}(\theta) := (\mathbf{x}(\theta)_i) \in \mathbb{R}^{nd \times 1}$  and  $\nabla_\theta \mathbf{X}(\theta) \in \mathbb{R}^{nd \times d_\theta}$  is the Jacobian matrix. Here we define  $F : \mathbf{X} \rightarrow L_{C(\mathbf{X})}^\epsilon(\frac{1}{n}\mathbf{1}_n, \mathbf{b})$ . Assuming  $(\mathcal{X}, \mathcal{Y})$  are convex subsets of  $\mathbb{R}^d$ , by chain rule we have

$$\begin{aligned} \nabla F(\mathbf{X}) &= \left[ \nabla_C L_C^\epsilon(\mathbf{a}, \mathbf{b})^T \nabla_{\mathbf{x}_i} C \right]_i \\ &= \left[ \sum_j^m P_{i,j,*} \nabla_{\mathbf{x}_i} c(\mathbf{x}_i, \mathbf{y}_j) \right]_{i=1}^n \in \mathcal{X}^n \end{aligned} \quad (30)$$

For instance, for  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ , for  $c(s, t) = \|s - t\|_2^2$ , one has

$$\nabla F(\mathbf{X}) = 2 \left[ \frac{1}{n} \mathbf{x}_i - \sum_j^m P_{i,j,*} \mathbf{y}_j \right]_{i=1}^n \in \mathcal{X}^n \subset \mathbb{R}^{nd} \quad (31)$$

Note that, up to a constant, this gradient is  $\text{Id} - T$ , where  $T$  is the **barycentric projection**.

## 2.2 Automatic differentiation

Note that for both Eulerian and Lagrangian discretization, the gradient in (26) and (29) requires knowledge on the exact optimal solution  $(\lambda_*, \mu_*)$  on dual problem and  $\mathbf{P}^*$  on primal problem. This can only be achieved with acceptable precision using a very large number of **Sinkhorn iterates**.

Given the limitation on time and computational resources, rather than approximating the gradient using the value obtained at a given iterate, it is usually better to *differentiate* directly the *output of Sinkhorn's algorithm*, using **reverse mode automatic differentiation**. In this case, we use the **Sinkhorn dual divergence**  $\mathfrak{D}_C^\epsilon(\mathbf{a}, \mathbf{b}) := \langle \lambda_*, \mathbf{a} \rangle + \langle \mu_*, \mathbf{b} \rangle$ , which incorporates the entropy of the regularized optimal transport, and differentiating it directly as a composition of simple maps using the inputs, either the histogram in the Eulerian case or the cost matrix in the Lagrangian cases.

Thus the objective for Eulerian and Lagrangian discretization becomes

$$\begin{aligned} (E) : & \quad \mathfrak{D}_C^\epsilon(\mathbf{a}(\theta), \mathbf{b}) \\ (L) : & \quad \mathfrak{D}_{C(\theta)}^\epsilon\left(\frac{1}{n}\mathbf{1}_n, \mathbf{b}\right) \end{aligned}$$

The cost for computing the **gradient** of functionals involving Sinkhorn dual divergences is *the same as* that of computation of the **functional** itself.

### 3 Wasserstein Barycenters and Clustering

#### 3.1 Discrete measures

Given input histogram  $\{\mathbf{b}_s\}_{s=1}^S$ , where  $\mathbf{b}_s \in \Delta_{n_s}$ , and weights  $\mathbf{w} \in \Delta_S$ , a Wasserstein barycenter is computed by minimizing

$$\min_{\mathbf{a} \in \Delta_n} \sum_{s=1}^S w_s L_{\mathbf{C}_s}(\mathbf{a}, \mathbf{b}_s) \quad (32)$$

where  $\mathbf{C}_s \in \mathbb{R}^{n \times n_s}$  is the cost matrix between  $\mathbf{a}$  and  $\mathbf{b}_s$ . Typically, we assume  $n_s = n$  and  $\mathbf{C}_s = \mathbf{C} = \mathbf{D}^p, \forall s$  as distance matrix, and the problem becomes

$$\min_{\mathbf{a} \in \Delta_n} \sum_{s=1}^S w_s \mathcal{W}_p^p(\mathbf{a}, \mathbf{b}_s)$$

For discrete measures, this barycenter problem (32) is in fact a *linear program*, since one can look for the  $S$  couplings  $(\mathbf{P}_s)_s$  between each input and the barycenter itself, which by construction must be constrained to share the same row marginal,

$$\begin{aligned} \min_{\substack{\mathbf{a} \in \Delta_n, \\ \{\mathbf{P}_s \in \mathbb{R}_+^{n \times n_s}, \forall s\}}} & \sum_{s=1}^S w_s \langle \mathbf{P}_s, \mathbf{C}_s \rangle \\ \text{s.t. } & \mathbf{P}_s \mathbf{1}_{n_s} = \mathbf{a}, \forall s = 1, \dots, S \\ & \mathbf{P}_s^T \mathbf{1}_n = \mathbf{b}_s, \forall s = 1, \dots, S \end{aligned}$$

This problem is a large-scale LP and can be solved using first order methods such as subgradient descent on the dual.

#### 3.2 Arbitrary measures

Given input measures  $\{\beta_s\}_{s=1}^S$ , where  $\beta_s \in \mathcal{M}_+(\mathcal{X})$ , and weights  $\mathbf{w} \in \Delta_S$ , a Wasserstein barycenter for arbitrary measures is

$$\min_{\alpha \in \mathcal{M}_+(\mathcal{X})} \sum_{s=1}^S w_s \mathcal{L}_c(\alpha, \beta_s) \quad (33)$$

In the case where  $\mathcal{X} = \mathbb{R}^d$  and  $c(x, y) = \|x - y\|_2^2$ , Agueh and Carlier show that if one of the input measures has a density, then this barycenter is **unique**. Moreover, under this condition, the **mean of the barycenter**  $\alpha_*$  is necessarily the **barycenter of the mean**, i.e.

$$\int_{\mathcal{X}} x d\alpha_* = \sum_{s=1}^S w_s \int_{\mathcal{X}} x d\beta_s$$

and the support of  $\alpha_*$  is located in the *convex hull* of the supports of the  $\{\beta_s\}_{s=1}^S$ .

Let us also note that it is possible to recast (33) as a **multimarginal OT** problem.



### 3.3 K-means as a Wasserstein variational problem

When  $\beta_s = \beta$  for all  $s$  and  $\mathcal{X} = \mathbb{R}^d$ ,  $c(x, y) = \|x - y\|_2^2$ , the barycenter problem becomes

$$\min_{\alpha \in \mathcal{M}_{k,1}(\mathcal{X})} \mathcal{L}_c(\alpha, \beta)$$

where  $\alpha = \sum_{i=1}^k a_i \delta_{\mathbf{c}_i}$  is constrained to be a **discrete measure** with a finite support of **size up to**  $k$ . This barycenter problem is equivalent to the usual **k-means** problem taking  $\beta$ . Indeed, one can easily show that the **centroids** output  $\{\mathbf{c}_i, i = 1, \dots, k\}$  by the k-means problem correspond to the **support** of the solution  $\alpha$  and that its **weights**  $a_i$  correspond to the **fraction** of points in  $\beta$  assigned to each centroid [Canas and Rosasco, 2012].

One can show that approximating  $\mathcal{L} \approx \mathcal{L}^\epsilon$  using entropic regularization results in smoothed out assignments that appear in **soft-clustering** variants of k-means, such as *mixtures of Gaussians*.

### 3.4 Infinite mixtures and non-parametric Bayesian prior

It is possible to generalize (33) to a possibly **infinite collection of measures**. This problem is described by considering a probability distribution  $M$  over the space  $\mathcal{M}_+^1(\mathcal{X})$  of probability distributions, i.e.  $M \in \mathcal{M}_+^1(\mathcal{M}_+^1(\mathcal{X}))$ . A **barycenter** is then a solution of

$$\min_{\alpha \in \mathcal{M}_+^1(\mathcal{X})} \mathbb{E}_{\beta \sim M} [\mathcal{L}_c(\alpha, \beta)] := \int_{\mathcal{M}_+^1(\mathcal{X})} \mathcal{L}_c(\alpha, \beta) dM(\beta). \quad (34)$$

where  $\beta$  is a *random measure* distributed according to  $M$ .  $M$  is a distribution over distributions, i.e. a **non-parametric Bayesian prior**. Then we can see that the problem (33) is a special case when  $M = \sum_s w_s \delta_{\beta_s}$  is a discrete measure with finite support on space of measures.

Drawing  $S$  measures uniformly randomly would results in  $\{\hat{\beta}_s\}_{s=1}^S$  with weight  $w_i = \frac{1}{S}$ . Define the unique solution of corresponding barycenter problem as  $\hat{\beta}_S = \arg \min_{\alpha} \sum_{s=1}^S \mathcal{L}_c(\alpha, \hat{\beta}_s)/S$ . It is shown that the barycenters  $\hat{\beta}_S$  has consistency, i.e.

$$\mathcal{L}_c(\hat{\beta}_S, \alpha_*) \xrightarrow{S \rightarrow \infty} 0$$

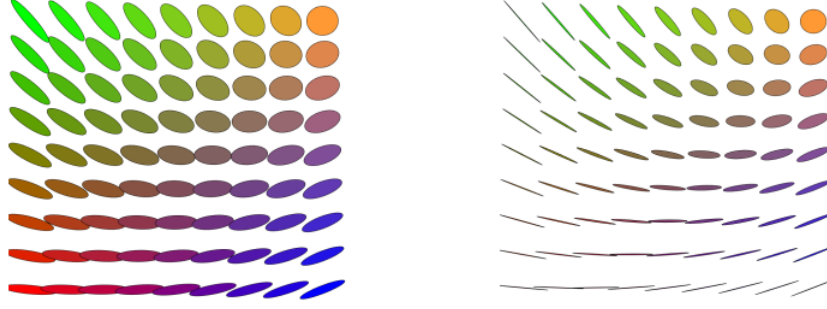
where  $\alpha_*$  is the optimal solution of infinite mixture barycenter problem (34). The convergence is in expectation or with high probability. This can be interpreted as a **law of large numbers** over the **Wasserstein space**.

### 3.5 Special cases

#### 3.5.1 Barycenter of Gaussian measures is Gaussian

When all the measures are Gaussian  $\beta_s = \mathcal{N}(\mathbf{m}_s, \Sigma_s)$  for all  $s$ , and  $c(x, y) = \|x - y\|_2^2$ , the **barycenter** (33) has a closed form solution and itself is a **Gaussian measure**  $\alpha_* = \mathcal{N}(\mathbf{m}_*, \Sigma_*)$ , where the **mean of barycenter** is the **barycenter of means** under Euclidean distance  $c$

$$\mathbf{m}_* = \sum_{s=1}^S w_s \mathbf{m}_s$$



**Figure 9.2:** Barycenters between four Gaussian distributions in 2-D. Each Gaussian is displayed using an ellipse aligned with the principal axes of the covariance, and with elongations proportional to the corresponding eigenvalues.

**Figure 2:** The smooth interpolation of Gaussian measures.

and the **covariance matrix of barycenter** is the **barycenter of covariance matrices** under *Bures metric*:

$$\Sigma_* = \arg \min_{\Sigma \in \mathcal{S}_+} \sum_{s=1}^S w_s \mathcal{B}(\Sigma, \Sigma_s)^2$$

This is a convex optimization problem and Bures metric is defined as

$$\mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 = \text{tr} \left( \Sigma_\alpha + \Sigma_\beta - 2 \left( \Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}} \right)^{\frac{1}{2}} \right).$$

$\Sigma_*$  is the fix-point solution of the following map

$$\begin{aligned} \Sigma_* &= \Phi(\Sigma_*) \\ \text{where } \Phi(\Sigma_*) &= \sum_{s=1}^S w_s \left( \Sigma_*^{\frac{1}{2}} \Sigma_s \Sigma_*^{\frac{1}{2}} \right)^{\frac{1}{2}} \end{aligned}$$

### 3.5.2 1-dimensional distributions on $\mathbb{R}$

For 1-D distributions, the  $\mathcal{W}_p$  barycenter can be computed almost in closed form using the fact that the transport is the monotone rearrangement. For empirical distribution  $\beta_s := \frac{1}{n} \sum_{i=1}^n \delta_{y_{i,s}}$ , where the points are assumed to be **ordered**  $y_{1,s} \leq y_{2,s} \leq \dots \leq y_{n,s}$ . The barycenter  $\alpha_w$  is also an empirical measure on  $n$  points

$$\alpha_w = \frac{1}{n} \sum_{i=1}^n \delta_{x_{i,w}}$$

where the mass location  $x_{i,w}$  is the result of barycenter map  $A_w$

$$x_{i,w} = A_w(y_{i,s}) := \arg \min_{x \in \mathbb{R}} \sum_{s=1}^S w_s \|x - y_{i,s}\|_p^p; \quad \forall i = 1, \dots, S.$$

For instance, for  $p = 2$ , one has  $x_{i,w} = \sum_{s=1}^S w_s y_{i,s}$  for all  $i$ .

In the general case, one needs to use the cumulative functions. Recall that the inverse of c.d.f, the generalized quantile function  $F_\alpha^{-1} : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty\}$  is defined as

$$F_\alpha^{-1}(r) := \min_x \{x \in \mathbb{R} \cup \{-\infty\} : F_\alpha(x) \geq r\}.$$

The  $\mathcal{W}_p$  distance is computed as

$$\mathcal{W}_p(\alpha, \beta)^p = \int_0^1 \left\| F_\alpha^{-1}(r) - F_\beta^{-1}(r) \right\|_p^p dr$$

So the **quantile** function of the *barycenter* is computed via the barycenter map

$$F_\alpha^{-1}(r) = A_w(F_{\beta_s}^{-1}(r))_{s=1}^S = \operatorname{argmin}_g \sum_s w_s \left\| g - F_{\beta_s}^{-1} \right\|_{L_p}^p$$

which can be used, for instance, to compute barycenters between discrete measures supported on less than  $n$  points in  $O(n \log(n))$  operations, using a simple sorting procedure.

### 3.5.3 Affine push-forward measures

Consider that all  $\beta_s = T_{r_s, u_s, \#} \alpha_0$  are push-forward measures from some base measure  $\alpha_0$ . The map  $T_{r,u} : x \rightarrow r x + u$  is **affine**, i.e. via scaling and translation. It can be shown that the barycenter  $\alpha_*$  of  $\{\beta_s\}$  is a **push-forward measure** from  $\alpha_0$

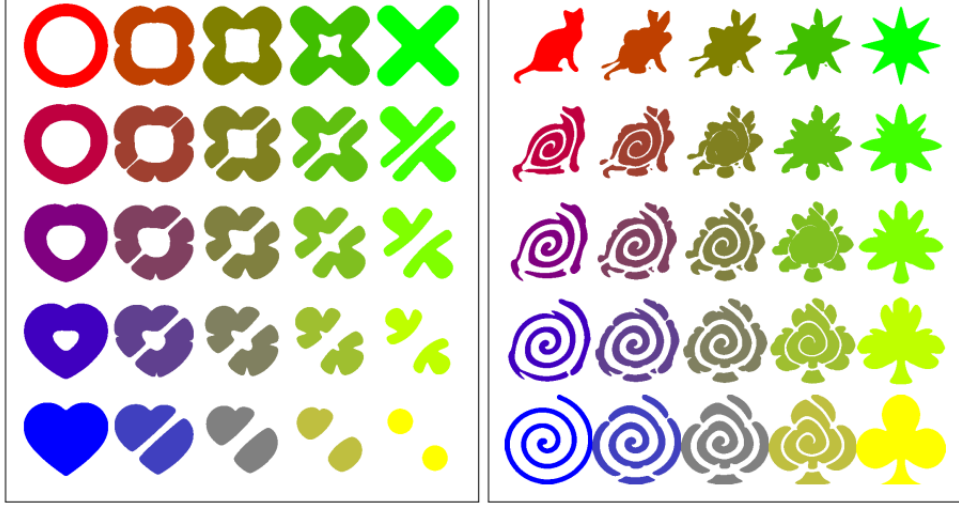
$$\begin{aligned} \alpha_* &= T_{r_*, u_*, \#} \alpha_0 \\ \text{where } r_* &= \left( \sum_s w_s / r_s \right)^{-1} \quad \text{i.e. the harmonic mean} \\ u_* &= \sum_s w_s u_s \quad \text{i.e. the arithmetic mean} \end{aligned}$$

## 3.6 Entropic approximation

One can use entropic smoothing and approximate the solution of (33) using

$$\min_{\alpha \in \mathcal{M}_+(\mathcal{X})} \sum_{s=1}^S w_s \mathcal{L}_c^\epsilon(\alpha, \beta_s) \tag{35}$$

for some  $\epsilon > 0$ . This is a **smooth convex** minimization problem, which can be tackled using gradient descent.



**Figure 9.3:** Barycenters between four input 2-D shapes using entropic regularization (9.15). To display a binary shape, the displayed images shows a thresholded density. The weights  $(\lambda_s)_s$  are bilinear with respect to the four corners of the square.

**Figure 3:** The barycenter shape obtained via entropic regularization [Peyr and Cuturi, 2019]

### 3.6.1 Primal solution and generalized Sinkhorn algorithm

The problem (35) can be reformulated using KL divergence

$$\begin{aligned} \min_{P_s \in \mathbb{R}_+^{n \times n_s}, \forall s} \sum_{s=1}^S w_s \text{KL}(P_s \| K_s) \\ \text{s.t. } P_s^T \mathbf{1}_n = \mathbf{b}_s, \quad s = 1, \dots, S \\ P_1 \mathbf{1}_{n_1} = \dots = P_S \mathbf{1}_{n_S} \end{aligned} \quad (36)$$

The optimal solution is

$$P_s = \text{diag}(\mathbf{u}_s) K_s \text{diag}(\mathbf{v}_s), \quad s = 1, \dots, S \quad (37)$$

and the scalings are sequentially updated via **generalized Sinkhorn algorithm**

$$\mathbf{v}_{s,t+1} \leftarrow \frac{\mathbf{b}}{K_s^T \mathbf{u}_{s,t}} := \mathbf{b} \oslash K_s^T \mathbf{u}_{s,t} \quad (38)$$

$$\mathbf{u}_{s,t+1} \leftarrow \frac{\mathbf{a}_{t+1}}{K_s \mathbf{v}_{s,t+1}} := \mathbf{a}_{t+1} \oslash K_s \mathbf{v}_{s,t+1} \quad (39)$$

where  $\mathbf{a}_{t+1} = \prod_s^S (K_s \mathbf{v}_{s,t+1})^{w_s}$  i.e. the geometric mean

### 3.6.2 Dual problem

The optimal  $(\mathbf{u}_{s,*}, \mathbf{v}_{s,*})$  appearing in (37) can be written as  $(\mathbf{u}_{s,*}, \mathbf{v}_{s,*}) = (\exp(\boldsymbol{\lambda}_{s,*}/\epsilon), \exp(\boldsymbol{\mu}_{s,*}/\epsilon))$ , where  $(\boldsymbol{\lambda}_{s,*}, \boldsymbol{\mu}_{s,*})$  are the solutions of the following dual program

$$\begin{aligned} \max_{(\boldsymbol{\lambda}_s, \boldsymbol{\mu}_s), \forall s} \quad & \sum_{s=1}^S w_s \{ \langle \boldsymbol{\mu}_s, \mathbf{b}_s \rangle - \epsilon \langle \exp(\boldsymbol{\lambda}_s/\epsilon), \mathbf{K}_s \exp(\boldsymbol{\mu}_s/\epsilon) \rangle \} \\ \text{s.t.} \quad & \sum_{s=1}^S w_s \boldsymbol{\lambda}_s = \mathbf{0} \end{aligned} \quad (40)$$

As with the original case, the generalized Sinkhorn algorithm (38) and (39) can be obtained by alternating minimization with respect to  $\boldsymbol{\mu}_s$  and  $\boldsymbol{\lambda}_s$ .

### 3.7 Wasserstein propagation

It is possible to generalize the barycenter problem (32), where one looks for distributions  $(\mathbf{b}_u)_{\mathcal{U}}$  at some given set  $\mathcal{U}$  of nodes in a **graph**  $\mathcal{G} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  given a set of fixed input distributions  $(\mathbf{b}_v)_{\mathcal{V}}$  on the *complementary* set  $\mathcal{V}$  of the nodes. The unknown are determined by minimizing the overall transportation distance between **all pairs of nodes**  $(u, v) \in \mathcal{E}$  forming edges in the graph

$$\min_{\{\mathbf{b}_u \in \Delta_{n_u}\}_{\mathcal{U}}} \sum_{(u,v) \in \mathcal{E}} L_{\mathbf{C}_{u,v}}(\mathbf{b}_u, \mathbf{b}_v)$$

where the cost matrices  $\mathbf{C}_{u,v} \in \mathbb{R}^{n_u \times n_v}$  need to be specified by the user.

The barycenter problem (32) is a special case of this problem where the considered graph  $\mathcal{G} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  is **star-shaped**, where  $\mathcal{U}$  is a single vertex connected to all the other vertices  $\mathcal{V}$  (the weight  $w_s$  associated to  $\mathbf{b}_s$  can be absorbed in the cost matrix). Introducing explicitly a coupling  $\mathbf{P}_{u,v} \in U(\mathbf{b}_u, \mathbf{b}_v)$  for each edge  $(u, v) \in \mathcal{E}$ , and using entropy regularization, one can rewrite this problem similarly as in (36), and one extends Sinkhorn iterations (38) and (39).

## 4 Minimum Kantorovich Estimators

### 4.1 Challenges for Maximum Likelihood Estimation

Given empirical measures  $\beta_n = \sum_{i=1}^n \frac{1}{n} \delta_{\mathbf{x}_i}$ ,  $\mathbf{x}_i \in \mathcal{X}$ , i.i.d generated from some unknown distribution  $\beta_*$ , the goal of model estimation is to fit a parametric model  $\theta \rightarrow \alpha_\theta \in \mathcal{M}(\mathcal{X})$  to the observed empirical input measure  $\beta_n$ . The standard **Maximum Likelihood Estimation (MLE)** minimize the *negative log-likelihood function*

$$\min_{\theta} \mathcal{L}_{MLE}(\alpha_\theta, \beta_n)$$

where  $\alpha_\theta$  has density  $\rho_\theta$  with respect to Lebesgue measure and  $\mathcal{L}_{MLE}(\alpha_\theta, \beta_n) = -\sum_i \log \rho_\theta(\mathbf{x}_i)$ . We know that the MLE convergences to KL divergence  $\mathcal{L}_{MLE}(\alpha, \beta_n) \rightarrow \mathbb{KL}(\alpha \| \beta_*)$  where  $\beta_*$  is the underlying distribution of data.

MLE has **challenges** to estimate  $\alpha_\theta$  that is *singular*, i.e. does not have density function  $\rho_\theta$  w.r.t. Lebesgue measure, or when the underlying distribution  $\beta_*$  is *singular*. In both cases, the KL divergence is not well defined since the support of  $\alpha_\theta$  and  $\beta$  may not overlap. Another issue is when  $\rho_\theta$  is complicated and hard to compute.

One typical case is when the task is to learn a generative model  $\alpha_\theta$  lies on a **low-dimensional sub-manifold**. Define the transform map  $h_\theta : \mathcal{Z} \rightarrow \mathcal{X}$  so that

$$\alpha_\theta = h_{\theta,\#}\zeta, \quad \zeta \in \mathcal{M}(\mathcal{Z})$$

where  $\mathcal{Z}$  is a low-dimensional subspace. Usually,  $\alpha_\theta$  does not have density w.r.t. Lebesgue measure on  $\mathcal{X}$  since its support does not cover the entire space  $\mathcal{X}$ . Furthermore, computing this density is usually intractable, while generating i.i.d. samples from  $\alpha_\theta$  is achieved by computing  $\mathbf{x}_i = h_\theta(\mathbf{z}_i)$ , where  $(\mathbf{z}_i)_i$  are i.i.d. samples from  $\zeta$ .

For instance, in low-rank covariance estimation,  $\mathbf{x} = \mathbf{A}\mathbf{z}$ , where  $\mathbf{A} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \Sigma_z)$ ,  $\Sigma_x = \mathbf{A}\Sigma_z\mathbf{A}^T$  has rank  $r \leq d$ . This is the degenerate case for covariance estimation.

## 4.2 Learning generative models with Minimum Kantorovich Estimators

We can use the Wasserstein distance as objective function in place of MLE objective. In particular, for Lipschitz function  $f$ , we can write  $\mathcal{W}_1$  in its dual form as

$$\mathcal{W}_1(\alpha, \beta) = \mathcal{L}_c(\alpha, \beta) := \sup_{f \in \text{Lip}_1} \left\{ \int_{\mathcal{X}} f(x) d\alpha(x) - \int_{\mathcal{X}} f(x) d\beta(x) \right\} \quad (41)$$

Note that since  $\alpha_\theta = h_{\theta,\#}\zeta$ , we have  $f(x) d\alpha_\theta(x) = f(x) d(h_{\theta,\#}\zeta)(x) = f(h_\theta(z)) d\zeta(z)$ . The estimation problem becomes

$$\min_{\theta} \mathcal{L}_c(\alpha_\theta, \beta_n) = \mathcal{L}_c(h_{\theta,\#}\zeta, \beta_n) \quad (42)$$

$$\Rightarrow \min_{\theta} \inf_{\pi \in U(\zeta, \beta_n)} \mathbb{E}_{\pi} [c(h_\theta(Z), Y)] \quad (\text{primal}) \quad (43)$$

$$\Rightarrow \min_{\theta} \sup_{f \in \text{Lip}_1} \{ \mathbb{E}_{Z \sim \zeta} [f(h_\theta(Z))] - \mathbb{E}_{X \sim \beta_n} [f(X)] \} \quad (\text{dual}) \quad (44)$$

The optimal transportation distances are weaker than MLE since they are not invariant to important families of invariances, such as rescaling, translation or rotations.

For  $\mathcal{E}(\theta) := \mathcal{L}_c(\alpha_\theta, \beta_n)$ , we can compute the subgradient function based on dual (41) when the **dual potential** function  $f$  is available,

$$\nabla_{\theta} \mathcal{E}(\theta) = \int_{\mathcal{Z}} [\nabla_{\theta} h_{\theta}(z)]^T \nabla f(h_{\theta}(z)) d\zeta(z) \quad (45)$$

Here  $\nabla f(x)$  is gradient w.r.t.  $x = h_{\theta}(z)$  and  $\nabla_{\theta} h_{\theta}(z)$  is the differential (with respect to  $\theta$ ) of  $h_{\theta}$ .

This formula is hard to use numerically, first because it requires first computing a continuous function  $f$ , which is a solution to a semi-discrete problem (16). For OT loss, this can be achieved using stochastic optimization, but this is hardly applicable in high dimension. Another option is to impose a *parametric* form for this **potential**, for instance expansion in an *RKHS* or a **deep-network approximation**. (e.g. the **Wasserstein GAN** [Arjovsky et al., 2017]). A last issue is

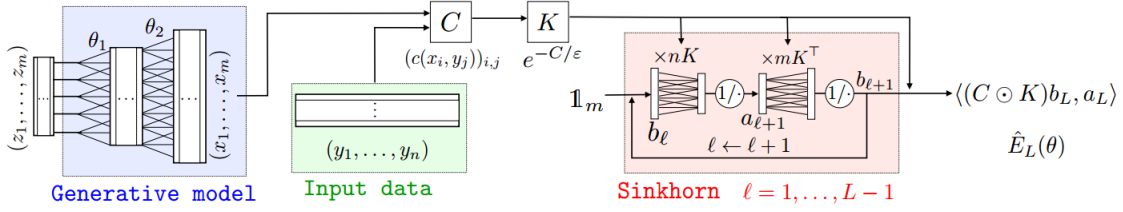


Figure 1: For a given fixed set of samples  $(z_1, \dots, z_m)$ , and input data  $(y_1, \dots, y_n)$ , flow diagram for the computation of Sinkhorn loss function  $\theta \mapsto \hat{E}_\varepsilon^{(L)}(\theta)$ . This function is the one on which automatic differentiation is applied to perform parameter learning. The display shows a simple 2-layer neural network  $g_\theta : z \mapsto x$ , but this applies to any generative model.

Figure 4: Learning generative model with Sinkhorn divergence [Genevay et al., 2018]

that it is **unstable** numerically because it requires the computation of the gradient  $\nabla f(x)$  of the dual potential  $f$ .

We can also solve the primal problem directly, which is a semi-discrete problem. Note that  $\beta_n = \sum_{i=1}^n \frac{1}{n} \delta_{\mathbf{x}_i}$ ,  $\mathbf{x}_i \subset \mathcal{X}$ , thus the objective becomes

$$\begin{aligned} \min_{\{\pi_i \in \mathcal{M}_+(\mathcal{Z})\}_i^n} \quad & \sum_i^n \int_{\mathcal{Z}} c(h_\theta(z), \mathbf{x}_i) d\pi_i(z) \\ \text{where} \quad & \sum_i d\pi_i(z) = \frac{1}{n} \\ & \sum_i \pi_i(z) = \zeta \end{aligned} \quad (46)$$

Based on (46), we can obtain gradient

$$\nabla_\theta \mathcal{E}(\theta) = \sum_i^n \int_{\mathcal{Z}} [\nabla_\theta h_\theta(z)]^T \nabla_1 c(h_\theta(z), \mathbf{x}_i) d\pi_i(z), \quad (47)$$

where  $\nabla_1 c(x, y)$  is the gradient of  $c$  w.r.t. first argument  $x$ . Note that as opposed to (45), this formula does not involve computing the gradient of the potentials being solutions of the dual OT problem.

The class of estimators obtained using  $\mathcal{L} = \mathcal{L}_c$ , often called **minimum Kantorovich estimators**, was initially introduced in [Bassetti et al., 2006], also [Canas and Rosasco, 2012]; It has been used in the context of generative models by [Montavon et al., 2016] to train restricted Boltzmann machines and in [Bernton et al., 2017] in conjunction with approximate Bayesian computations. Approximations of these computations using Deep Network are used to train deep generative models for both **GAN** [Arjovsky et al., 2017] and **Variational Auto-Encoder (VAE)** [Tolstikhin et al., 2018]; see also [Genevay et al., 2017, 2018, Salimans et al., 2018]. Note that the use of **Sinkhorn divergences** for parametric model fitting is used routinely for **shape matching** and **registration**, see [Gold et al., 1998, Rangarajan and Chui, 2000, Myronenko and Song, 2010, Feydy et al., 2017].

OT can also be applied to **metric learning** to learn cost function  $c$  [Wang and Guibas, 2012, Cuturi and Avis, 2014, Zen et al., 2014, Huang et al., 2016]; or to **domain adaptation** [Flamary et al., 2016] and transfer learning [Pan and Yang, 2009].

## References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- Federico Bassetti, Antonella Bodini, and Eugenio Regazzini. On minimum kantorovich distance estimators. *Statistics & probability letters*, 76(12):1298–1302, 2006.
- Espen Bernton, Pierre E Jacob, Mathieu Gerber, and Christian P Robert. Inference in generative models using the wasserstein distance. *arXiv preprint arXiv:1701.05146*, 1(8):9, 2017.
- Guillermo Canas and Lorenzo Rosasco. Learning probability measures with respect to optimal transport metrics. *Advances in Neural Information Processing Systems*, 25, 2012.
- Marco Cuturi and David Avis. Ground metric learning. *The Journal of Machine Learning Research*, 15(1):533–564, 2014.
- Jean Feydy, Benjamin Charlier, François-Xavier Vialard, and Gabriel Peyré. Optimal transport for diffeomorphic registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 291–299. Springer, 2017.
- R Flamary, N Courty, D Tuia, and A Rakotomamonjy. Optimal transport for domain adaptation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 1, 2016.
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Gan and vae from an optimal transport point of view. *arXiv preprint arXiv:1706.01807*, 2017.
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–1617. PMLR, 2018.
- Steven Gold, Anand Rangarajan, Chien-Ping Lu, Suguna Pappu, and Eric Mjolsness. New algorithms for 2d and 3d point matching: pose estimation and correspondence. *Pattern recognition*, 31(8):1019–1031, 1998.
- Gao Huang, Chuan Guo, Matt J Kusner, Yu Sun, Fei Sha, and Kilian Q Weinberger. Supervised word mover’s distance. *Advances in neural information processing systems*, 29, 2016.
- Grégoire Montavon, Klaus-Robert Müller, and Marco Cuturi. Wasserstein training of restricted boltzmann machines. *Advances in Neural Information Processing Systems*, 29, 2016.
- Andriy Myronenko and Xubo Song. Point set registration: Coherent point drift. *IEEE transactions on pattern analysis and machine intelligence*, 32(12):2262–2275, 2010.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- Gabriel Peyr and Marco Cuturi. Computational optimal transport: With applications to data science. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019. ISSN 1935-8237.
- Anand Rangarajan and Haili Chui. Applications of optimizing neural networks in medical image registration. In *Artificial Neural Networks in Medicine and Biology*, pages 99–104. Springer, 2000.
- Tim Salimans, Dimitris Metaxas, Han Zhang, and Alec Radford. Improving gans using optimal transport. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.



- Filippo Santambrogio. *Optimal transport for applied mathematicians*, volume 55. Springer, 2015.
- I Tolstikhin, O Bousquet, S Gelly, and B Schölkopf. Wasserstein auto-encoders. In *6th International Conference on Learning Representations (ICLR 2018)*. OpenReview. net, 2018.
- Fan Wang and Leonidas J Guibas. Supervised earth movers distance learning and its computer vision applications. In *European Conference on Computer Vision*, pages 442–455. Springer, 2012.
- Gloria Zen, Elisa Ricci, and Nicu Sebe. Simultaneous ground metric learning and matrix factorization with earth mover’s distance. In *2014 22nd International Conference on Pattern Recognition*, pages 3690–3695. IEEE, 2014.