# Differential Properties of Sinkhorn Approximation for Learning with Wasserstein Distance

# Giulia Luise<sup>1</sup> Alessandro Rudi<sup>2</sup> Massimiliano Pontil<sup>1,3</sup> Carlo Ciliberto

<sup>1</sup>Department of Computer Science, University College London, London, UK. <sup>2</sup>INRIA- Sierra-Project team ENS, Paris, France.

<sup>3</sup>Computational Statistics and Machine Learning, Istituto Italiano di Tecnologia, Genova, Italy.







ISTITUTO ITALIANO DI TECNOLOGIA

## Abstract

Applications of optimal transport have recently received remarkable attention thanks to the computational advantages of entropic regularization. However, in most situations the Sinkhorn approximation of the Wasserstein distance is replaced by a regularized version that is less accurate but easy to differentiate. In this work we characterize the differential properties of the original Sinkhorn distance, proving that it enjoys the same smoothness as its regularized version and we explicitly provide an efficient algorithm to compute its gradient. We show that this result benefits both theory and applications: on one hand, high order smoothness confers statistical guarantees to learning with Wasserstein approximations. On the other hand, the gradient formula allows us to efficiently solve learning and optimization problems in practice.

In the following, the barycenters computed with  $S_{\lambda}$  and  $\tilde{S}_{\lambda}$  constitute an example of how  $S_{\lambda}$  is not affected by the same oversmoothing effect as  $\tilde{S}_{\lambda}$ .



#### **1. BACKGROUND: Entropic regularizations of Wasserstein distance**

Optimal transport theory investigates how to compare probability measures over a metric space X.

**Discrete Setting:** discrete probability measures of the form  $\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$  with  $(x_i)_{i=1}^{n}$  and the vector weight  $a = (a_1, \ldots, a_n)^{\top} \in \Delta_n$  in the *n*-dimensional simplex  $\Delta_n = \left\{ p \in \mathbb{R}^n_+ \mid p^{\top} \mathbb{1}_n = 1 \right\}$ . Given  $\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$  and  $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$ , the Wasserstein distance between  $\mu$  and  $\nu$  is defined as

$$\mathbf{W}_{p}^{p}(\mu,\nu) = \min_{T \in \Pi(a,b)} \langle T, M \rangle$$
(1)

where  $M \in \mathbb{R}^{n \times m}$  is the *cost matrix* with entries  $M_{ij} = \mathsf{d}(x_i, y_j)^p$  and  $\Pi(a, b)$  denotes the *transportation polytope* 

$$\Pi(a,b) = \{ T \in \mathbb{R}^{n \times m}_+ \mid T \mathbb{1}_m = a, \quad T^\top \mathbb{1}_n = b \}.$$

**Regularization of Wasserstein distance** 

Set  $h(T) := -\sum_{i,j=1}^{n,m} T_{ij}(\log T_{ij} - 1)$  and  $T_{\lambda} = \operatorname{argmin}_{T \in \Pi(a,b)} \langle T, M \rangle - \frac{1}{\lambda}h(T)$ 

**Definition** Given  $\mu$  and  $\nu$  as above, entropic regularizations of the Wasserstein distance, referred to as Sinkhorn distances [2] are defined as

$$\widetilde{\mathbf{S}}_{\lambda}(a,b) = \langle T_{\lambda}, M \rangle - \frac{1}{\lambda} h(T_{\lambda}) \text{ and } \mathbf{S}_{\lambda} = \langle T_{\lambda}, M \rangle.$$
 (2)

**Proposition** Let  $\lambda > 0$ . For any pair of discrete measures  $\mu, \nu \in \mathcal{P}(X)$  with respective weights  $a \in \Delta_n$  and  $b \in \Delta_m$ , we have

 $\left|S_{\lambda}(\mu,\nu)-W(\mu,\nu)\right|\leq c_{1}e^{-\lambda}$   $\left|\widetilde{S}_{\lambda}(\mu,\nu)-W(\mu,\nu)\right|\leq c_{2}\lambda^{-1},$ 

where  $c_1, c_2$  are constants independent of  $\lambda$ , depending on the support of  $\mu$  and  $\nu$ .

Figure: Barycenter of Nested Ellipses: (Left) Sample input data. (Middle) Barycenter with  $S_{\lambda}$ . (Right) Barycenter with  $S_{\lambda}$ .

#### **3. LEARNING WITH SINKHORN LOSS: SETTING**

**Problem Setting:**  $\mathcal{X}$  input space,  $\mathcal{Y} = \Delta_n$  a set of histograms (output space). **Goal:** approximate a minimizer of the *expected risk* 

$$\min_{f:\mathcal{X}\to\mathcal{Y}} \mathcal{E}(f), \qquad \qquad \mathcal{E}(f) = \int_{\mathcal{X}\times\mathcal{Y}} \mathcal{S}(f(x), y) \, d\rho(x, y) \tag{4}$$

given a finite number of training points  $(x_i, y_i)_{i=1}^{\ell}$  independently sampled from the unknown distribution  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$ . The loss function  $\mathcal{S} : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$  measures prediction errors and in our setting corresponds to either  $S_{\lambda}$  or  $\tilde{S}_{\lambda}$ .

**Structured Prediction Estimator.** Given a training set  $(x_i, y_i)_{i=1}^{\ell}$ , we consider  $\hat{f} : \mathcal{X} \to \mathcal{Y}$  the structured prediction estimator proposed in [1], defined as

$$\hat{f}(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmin}} \sum_{i=1}^{\ell} \alpha_i(x) \,\mathcal{S}(y, y_i)$$
(5)

for any  $x \in \mathcal{X}$ .

(3)

The weights  $\alpha_i(x)$  are learned from the data and can be interpreted as scores suggesting the candidate output distribution y to be close to a specific output distribution  $y_i$  observed in training *according to the metric* S.

 $\alpha$  are obtained via a kernel-based approach: given a positive definite kernel  $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , we have

**Question:** Is  $S_{\lambda}$  a more natural approximation of the Wasserstein distance W?



Figure: Comparison of the sharp (Blue) and regularized (Orange) barycenters of two Dirac's deltas (Black) centered in 0 and 20 for different values of  $\lambda$ .

#### **2. DIFFERENTIAL PROPERTIES**

#### Sinkhorn maps are $C^{\infty}$ smooth in the interior of $\Delta_n$ .

The results are obtained leveraging the Implicit Function Theorem via a proof technique analogous to that in [3].

**Theorem** For any  $\lambda > 0$ , the Sinkhorn distances  $S_{\lambda}$  and  $S_{\lambda} : \Delta_n \times \Delta_n \to \mathbb{R}$  are  $C^{\infty}$  in the interior of their domain.

**Sketch of the proof.** Step 1: since  $S_{\lambda}$  and  $\tilde{S}_{\lambda}$  are smooth as functions of  $T^{\lambda}$ , it is enough to show that  $T^{\lambda}$  is a smooth function of a, b.

Step 2: by Sinkhorn's scaling theorem,  $T^{\lambda} = \text{diag}(e^{\lambda \alpha^{\lambda}})e^{-\lambda M}\text{diag}(e^{\lambda \beta^{\lambda}})$ . Then,  $T^{\lambda}$  is smooth if  $(\alpha^{\lambda}, \beta^{\lambda})$  is smooth as a function of (a, b). Step 3: let us set  $\alpha(x) = (\alpha_1(x), \dots, \alpha(x))^\top = (K + \gamma \ell I)^{-1} K_x$ (6)

where  $\gamma > 0$  is a regularization parameter while  $K \in \mathbb{R}^{\ell \times \ell}$  and  $K_x \in \mathbb{R}^n$  are respectively the empirical kernel matrix with entries  $K_{ij} = k(x_i, x_j)$  and the evaluation vector with entries  $(K_x)_i = k(x, x_i)$ , for any  $i, j = 1, ..., \ell$ .

#### **4. LEARNING WITH SINKHORN LOSS: STATISTICAL ANALYSIS**

#### Thanks to the smoothness of $S_{\lambda}$ we can prove consistency and learning rates of the estimator. Thanks to the gradient we can solve the problem in practice!

**Theorem**[Universal Consistency] Let  $\mathcal{Y} = \Delta_n^{\epsilon}$ ,  $\lambda > 0$  and  $\mathcal{S}$  be either  $S_{\lambda}$  or  $S_{\lambda}$ . Let k be a bounded continuous universal kernel on  $\mathcal{X}$ . For any  $\ell \in \mathbb{N}$  and any distribution  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$  let  $\hat{f}_{\ell} : \mathcal{X} \to \mathcal{Y}$  be the estimator in (5) trained with  $(x_i, y_i)_{i=1}^{\ell}$  points independently sampled from  $\rho$  and  $\gamma_{\ell} = \ell^{-1/4}$ . Then

$$\lim_{\ell \to \infty} \mathcal{E}(\widehat{f}_{\ell}) = \min_{f: \mathcal{X} \to \mathcal{Y}} \mathcal{E}(f) \quad with \text{ probability } 1.$$
(7)

**Theorem**[Learning Rates](informal) Let  $\mathcal{Y} = \Delta_n^{\epsilon}$ ,  $\lambda > 0$  and  $\mathcal{S}$  be either  $\tilde{S}_{\lambda}$  or  $S_{\lambda}$ . Under suitable regularity assumption on  $\rho$ , the estimator  $\hat{f}_{\ell} : \mathcal{X} \to \mathcal{Y}$  with  $\ell$  training points independently sampled from  $\rho$  and with  $\gamma = \ell^{-1/2}$ 

$$\mathcal{E}(f) - \min_{f:\mathcal{X}\to\mathcal{Y}} \mathcal{E}(f) = O(\ell^{-1/4})$$
(8)

holds with high probability with respect to the sampling of training data.

**Role of the smoothness in the statistical analysis:** the proof is technical but a key role is played by the smoothness of Sinkhorn maps shown before. This is the first universal consistency result for learning with Sinkhorn loss!

 $\mathcal{L}(a,b;\alpha,\beta) = \alpha^{\top} a + \beta^{\top} b - \frac{1}{\lambda} \sum_{i,j=1}^{n,m} e^{-\lambda(M_{ij} - \alpha_i - \beta_j)}$ 

and recall that  $(\alpha^{\lambda}, \beta^{\lambda}) = \operatorname{argmax}_{\alpha,\beta} \mathcal{L}(a, b; \alpha, \beta)$ . The smoothness of  $(\alpha^{\lambda}, \beta^{\lambda})$  is proved using the Implicit Function theorem and follows from the smoothness and strong convexity in  $\alpha, \beta$  of the function  $\mathcal{L}$ .

The Implicit Function Theorem also provides a formula for the gradient of  $S_{\lambda}$ :

Input: 
$$a \in \Delta_n$$
,  $b \in \Delta_m$ , cost matrix  $M \in \mathbb{R}^{n,m}_+$ ,  $\lambda > 0$ .  
 $T = \text{SINKHORN}(a, b, M, \lambda)$ ,  $\overline{T} = T_{1:n,1:(m-1)}$   
 $L = T \odot M$ ,  $\overline{L} = L_{1:n,1:(m-1)}$   
 $D_1 = \text{diag}(T \mathbb{1}_m)$ ,  $D_2 = \text{diag}(\overline{T}^\top \mathbb{1}_n)^{-1}$   
 $H = D_1 - \overline{T}D_2\overline{T}^\top$ ,  $\mathbf{f} = -L\mathbb{1}_m + \overline{T}D_2\overline{L}^\top\mathbb{1}_n$ ,  
 $\mathbf{g} = H^{-1}$  f  
Return:  $\mathbf{g} - \mathbb{1}_n (\mathbf{g}^\top \mathbb{1}_n)$   
Algorithm 1: Gradient of  $S_\lambda$ 

The routine for the gradient is used to implement optimization problems with  $S_{\lambda}$  as loss. While solutions of optimization with  $\tilde{S}_{\lambda}$  are often 'blurry',  $S_{\lambda}$  preserves the *sharpness* of the data.

**5. EXPERIMENTS** 

We evaluated the Sinkhorn distances with the estimator (5) in an image reconstruction problem: given an image depicting a drawing, the goal is to learn how to reconstruct the lower half of the image (output) given the upper half (input).

	Keconst	ruction	Error on	QuickDraw (%)
# Classes	$\mathrm{S}_{\lambda}$	$\widetilde{\mathbf{S}}_{oldsymbol{\lambda}}$	Hell	KDE
2	$3.7\pm0.6$	$4.9\pm0.9$	$8.0 \pm 2.4$	$12.0 \pm 4.1$
4	$22.2 \pm 0.9$	$31.8 \pm 1.1$	$29.2 \pm 0.8$	$40.8 \pm 4.2$
<b>10</b>	$38.9 \pm 0.9$	$44.9 \pm 2.5$	$48.3 \pm 2.4$	$64.9 \pm 1.4$



Figure: Average reconstruction errors of the Sinkhorn, Hellinger, and KDE estimators on the Google QuickDraw reconstruction problem. On the right, example of reconstruction on MNIST.

### REFERENCES

- [1] Carlo Ciliberto, Lorenzo Rosasco, and Alessandro Rudi. "A Consistent Regularization Approach for Structured Prediction". In: *Advances in Neural Information Processing Systems 29*.
- [2] Marco Cuturi. "Sinkhorn Distances: Lightspeed Computation of Optimal Transport". In: *Advances in Neural Information Processing Systems 26*. Ed. by C. J. C. Burges et al. Curran Associates, Inc., 2013, pp. 2292–2300.
- [3] Rémi Flamary et al. "Wasserstein discriminant analysis". In: *Machine Learning* (2018).