# Differential Properties of Sinkhorn Approximation for Learning with Wasserstein Distance

# **CONTRIBUTIONS IN A NUTSHELL**

**Goal:** In a supervised learning setting, learn functions  $f : \mathcal{X} \to \mathcal{Y}$ where the **output** space  $\mathcal{Y}$  is a set of discrete **probability distributions**.

### **Tool:** Optimal Transport

**Approach:** Use Sinkhorn approximations as loss functions

### **Contributions:**

- Characterise the **differential properties** of Sinkhorn approximations.
- Provide learning bounds for learning with Sinkhorn loss(es), adopting a structured prediction perspective.

# BACKGROUND

Optimal transport theory compares probability measures over a metric space. Wasserstein distance (discrete setting):

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



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**

**Definition** Given  $\mu$  and  $\nu$  as above, entropic regularizations of the Wasserstein distance, referred to as Sinkhorn distances [1] are defined as

$$\tilde{\mathsf{S}}_{\lambda}(a,b) = \langle T_{\lambda}, M \rangle - \frac{1}{\lambda}h(T_{\lambda}) \text{ and } \mathsf{S}_{\lambda}(a,b) = \langle T_{\lambda}, M \rangle$$

where

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

**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

$$\left| S_{\lambda}(\mu,\nu) - W(\mu,\nu) \right| \le c_1 e^{-\lambda} \qquad \left| \tilde{S}_{\lambda}(\mu,\nu) - W(\mu,\nu) \right| \le c_2/2$$

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

**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  $\lambda$ .

Giulia Luise <sup>1</sup> Alessandro Rudi <sup>2</sup> Massimiliano Pontil <sup>1,3</sup> Carlo Ciliberto <sup>1,4</sup>

<sup>1</sup>Department of Computer Science, University College of London, UK. <sup>2</sup>INRIA - Sierra-Project team, École Normale Supérieure, PSL Research, Paris, France <sup>3</sup>Computational Statistics and Machine Learning, Istituto Italiano di Tecnologia, Genova, Italy. <sup>4</sup>Department of Electrical and Electronic Engineering, Imperial College London, UK.









lambda,

# **DIFFERENTIAL PROPERTIES**

We characterise regularity properties of Sinkhorn maps.

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

**Proof (sketch).** The proof is organized in the following steps: Step 1:  $S_{\lambda}$  and  $\tilde{S}_{\lambda}$  are smooth as functions of  $T^{\lambda} \rightarrow sufficient$  to show that  $T^{\lambda}$  is smooth in a, b. Step 2: Set  $(\alpha^*, \beta^*) = \operatorname{argmax}_{\alpha, \beta} \mathcal{L}(a, b; \alpha, \beta)$ , with

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

By Sinkhorn's scaling theorem,  $T^{\lambda} = \text{diag}(e^{\lambda \alpha^{\star}})e^{-\lambda M}\text{diag}(e^{\lambda \beta^{\star}}) \rightarrow T^{\lambda}$  is smooth if  $(\alpha^{\star}, \beta^{\star})$  is smooth as a function of (a, b).

Step 3: The smoothness of  $(\alpha^{\star}, \beta^{\star})$  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  $M \in \mathbb{R}^{n,m}_+$ ,  $\lambda > 0$ .  $T = \text{SINKHORN}(a, b, M, \lambda), \qquad T = T_{1:n,1:(m-1)}$  $L = T \odot M$ ,  $L = L_{1:n,1:(m-1)}$  $D_1 = \text{diag}(T1_m), \quad D_2 = \text{diag}(\bar{T}^{\top}1_n)^{-1}$  $H = D_1 - TD_2T^{\top}$ ,  $f = -L1_m + \bar{T}D_2\bar{L}^{\top}1_n$ ,  $g = H^{-1}$ **Return:**  $g - 1_n(g^{\top}1_n)$ **Algorithm 1:** Gradient of  $S_{\lambda}$ 

**Synthetic experiment.** Find the barycenter of nested ellipses.



**Figure:** (Left) Sample input data. (Middle) Barycenter with  $\tilde{S}_{\lambda}$ . (Right) Barycenter with  $S_{\lambda}$ . While solutions of optimization with  $\tilde{S}_{\lambda}$  are often 'blurry',  $S_{\lambda}$  preserves the sharpness of the data.





# LEARNING WITH SINKHORN LOSS: SETTING

**Goal:** approximate a minimizer of the *expected risk* 

given a training set  $(x_i, y_i)_{i=1}^\ell$  independently sampled from  $\rho$ . The loss function  $\mathcal{S} : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ in our setting is either  $S_{\lambda}$  or  $S_{\lambda}$ .

 $\hat{f}(x) = \operatorname*{argmin}_{u \in \mathcal{V}} \sum_{i=1}^{c} \alpha_i(x) \mathcal{S}(y, y_i), \quad \text{for any } x \in \mathcal{X}.$ (1)

# STATISTICAL ANALYSIS

**Theorem** (Universal Consistency) Let  $\mathcal{Y} = \Delta_n^{\epsilon}$ ,  $\lambda > 0$  and  $\mathcal{S}$  be either  $\tilde{S}_{\lambda}$  or  $S_{\lambda}$ . Let k be a bounded continuous universal kernel on  $\mathcal{X}$ . For any  $\ell \in \mathbb{N}$  and any distribution ho on  $\mathcal{X} imes \mathcal{Y}$  let  $f_\ell : \mathcal{X} o \mathcal{Y}$ be the estimator in (1) trained with  $\ell$  points sampled from  $\rho$ . Then  $\lim_{\ell \to \infty} \mathcal{E}(\widehat{f_{\ell}}) = \min_{f: \mathcal{X} \to \mathcal{V}} \mathcal{E}(f) \quad with \text{ probability } 1.$ 

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

Role of the smoothness : the proof is technical but essentially allows to embed the problem into a Hilbert setting. This is the first universal consistency result for learning with Sinkhorn loss!

# EXPERIMENTS

Image Reconstruction

**Goal:** given the upper half of Google QuickDraw images, predict their bottom half. ruction Error (%) KEC KDE Hell  $8.0 \pm 2.4$  $12.0 \pm 4.1$ 0.9Ø  $\pm 1.1 \quad 29.2 \pm 0.8 \quad 40.8 \pm 4.2$  $\pm 2.5 \quad 48.3 \pm 2.4 \quad 64.9 \pm 1.4$ 

	Reconstr	
# Cls.	$S_{\boldsymbol{\lambda}}$	$\tilde{S}_{\boldsymbol{\lambda}}$
2	$3.7\pm0.6$	$4.9 \pm$
4	$22.2 \pm 0.9$	$31.8 \pm$
<b>10</b>	$38.9 \pm 0.9$	$44.9 \pm$

Figure: (Left) Reconstruction error of Sinkhorn, Hellinger and KDE. Misclassification rate of the base SVM classifier: 0.02, 0.07, 0.17. (Right) Examples of training and reconstructed data.

## REFERENCES





**Problem Setting:**  $\mathcal{X}$  input space,  $\mathcal{Y} = \Delta_n$  a set of normalized histograms (output space).

$$\mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{S}(f(x), y) \, d\rho(x, y)$$

**Structured Prediction Estimator.** Given a training set  $(x_i, y_i)_{i=1}^\ell$ , we consider  $\widehat{f}: \mathcal{X} \to \mathcal{Y}$  the structured prediction estimator proposed in [2], defined as

The weights  $\alpha_i(x) \rightarrow$  Are scores measuring similarity of test point and training points  $\rightarrow$  Are obtained via Kernel Ridge Regression

We use the smoothness of  $S_{\lambda}$  to prove consistency and learning rates of the estimator

**Theorem** (Learning Rates -informal) Let  $\mathcal{Y} = \Delta_n^{\epsilon}$ ,  $\lambda > 0$  and  $\mathcal{S}$  and  $\hat{f}_{\ell}$  as above. Then,  $\mathcal{E}(\widehat{f}_{\ell}) - \min_{f \colon \mathcal{X} \to \mathcal{Y}} \mathcal{E}(f) = O(\ell^{-1/4})$