Differential Properties of Sinkhorn Approximation for Learning with Wasserstein Distance

Giulia Luise 1 Alessandro Rudi 2 Massimiliano Pontil 1,3 Carlo Ciliberto 1,4





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

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

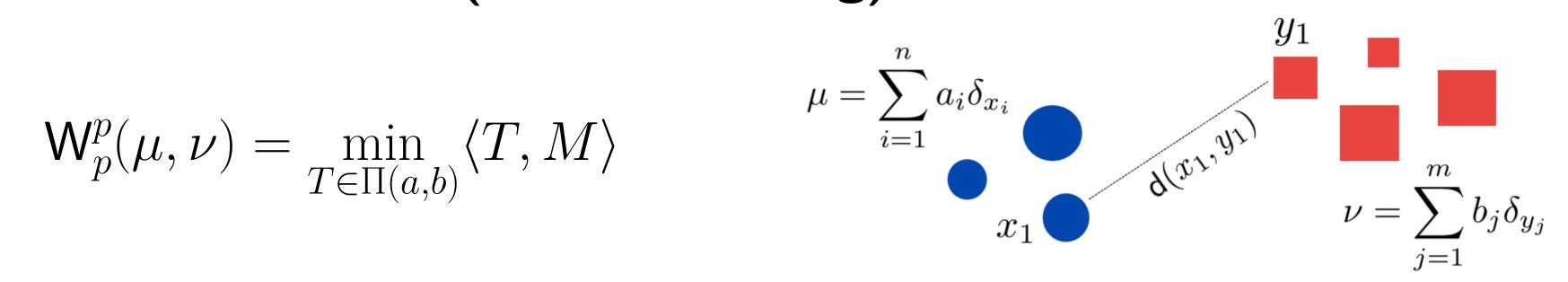
Contributions:

- 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):

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



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

$$\Pi(a,b) = \{ T \in \mathbb{R}_+^{n \times m} \mid T1_m = a, T^{\mathsf{T}}1_n = b \}.$$

Regularization of Wasserstein distance

Definition Given μ and ν as above, entropic regularizations of the Wasserstein distance, referred to as Sinkhorn distances [1] are defined as

$$ilde{\mathsf{S}}_{\lambda}(a,b) \; = \; \langle T_{\lambda},M
angle - rac{1}{\lambda} h(T_{\lambda}) \quad ext{and} \quad \mathsf{S}_{\lambda}(a,b) \; = \; \langle T_{\lambda},M
angle \; ,$$

where

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

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

$$\left| S_{\lambda}(\mu,\nu) - W(\mu,\nu) \right| \le c_1 e^{-\lambda}$$
 $\left| \tilde{S}_{\lambda}(\mu,\nu) - W(\mu,\nu) \right| \le c_2/lambda$, with c_1,c_2 constants independent of λ , depending on the support of μ and ν .

Question: Is S_{λ} a more natural approximation of the Wasserstein distance W?

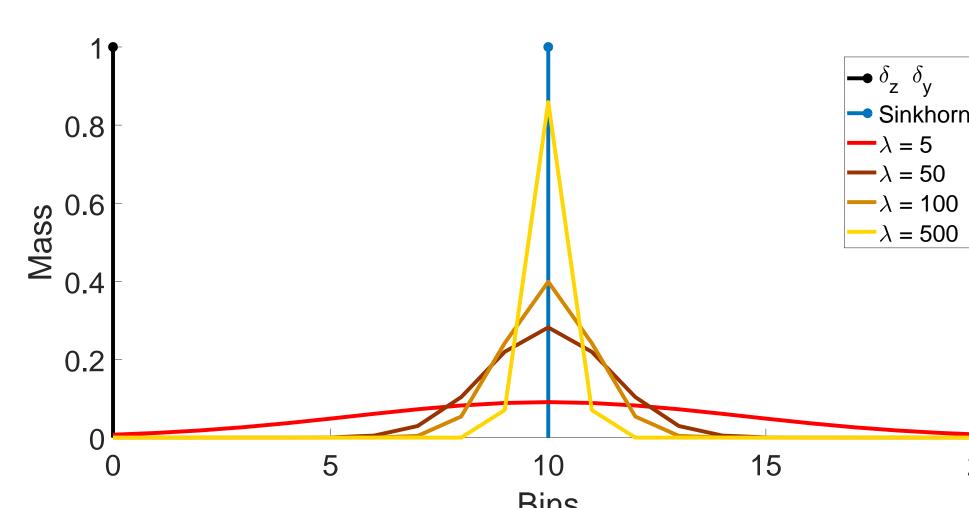


Figure: Comparison of the sharp (Blue) and regularized (Orange) barycenters of two Dirac's deltas (Black) centered in 0 and 20 for different values of λ .

DIFFERENTIAL PROPERTIES

We characterise regularity properties of Sinkhorn maps.

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

Proof (sketch). The proof is organized in the following steps:

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

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

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

Step 3: The smoothness of (α^*, β^*) is proved using the Implicit Function theorem and follows from the smoothness and strong convexity in α, β of the function \mathcal{L} .

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

Input:
$$a \in \Delta_n, \ b \in \Delta_m$$
, cost matrix $M \in \mathbb{R}^{n,m}_+$, $\lambda > 0$. $T = \text{SINKHORN}(a, b, M, \lambda), \quad \bar{T} = T_{1:n,1:(m-1)}$ $L = T \odot M, \quad \bar{L} = L_{1:n,1:(m-1)}$ $D_1 = \text{diag}(T1_m), \quad D_2 = \text{diag}(\bar{T}^\top 1_n)^{-1}$ $H = D_1 - \bar{T}D_2\bar{T}^\top$, $f = -L1_m + \bar{T}D_2\bar{L}^\top 1_n$, $g = H^{-1}$ f Return: $g - 1_n (g^\top 1_n)$

Algorithm 1: Gradient of S_{λ}

Synthetic experiment. Find the barycenter of nested ellipses.

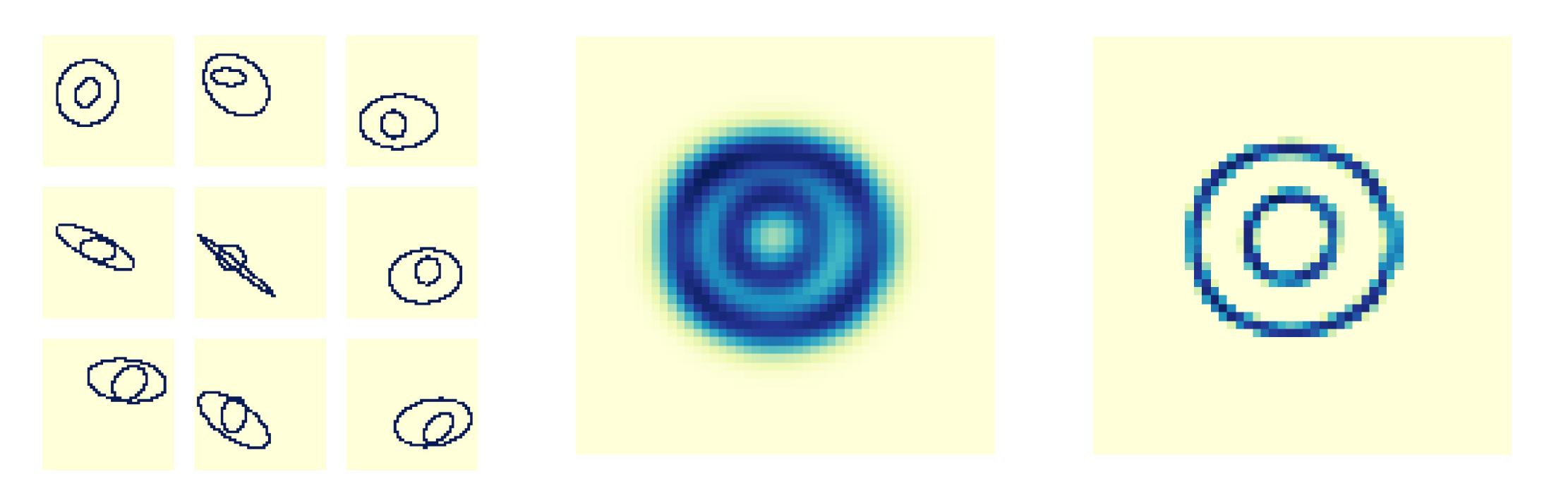


Figure: (Left) Sample input data. (Middle) Barycenter with \tilde{S}_{λ} . (Right) Barycenter with S_{λ} . While solutions of optimization with \hat{S}_{λ} are often 'blurry', S_{λ} preserves the sharpness of the data.

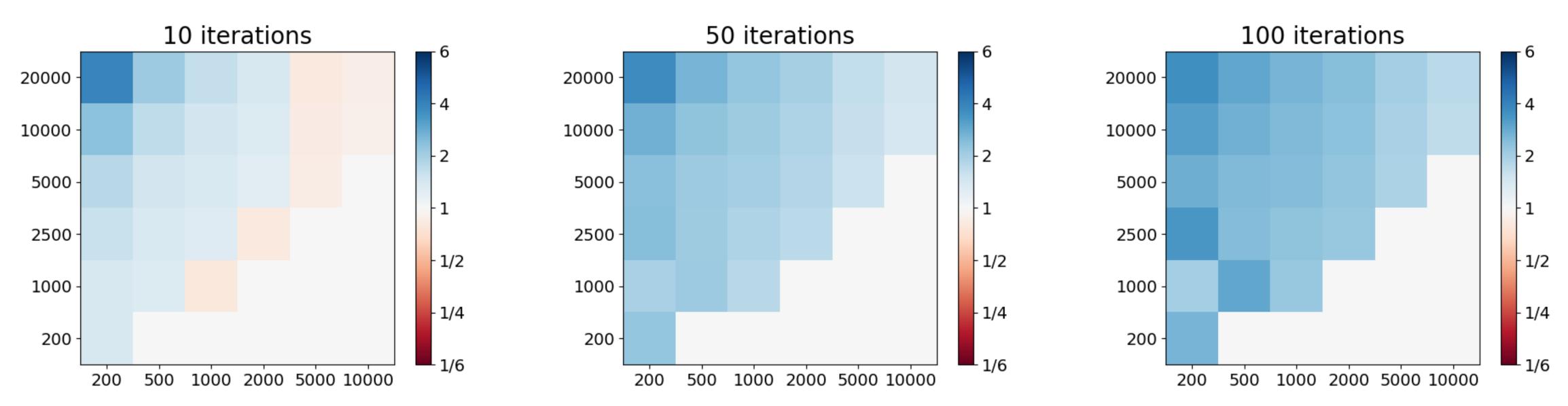


Figure: Ratio of time(autodiff) / time(Alg.1) for 10, 50, and 100 iterations of Sinkhorn algorithm

LEARNING WITH SINKHORN LOSS: SETTING

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

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

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

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

$$\hat{f}(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmin}} \sum_{i=1}^{\ell} \alpha_i(x) \, \mathcal{S}(y, y_i), \quad \text{for any } x \in \mathcal{X}.$$
 (1)

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

STATISTICAL ANALYSIS

We use the smoothness of S_{λ} to prove consistency and learning rates of the estimator

Theorem (Universal Consistency) Let $\mathcal{Y} = \Delta_n^{\epsilon}$, $\lambda > 0$ and \mathcal{S} be either \tilde{S}_{λ} or S_{λ} . 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 $\widehat{f_\ell}:\mathcal X o\mathcal Y$ be the estimator in (1) trained with ℓ points sampled from ρ . Then

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

Theorem (Learning Rates -informal) Let $\mathcal{Y}=\Delta_n^\epsilon$, $\lambda>0$ and \widehat{f}_ℓ as above. Then,

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

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.

Reconstruction Error (%)			
λ	$\widetilde{S}_{oldsymbol{\lambda}}$	Hell	KDE
<u> </u>	10 ± 00	$\frac{20 + 21}{}$	12 0 +

 38.9 ± 0.9 44.9 ± 2.5 48.3 ± 2.4 64.9 ± 1.4

 22.2 ± 0.9 31.8 ± 1.1 29.2 ± 0.8 40.8 ± 4.2



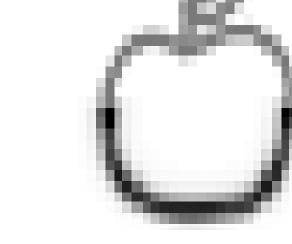


Figure: (Left) Reconstruction error of Sinkhorn, Hellinger and KDE. Misclassification rate of the base SVM classifier: 0.02, 0.07, 0.17. (Right) Examples of training and reconstructed data.

REFERENCES

- M. Cuturi, Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances, NIPS 2013
- C. Ciliberto et al, A Consistent Regularization Approach for Structured Prediction, NIPS 2016
- R. Flamary et al, Wasserstein Discriminative Analisys, Journal of Machine Learning 2018