



Objectives

- Introduce a novel strategy for efficiently approximating the optimal transport (OT) distance between two discrete measures.
- Propose the SCREENKHORN algorithm: solve a smaller Sinkhorn problem while ensuring approximation.
- Illustrate the efficiency of SCREENKHORN on complex tasks such as domain adaptation with regularized OT.

Introduction

• OT is a method for comparing probability distributions with the ability to incorporate spatial information.



(courtesy of M. Cuturi)

• Given $\boldsymbol{\mu} = \sum_{i=1}^{n} \boldsymbol{\mu}_i \delta_{\boldsymbol{x}_i}$ and $\boldsymbol{\nu} = \sum_{j=1}^{m} \boldsymbol{\nu}_j \delta_{\boldsymbol{x}_j}$ two discrete probability distributions and a nonnegative cost matrix $C = (C_{ij}) \in \mathbb{R}^{n \times m}_+$, the OT (Wasserstein) distance writes as

 $\mathcal{S}(\boldsymbol{\mu},\boldsymbol{\nu}) = \min_{\boldsymbol{P}\in\boldsymbol{\Pi}(\boldsymbol{\mu},\boldsymbol{\nu})} \langle \boldsymbol{C},\boldsymbol{P} \rangle,$

where

 $\Pi(\boldsymbol{\mu},\boldsymbol{\nu}) = \{\boldsymbol{P} \in \mathbb{R}^{n \times m}, \boldsymbol{P} \boldsymbol{1}_m = \boldsymbol{\mu}, \boldsymbol{P}^\top \boldsymbol{1}_n = \boldsymbol{\nu}\}.$

• Sinkhorn Divergence: The entropic regularization of OT distances relies on the addition of a penalty term:

$$\mathcal{S}_{\eta}(\boldsymbol{\mu},\boldsymbol{\nu}) = \min_{\boldsymbol{P}\in\boldsymbol{\Pi}(\boldsymbol{\mu},\boldsymbol{\nu})} \{ \langle \boldsymbol{C}, \boldsymbol{P} \rangle - \eta H(\boldsymbol{P}) \},\$$

where $H(\mathbf{P}) = -\sum_{i,j} \mathbf{P}_{ij} \log(\mathbf{P}_{ij})$ is the negative entropy and $\eta > 0$ is a regularization parameter.

Sinkhorn Algorithm fo Mokhtar Z. Alaya, Maxime Bérar, LITIS EA4108, Unive	D r , Gille
UNIVERSITÉ DE ROUEN	itis
• The dual of Sinkhorn divergence is given by	• (
$\mathcal{S}_{\eta}^{d}(\boldsymbol{\mu},\boldsymbol{\nu}) = \min_{\boldsymbol{u}\in\mathbb{R}^{n},\boldsymbol{v}\in\mathbb{R}^{m}} \{\Psi(\boldsymbol{u},\boldsymbol{v})\},\$ where $\Psi(\boldsymbol{u},\boldsymbol{v}) := 1_{n}^{\top}B(\boldsymbol{u},\boldsymbol{v})1_{m} - \langle \boldsymbol{u},\boldsymbol{\mu}\rangle - \langle \boldsymbol{v},\boldsymbol{\nu}\rangle,\$ $B(\boldsymbol{u},\boldsymbol{v}) := \operatorname{diag}(e^{\boldsymbol{u}}) \boldsymbol{K}\operatorname{diag}(e^{\boldsymbol{v}}), \boldsymbol{K} := e^{-\boldsymbol{C}/\eta}.$ • The optimal solution \boldsymbol{P}^{\star} of the primal problem $\mathcal{S}_{\eta}^{d}(\boldsymbol{\mu},\boldsymbol{\nu})$ takes the form $\boldsymbol{P}^{\star} = \operatorname{diag}(e^{\boldsymbol{u}^{\star}}) \boldsymbol{K}\operatorname{diag}(e^{\boldsymbol{v}^{\star}}),\$ where $(\boldsymbol{u}^{\star},\boldsymbol{v}^{\star}) = \operatorname{argmin}_{\boldsymbol{u},\boldsymbol{v}}\{\Psi(\boldsymbol{u},\boldsymbol{v})\}.$	S • E C

Toy Example



Let $(\mathbf{u}^*, \mathbf{v}^*)$ be an optimal solution of $\mathcal{S}_n^{\mathsf{ad}}(\boldsymbol{\mu}, \boldsymbol{\nu})$ then $e^{\mathbf{u}_i^*} = \frac{\varepsilon}{\kappa}$ $\boldsymbol{I}_{\varepsilon,\kappa} = \left\{ i = 1, \dots, n : \boldsymbol{\mu}_i \geq \frac{\varepsilon^2}{\kappa} r_i(\boldsymbol{K}) \right\} \text{ and } \boldsymbol{J}_{\varepsilon,\kappa}$

Screening with a Fixed Number Budget of Points

- Let $\boldsymbol{\xi} = \boldsymbol{\mu} \oslash r(\boldsymbol{K}) \in \mathbb{R}^n$ and $\boldsymbol{\zeta} = \boldsymbol{\nu} \oslash c(\boldsymbol{K}) \in \mathbb{R}^m$ sorted in descending order. To keep only n_b -budget and m_b -budget of points, the parameters κ and ε satisfy $\frac{\varepsilon^2}{\kappa} = \boldsymbol{\xi}_{n_b}$ and $\varepsilon^2 \kappa = \boldsymbol{\zeta}_{m_b}$.
- We restrict the constraints feasibility to the *screened domain* defined by $\mathcal{U}^{sc} \cap \mathcal{V}^{sc}$, where $\mathcal{U}^{sc} = \{ u \in \mathbb{R}^{n_b} : v \in \mathbb{R}^{n_b} \}$ $e^{\boldsymbol{u}_i} \geq \frac{\varepsilon}{\kappa}$ and $\mathcal{V}^{sc} = \{\boldsymbol{v} \in \mathbb{R}^{m_b} : e^{\boldsymbol{v}_j} \geq \varepsilon \kappa\}$, then we derive the screened dual of Sinkhorn divergence problem as

 $\mathcal{S}_{\eta}^{\mathsf{scd}}(\boldsymbol{\mu},\boldsymbol{\nu}) = \min_{\boldsymbol{u}\in\mathcal{U}_{\mathsf{sc}},\boldsymbol{v}\in\mathcal{V}_{\mathsf{sc}}} \{\Psi_{\varepsilon,\kappa}(\boldsymbol{u},\boldsymbol{v})\}.$

Regularized Optimal Transport

es Gasso, Alain Rakotomamonjy

f Rouen Normandy



INSTITUT NATIONAL **DES SCIENCES** APPLIQUÉES **ROUEN NORMANDIE**

Static Screening Test

OT plans present a large number of neglectable values. This favorites the using of *static screening test* like in supervised learning (Lasso).

Based on this idea, we define a so-called **approximate** dual of Sinkhorn divergence

 $\mathcal{S}^{\mathsf{ad}}_{\eta}(\boldsymbol{\mu},\boldsymbol{\nu}) = \min_{\boldsymbol{u}\in\mathcal{C}^n_{\varepsilon},\boldsymbol{v}\in\mathcal{C}^m_{\varepsilon\kappa}} \{\Psi_{\kappa}(\boldsymbol{u},\boldsymbol{v})\},\$

 $\Psi_{\kappa}(\boldsymbol{u},\boldsymbol{v}) := \mathbf{1}_{n}^{\top}B(\boldsymbol{u},\boldsymbol{v})\mathbf{1}_{m} - \langle \kappa \boldsymbol{u},\boldsymbol{\mu} \rangle - \langle \frac{\boldsymbol{v}}{\kappa},\boldsymbol{\nu} \rangle, \text{ and }$ $\mathcal{C}_{\alpha}^{r} = \{ w \in \mathbb{R}^{r} : e^{w_{i}} \geq \alpha \}$, for $\alpha > 0$.

Screening Preprocessing

return $B(\mathbf{u}^{sc}, \mathbf{v}^{sc})$.

and $e^{\mathbf{v}_j^*} = \varepsilon \kappa$ for all $i \in I_{\varepsilon,\kappa}^{\complement}$ and $j \in J_{\varepsilon,\kappa}^{\complement}$ where	
$= \{ j = 1, \dots, m : \boldsymbol{\nu}_j \geq \kappa \varepsilon^2 c_j(\boldsymbol{K}) \}.$	

Screenkhorn Algorithm

_					
	SCREENKHORN($C, \eta, \mu, \nu, n_b, m_b$)				
1.	$K \leftarrow e^{-C/\eta};$)			
2.	$\boldsymbol{\xi} \leftarrow \texttt{sort}(\boldsymbol{\mu} \oslash r(\boldsymbol{\kappa})), \boldsymbol{\zeta} \leftarrow \texttt{sort}(\boldsymbol{\nu} \oslash c(\boldsymbol{\kappa})); //(\texttt{decreasing order})$				
3.	$\varepsilon \leftarrow (\boldsymbol{\xi}_{n_b} \boldsymbol{\zeta}_{m_b})^{1/4}, \kappa \leftarrow \sqrt{\boldsymbol{\zeta}_{m_b} / \boldsymbol{\xi}_{n_b}};$				
4.	$I_{\varepsilon,\kappa} \leftarrow \{i=1,\ldots,n: \mu_i \geq \varepsilon^2 \kappa^{-1} r_i(K)\};$				
5.	$J_{\varepsilon,\kappa} \leftarrow \{j = 1, \ldots, m : \mathbf{\nu}_j \geq \varepsilon^2 \kappa c_j(\mathbf{K})\};$				
6.	$K_{\min} = \min_{I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa}} K_{ij};$	Step 1: Screening			
7.	$\underline{\boldsymbol{\mu}} \leftarrow \min_{i \in \boldsymbol{I}_{\varepsilon,\kappa}} \boldsymbol{\mu}_{i}, \bar{\boldsymbol{\mu}} \leftarrow \max_{i \in \boldsymbol{I}_{\varepsilon,\kappa}} \boldsymbol{\mu}_{i}; \underline{\boldsymbol{\nu}} \leftarrow \min_{j \in \boldsymbol{J}_{\varepsilon,\kappa}} \boldsymbol{\mu}_{i}, \bar{\boldsymbol{\nu}} \leftarrow \max_{j \in \boldsymbol{J}_{\varepsilon,\kappa}} \boldsymbol{\mu}_{i};$				
8.	$\underline{\underline{u}} \leftarrow \log\left(\frac{\varepsilon}{\kappa} \lor \frac{\underline{\mu}}{\varepsilon(m-m_b)+\varepsilon \lor \frac{\overline{\nu}}{n\varepsilon\kappa} K_{\min}}m_b}\right), \overline{\underline{u}} \leftarrow \log\left(\frac{\varepsilon}{\kappa} \lor \frac{\overline{\mu}}{m\varepsilon} K_{\min}\right);$				
9.	$\underline{\underline{v}} \leftarrow \log\left(\varepsilon\kappa \lor \frac{\underline{\underline{v}}}{\varepsilon(n-n_b)+\varepsilon \lor \frac{\kappa\overline{\underline{\mu}}}{m\varepsilon} \frac{n_b}{K_{\min}} n_b}\right), \overline{\underline{v}} \leftarrow \log\left(\varepsilon\kappa \lor \frac{\overline{\underline{v}}}{n\varepsilon} \frac{\overline{\underline{v}}}{K_{\min}}\right);$				
0.	$\bar{\boldsymbol{\theta}} \leftarrow \texttt{stack}(\bar{\boldsymbol{u}} 1_{n_b}, \bar{\boldsymbol{v}} 1_{m_b}), \underline{\boldsymbol{\theta}} \leftarrow \texttt{stack}(\underline{\boldsymbol{u}} 1_{n_b}, \underline{\boldsymbol{v}} 1_{m_b});$	J			
1.	$\mathbf{u}^{(0)} \leftarrow \log(\varepsilon \kappa^{-1}) 1_{n_k}, \mathbf{v}^{(0)} \leftarrow \log(\varepsilon \kappa) 1_{m_k};)$				
2.	$\boldsymbol{\theta}^{(0)} \leftarrow \operatorname{stack}(\boldsymbol{u}^{(0)}, \boldsymbol{v}^{(0)});$				
3.	$\boldsymbol{\theta} \leftarrow \text{L-BFGS-B}(\ \boldsymbol{\theta}^{(0)}, \underline{\boldsymbol{\theta}}, \ \bar{\boldsymbol{\theta}});$ Step 2: L-BFGS-B				
4.	$\boldsymbol{\theta}_{u} \leftarrow (\boldsymbol{\theta}_{1}, \ldots, \boldsymbol{\theta}_{n_{b}})^{\top};$				
5.	$\boldsymbol{\theta}_{v} \leftarrow (\boldsymbol{\theta}_{n_{b}+1}, \ldots, \boldsymbol{\theta}_{n_{b}+m_{b}})^{\top}; $				
ò.	$\boldsymbol{u}_{i}^{\mathrm{sc}} \leftarrow (\boldsymbol{\theta}_{u})_{i} \text{ if } i \in \boldsymbol{I}_{\varepsilon,\kappa} \text{ and } \boldsymbol{u}_{i}^{\mathrm{sc}} \leftarrow \log(\varepsilon \kappa^{-1}) \text{ if } i \in \boldsymbol{I}_{\varepsilon,\kappa}^{C};$				
7.	$\mathbf{v}_{j}^{\mathrm{sc}} \leftarrow (\mathbf{\theta}_{v})_{j} \text{ if } j \in \mathbf{J}_{\varepsilon,\kappa} \text{ and } \mathbf{v}_{j}^{\mathrm{sc}} \leftarrow \log(\varepsilon\kappa) \text{ if } j \in \mathbf{J}_{\varepsilon,\kappa}^{\mathbf{G}};$				

M. Cuturi. Sinkhorn distances: Light speed computation of optimal transport, NIPS 2013. R. Flamary and N. Courty. POT: Python optimal
 Aligned Provide Action Provided Action Provided Action
 Aligned Provide Action Provided Action
 Aligned Provided Action Provided Action Provided Action
 Aligned Provided Action Provided Action
 Aligned Provided Action Provided Action
 Aligned Provided Action
 Al transport library, 2017.

This work was supported by grants from the Normandie Projet GRR-DAISI, European funding FEDER DAISI and OATMIL ANR-17-CE23-0012 Project of the French National Research Agency (ANR).









OT Domain Adaptation (OTDA): MNIST to **USPS**

Acknowledgments



