
Screening Sinkhorn Algorithm via Dual Projections

Anonymous Author(s)

Affiliation

Address

email

Abstract

This paper deals with the problem of approximating optimal transport (OT) distance between two discrete measures. Our proposed approach involves a convex projection of the *dual of Sinkhorn divergence*, allowing to formulate two appropriate active indices sets for the potential variables. These indices sets depend on two parameters acting like a threshold and a scaling factor and they are both directly linked to a priori fixed number budget of points from the supports of the given discrete measures. This new analysis induces a screened version of the dual of Sinkhorn divergence and suggests the *Screenkhorn* algorithm. We illustrate the favorable performance of Screenkhorn in practice with numerical experiments on synthetic and real datasets.

1 Introduction

Computing OT distances between pairs of probability measures or histograms, such as the earth mover’s distance [32, 28] and Monge-Kantorovich or Wasserstein distance [9], are currently generating an increasing attraction in different machine learning tasks [31, 23, 4, 18], statistics [16, 25, 13, 6, 15], and computer vision [8, 28, 30], among other applications [22, 27]. In many of these problems, OT exploits the geometric features of the objects at hand in the underlying spaces to be leveraged in comparing probability measures. This effectively leads to improve performance of methods that are oblivious to the geometry, for example the chi-squared distances or the Kullback-Leibler divergence. Unfortunately, this advantage comes at the price of an enormous computational cost of solving the OT problem, that can be prohibitive in large scale applications. For instance, the OT between two histograms with supports of equal size n can be formulated as a linear programming problem that requires generally super $\mathcal{O}(n^3)$ [26] arithmetic operations, which is problematic when n is larger than 10^3 .

A remedy to the heavy computation burden of OT lies in a prevalent approach referred to as regularized OT [10] and operates by adding an entropic regularization penalty to the original problem. Such a regularization guarantees a unique solution, since the objective function is strongly convex, and a greater computational stability. Furthermore, [10] proposed the so-called dual of Sinkhorn divergence as the dual of the entropic problem and noticed that finding the dual solution was equivalent to finding two diagonal matrices that made a full matrix bistochastic. Therefore, the OT can be solved efficiently with celebrated matrix scaling algorithms, such as Sinkhorn’s fixed point iteration method [29, 21, 19].

Sinkhorn scaling for computing OT distances is a well studied problem in many recent works. The main idea is to improve the matrix-vector operations that are the true computational bottleneck of Sinkhorn’s algorithm. [3] proposed the Greenkhorn algorithm, a greedy version of Sinkhorn algorithm that selects columns and rows to be updated that most violate the constraints. [2] provided the Nys-Sink algorithm which is based on low-rank approximation of the cost matrix using Nystrom method. Other classical optimization algorithms have been considered to approximate the OT, for instance accelerated gradient descent [12, 24], quasi Newton methods [7, 11] and stochastic gradient descent [17, 1].

We give a new algorithm to approximate the regularized OT distance between discrete measures. Our algorithmic analysis is based on an approximate of the dual of Sinkhorn divergence by adding new constraints feasibility. These constraints are defined through a convex set which depends on two parameters, acting like threshold and scaling factor. We prove that dual solution of this approximation guarantees the existence of two active indices sets for the potential variables. These active sets are both directly linked to a priori fixed number budget of points from the supports of the given discrete measures. We then restrict the constraints feasibility with respect to the active sets to get a “screened” version of the dual of Sinkhorn divergence. The Screenkhorn algorithm developed in this paper relies on two steps; the first one consists of an initialization step devoted to determine the active sets while the second is a constrained L-BFGS solver [33, 11, 7].

The remainder of the paper is organized as follow. In Section 2 we briefly review the basic setup of regularized discrete OT. Section 3 contains our main contribution, that is, the Screenkhorn algorithm. Section 4 devotes to theoretical guarantees for the marginal violations of Screenkhorn. In Section 5 we present numerical results for the proposed algorithm, compared with the state-of-art Sinkhorn algorithm as implemented in [14]. The proofs of theoretical results are postponed to the supplementary material.

Notation. For any positive matrix $T \in \mathbb{R}^{n \times m}$, we define its negative entropy as $H(T) = -\sum_{i,j} T_{ij} \log(T_{ij})$. Let $r(T) = T\mathbf{1}_m \in \mathbb{R}^n$ and $c(T) = T^\top \mathbf{1}_n \in \mathbb{R}^m$ denote the rows and columns sums of T respectively. The coordinates $r_i(T)$ and $c_j(T)$ denote the i -th row sum and the j -th column sum of T , respectively. The scalar product between two matrices denotes the usual inner product, that is $\langle T, W \rangle = \text{tr}(T^\top W) = \sum_{i,j} T_{ij} W_{ij}$, where T^\top is the transpose of T . We write $\mathbf{1}$ (resp. $\mathbf{0}$) the vector having all coordinates equal to one (resp. zero). $\Delta(w)$ denotes the diag operator, such that if $w \in \mathbb{R}^n$, then $\Delta(w) = \text{diag}(w_1, \dots, w_n) \in \mathbb{R}^{n \times n}$. For a set of indices $L = \{i_1, \dots, i_k\} \subseteq \{1, \dots, n\}$ satisfying $i_1 < \dots < i_k$, we denote the complementary set of L by $L^c = \{1, \dots, n\} \setminus L$. We also denote $|L|$ the cardinality of L . Given a vector $w \in \mathbb{R}^n$, we denote $w_L = (w_{i_1}, \dots, w_{i_k})^\top \in \mathbb{R}^k$ and its complementary $w_{L^c} \in \mathbb{R}^{n-k}$. The notation is similar for matrices; given another subset of indices $S = \{j_1, \dots, j_l\} \subseteq \{1, \dots, m\}$ with $j_1 < \dots < j_l$, and a matrix $T \in \mathbb{R}^{n \times m}$, we use $T_{(L,S)}$, to denote the submatrix of T , namely the rows and columns of $T_{(L,S)}$ are indexed by L and S respectively. When applied to matrices and vectors, \odot and \oslash (Hadamard product and division) and exponential notations refer to elementwise operators. Given two real numbers a and b , we write $a \vee b = \max(a, b)$ and $a \wedge b = \min(a, b)$.

2 Regularized discrete OT

We briefly present in this section the setup of OT between two discrete measures. We then consider the case when those distributions are only available through a finite number of samples, that is $\mu = \sum_{i=1}^n \mu_i \delta_{x_i} \in \Sigma_n$ and $\nu = \sum_{j=1}^m \nu_j \delta_{y_j} \in \Sigma_m$, where Σ_n is the probability simplex with n bins, namely the set of probability vectors in \mathbb{R}_+^n , i.e., $\Sigma_n = \{w \in \mathbb{R}_+^n : \sum_{i=1}^n w_i = 1\}$. We denote their probabilistic couplings set as $\Pi(\mu, \nu) = \{P \in \mathbb{R}_+^{n \times m}, P\mathbf{1}_m = \mu, P^\top \mathbf{1}_n = \nu\}$.

Sinkhorn divergence. Approximating OT distance between the two discrete measures μ and ν amounts to solving a linear problem [20] given by

$$\mathcal{S}(\mu, \nu) = \min_{P \in \Pi(\mu, \nu)} \langle C, P \rangle, \quad (1)$$

where $P = (P_{ij}) \in \mathbb{R}^{n \times m}$ is called the transportation plan, namely each entry P_{ij} represents the fraction of mass moving from x_i to y_j , and $C = (C_{ij}) \in \mathbb{R}^{n \times m}$ is a cost matrix comprised of nonnegative elements and related to the energy needed to move a probability mass from x_i to y_j . The entropic regularization of OT distances [10] relies on the addition of a penalty term as follows:

$$\mathcal{S}_\eta(\mu, \nu) = \min_{P \in \Pi(\mu, \nu)} \{\langle C, P \rangle - \eta H(P)\}, \quad (2)$$

where $\eta > 0$ is a regularization parameter. We refer to $\mathcal{S}_\eta(\mu, \nu)$ as the *Sinkhorn divergence* [10].

Dual of Sinkhorn divergence. Below we provide the derivation of the dual problem for the regularized OT problem (2). Towards this end, we begin with writing its Lagrangian dual function :

$$\mathcal{L}(P, y, z) = \langle C, P \rangle + \eta \langle \log P, P \rangle + \langle y, P\mathbf{1}_m - \mu \rangle + \langle z, P^\top \mathbf{1}_n - \nu \rangle.$$

85 The dual of Sinkhorn divergence can be derived by solving $\min_{P \in \mathbb{R}_+^{n \times m}} \mathcal{L}(P, y, z)$. It is easy to
 86 check that objective function $P \mapsto \mathcal{L}(P, y, z)$ is strongly convex and differentiable. Hence, one can
 87 solve the latter minimum by setting $\nabla_P \mathcal{L}(P, y, z)$ to $\mathbf{0}_{n \times m}$. Therefore, we get

$$P_{ij}^* = \exp \left(-\frac{1}{\eta} (y_i + z_j + C_{ij}) - 1 \right), \quad (3)$$

88 for all $i = 1, \dots, n$ and $j = 1, \dots, m$. Plugging this solution, and setting the change of variables
 89 $u = -y/\eta - 1/2$ and $v = -z/\eta - 1/2$, the dual problem is given by

$$\min_{u \in \mathbb{R}^n, v \in \mathbb{R}^m} \{ \Psi(u, v) := \mathbf{1}_n^\top B(u, v) \mathbf{1}_m - \langle u, \mu \rangle - \langle v, \nu \rangle \}, \quad (4)$$

90 where $B(u, v) = \Delta(e^u) K \Delta(e^v)$ and $K = e^{-C/\eta}$ stands for the Gibbs kernel associated to the
 91 cost matrix C . We refer to problem (4) to the *dual of Sinkhorn divergence*. Therefore, the optimal
 92 solution P^* of Sinkhorn divergence takes the form $P^* = \Delta(e^{u^*}) K \Delta(e^{v^*})$ where the couple (u^*, v^*)
 93 satisfies:

$$(u^*, v^*) = \underset{u \in \mathbb{R}^n, v \in \mathbb{R}^m}{\operatorname{argmin}} \{ \Psi(u, v) \}.$$

94 Note that the matrices $\Delta(e^{u^*})$ and $\Delta(e^{v^*})$ are unique up to a constant factor [29]. Furthermore, P^*
 95 can be solved efficiently by iterative Bregman projections [5] referred as Sinkhorn iterations, and
 96 the method is referred as Sinkhorn algorithm which, recently, is proven to achieve a near- $\mathcal{O}(n^2)$
 97 complexity.

98 3 Screened dual of Sinkhorn divergence

99 For a fixed $\varepsilon > 0$ and $\kappa > 0$ we define an *approximate dual of Sinkhorn divergence* as follows:

$$\min_{u \in \mathcal{C}_{\frac{\varepsilon}{\kappa}}^n, v \in \mathcal{C}_{\varepsilon \kappa}^m} \{ \Psi_\kappa(u, v) := \mathbf{1}_n^\top B(u, v) \mathbf{1}_m - \langle \kappa u, \mu \rangle - \langle \frac{v}{\kappa}, \nu \rangle \}, \quad (5)$$

100 where $\mathcal{C}_\alpha^r \subseteq \mathbb{R}^r$, for $r \in \mathbb{N}$ and $\alpha > 0$, is a convex set given by $\mathcal{C}_\alpha^r = \{w \in \mathbb{R}^r : \min_{1 \leq i \leq r} e^{w_i} \geq \alpha\}$.

101 The objective function Ψ_κ is convex with respect to (u, v) , then the set of optima of problem (5) is
 102 non empty. The κ -parameter plays a role of scaling factor, namely it allows to get a closed order
 103 of the potential variables e^u and e^v , while the ε -parameter acts like a threshold for e^u and e^v . Note
 104 that the approximate dual of Sinkhorn divergence coincides with the dual of Sinkhorn divergence (4)
 105 in the setting of $\varepsilon = 0$ and $\kappa = 1$. The screening procedure presented in this work is based on
 106 constructing two active sets $I_{\varepsilon, \kappa}$ and $J_{\varepsilon, \kappa}$ throughout the dual problem of (5) in the following way:

107 **Lemma 1.** *Let (u^*, v^*) be an optimal solution of the problem (5). Define*

$$I_{\varepsilon, \kappa} = \{i = 1, \dots, n : \mu_i \geq \frac{\varepsilon^2}{\kappa} r_i(K)\}, J_{\varepsilon, \kappa} = \{j = 1, \dots, m : \nu_j \geq \kappa \varepsilon^2 c_j(K)\} \quad (6)$$

108 *Then one has $e^{u_i^*} = \varepsilon \kappa^{-1}$ and $e^{v_j^*} = \varepsilon \kappa$ for all $i \in I_{\varepsilon, \kappa}^c$ and $j \in J_{\varepsilon, \kappa}^c$.*

109 First order conditions applied to (u^*, v^*) ensure that if $e^{u_i^*} > \varepsilon \kappa^{-1}$ then $e^{u_i^*} (K e^{v^*})_i = \kappa \mu_i$ and if
 110 $e^{v_j^*} > \varepsilon \kappa$ then $e^{v_j^*} (K^\top e^{u^*})_j = \kappa^{-1} \nu_j$ which correspond to the Sinkhorn marginal conditions up to
 111 the scaling factor κ .

112 **Screening with a fixed number budget of points.** Recall that the approximate dual of Sinkhorn
 113 divergence is defined with respect to ε and κ . The explicit determination of its values depends on
 114 a priori *fixed number budget of points* from the supports of μ and ν . In the sequel of the paper, we
 115 denote by $n_b \in \{1, \dots, n\}$ and the $m_b \in \{1, \dots, m\}$ the number budget of points to be given for
 116 resolving problem (5).

117 Let us define $\xi \in \mathbb{R}^n$ and $\zeta \in \mathbb{R}^m$ to be the ordered decreasing vectors of $\mu \odot r(K)$ and $\nu \odot c(K)$
 118 respectively, that is $\xi_1 \geq \xi_2 \geq \dots \geq \xi_n$ and $\zeta_1 \geq \zeta_2 \geq \dots \geq \zeta_m$. To keep only n_b -budget and
 119 m_b -budget of points, the parameters κ and ε satisfy $\varepsilon^2 \kappa^{-1} = \xi_{n_b}$ and $\varepsilon^2 \kappa = \zeta_{m_b}$. Hence

$$\varepsilon = (\xi_{n_b} \zeta_{m_b})^{1/4} \text{ and } \kappa = \sqrt{\frac{\zeta_{m_b}}{\xi_{n_b}}}. \quad (7)$$

120 Note that $|I_{\varepsilon,\kappa}| = n_b$ and $|J_{\varepsilon,\kappa}| = m_b$. Using the previous analysis, any solution (u^*, v^*) of
 121 problem (5) satisfy $e^{u_i^*} \geq \varepsilon\kappa^{-1}$ and $e^{v_j^*} \geq \varepsilon\kappa$ for all $(i, j) \in (I_{\varepsilon,\kappa} \times J_{\varepsilon,\kappa})$, and $e^{u_i^*} = \varepsilon\kappa^{-1}$ and
 122 $e^{v_j^*} = \varepsilon\kappa$ for all $(i, j) \in (I_{\varepsilon,\kappa}^c \times J_{\varepsilon,\kappa}^c)$.

123 Basing on that facts we restrict the constraints feasibility $\mathcal{C}_{\frac{\varepsilon}{\kappa}}^n \cap \mathcal{C}_{\varepsilon\kappa}^m$ in problem (5) to the screened
 124 domain $\mathcal{U}_{\text{sc}} \cap \mathcal{V}_{\text{sc}}$ where

$$\mathcal{U}_{\text{sc}} = \{u \in \mathbb{R}^n : e^{u_{I_{\varepsilon,\kappa}}} \succeq \frac{\varepsilon}{\kappa} \mathbf{1}_{n_b}, \text{ and } e^{u_{I_{\varepsilon,\kappa}^c}} = \frac{\varepsilon}{\kappa} \mathbf{1}_{n-n_b}\},$$

125 and

$$\mathcal{V}_{\text{sc}} = \{v \in \mathbb{R}^m : e^{v_{J_{\varepsilon,\kappa}}} \succeq \varepsilon\kappa \mathbf{1}_{m_b}, \text{ and } e^{v_{J_{\varepsilon,\kappa}^c}} = \varepsilon\kappa \mathbf{1}_{m-m_b}\}.$$

126 where the vector comparison \succeq has to be understood elementwise. Now, we are ready to define the
 127 *screened dual of Sinkhorn divergence* as

$$\min_{u \in \mathcal{U}_{\text{sc}}, v \in \mathcal{V}_{\text{sc}}} \{\Psi_{\varepsilon,\kappa}(u, v)\} \quad (8)$$

128 where

$$\begin{aligned} \Psi_{\varepsilon,\kappa}(u, v) = & (e^{u_{I_{\varepsilon,\kappa}}})^\top K_{(I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa})} e^{v_{J_{\varepsilon,\kappa}}} + \varepsilon\kappa (e^{u_{I_{\varepsilon,\kappa}}})^\top K_{(I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa}^c)} \mathbf{1}_{m_b} + \varepsilon\kappa^{-1} \mathbf{1}_{n_b}^\top K_{(I_{\varepsilon,\kappa}^c, J_{\varepsilon,\kappa})} e^{v_{J_{\varepsilon,\kappa}}} \\ & - \kappa \mu_{I_{\varepsilon,\kappa}}^\top u_{I_{\varepsilon,\kappa}} - \kappa^{-1} \nu_{J_{\varepsilon,\kappa}}^\top v_{J_{\varepsilon,\kappa}} + \Xi \end{aligned}$$

129 with $\Xi = \varepsilon^2 \sum_{i \in I_{\varepsilon,\kappa}^c, j \in J_{\varepsilon,\kappa}^c} K_{ij} - \kappa \log(\varepsilon\kappa^{-1}) \sum_{i \in I_{\varepsilon,\kappa}^c} \mu_i - \kappa^{-1} \log(\varepsilon\kappa) \sum_{j \in J_{\varepsilon,\kappa}^c} \nu_j$.

130 Pseudocode of our proposed algorithm is given in Algorithm 1. Screenkhorn consists of two steps: the
 131 first one is an initialization where we calculate the active sets $I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa}$. The second is a constrained
 132 L-BFGS solver [33] for the stacked variable $\theta = (u_{I_{\varepsilon,\kappa}}, v_{J_{\varepsilon,\kappa}})$. It is worth to note that Screenkhorn
 133 uses only the restricted parts $K_{(I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa})}$, $K_{(I_{\varepsilon,\kappa}, J_{\varepsilon,\kappa}^c)}$, and $K_{(I_{\varepsilon,\kappa}^c, J_{\varepsilon,\kappa})}$ of the Gibbs kernel K , in
 134 contrast to Sinkhorn algorithm which performs alternating updates of all rows and columns of K .

135 The following lemma expresses upper and lower bounds to be respected in Screenkhorn.

136 **Lemma 2.** *Let $(u^{\text{sc}}, v^{\text{sc}})$ be an optimal solution of problem (8). Then, one has*

$$\frac{\varepsilon}{\kappa} \vee \frac{\min_{i \in I_{\varepsilon,\kappa}} \mu_i}{\varepsilon(m - m_b) + \varepsilon \vee \frac{\max_{j \in J_{\varepsilon,\kappa}} \nu_j}{n\varepsilon \min_{i,j} K_{ij}} m_b} \leq e^{u_i^{\text{sc}}} \leq \frac{\varepsilon}{\kappa} \vee \frac{\max_{i \in I_{\varepsilon,\kappa}} \mu_i}{m\varepsilon \min_{i,j} K_{ij}}, \quad (9)$$

137 and

$$\varepsilon\kappa \vee \frac{\min_{j \in J_{\varepsilon,\kappa}} \nu_j}{\varepsilon(n - n_b) + \varepsilon \vee \frac{\kappa \max_{i \in I_{\varepsilon,\kappa}} \mu_i}{m\varepsilon \min_{i,j} K_{ij}} n_b} \leq e^{v_j^{\text{sc}}} \leq \varepsilon\kappa \vee \frac{\max_{j \in J_{\varepsilon,\kappa}} \nu_j}{n\varepsilon \min_{i,j} K_{ij}} \quad (10)$$

138 for all $i \in I_{\varepsilon,\kappa}$ and $j \in J_{\varepsilon,\kappa}$.

139 4 Analysis of marginal violations

140 This section is devoted to study the marginal violations of Screenkhorn. Towards this end, let us
 141 define the screened marginals $\mu^{\text{sc}} = B(u^{\text{sc}}, v^{\text{sc}}) \mathbf{1}_m$ and $\nu^{\text{sc}} = B(u^{\text{sc}}, v^{\text{sc}})^\top \mathbf{1}_n$. Lemma 3 expresses
 142 an upper bound with respect to ℓ_1 -norm of μ^{sc} and ν^{sc} .

143 **Lemma 3.** *Let $(u^{\text{sc}}, v^{\text{sc}})$ be an optimal solution of problem (8). Then one has*

$$\|\mu^{\text{sc}}\|_1 \leq \kappa \|\mu_{I_{\varepsilon,\kappa}}\|_1 + (n - n_b) \left(\frac{m_b \max_{j \in J_{\varepsilon,\kappa}} \nu_j}{n\kappa \min_{i,j} K_{ij}} + (m - m_b) \varepsilon^2 \right)$$

144 and

$$\|\nu^{\text{sc}}\|_1 \leq \kappa^{-1} \|\nu_{J_{\varepsilon,\kappa}}\|_1 + (m - m_b) \left(\frac{n_b \kappa \max_{i \in I_{\varepsilon,\kappa}} \mu_i}{m \min_{i,j} K_{ij}} + (n - n_b) \varepsilon^2 \right).$$

145 The following Proposition gives also an upper bound of the marginal errors.

Algorithm 1: Sinkhorn($C, \eta, \mu, \nu, n_b, m_b$)

step 1: Initialization

1. $\xi \leftarrow \mu \odot r(K)$;
 2. $\zeta \leftarrow \nu \odot c(K)$;
 3. $\xi \leftarrow \text{sort}(\xi)$; //(decreasing order)
 4. $\zeta \leftarrow \text{sort}(\zeta)$; //(decreasing order)
 5. $\varepsilon \leftarrow (\xi_{n_b} \zeta_{m_b})^{1/4}$, $\kappa \leftarrow \sqrt{\zeta_{m_b} / \xi_{n_b}}$;
 6. $I_{\varepsilon, \kappa} \leftarrow \{i = 1, \dots, n : \mu_i \geq \varepsilon^2 \kappa^{-1} r_i(K)\}$;
 7. $J_{\varepsilon, \kappa} \leftarrow \{j = 1, \dots, m : \nu_j \geq \varepsilon^2 \kappa c_j(K)\}$;
 8. $K_{\min} \leftarrow \min_{i \in I_{\varepsilon, \kappa}, j \in J_{\varepsilon, \kappa}} K_{ij}$;
 9. $\underline{\mu} \leftarrow \min_{i \in I_{\varepsilon, \kappa}} \mu_i$, $\bar{\mu} \leftarrow \max_{i \in I_{\varepsilon, \kappa}} \mu_i$;
 10. $\underline{\nu} \leftarrow \min_{j \in J_{\varepsilon, \kappa}} \nu_j$, $\bar{\nu} \leftarrow \max_{j \in J_{\varepsilon, \kappa}} \nu_j$;
 11. $\underline{u} \leftarrow \log\left(\frac{\varepsilon}{\kappa} \vee \frac{\mu}{\varepsilon(m-m_b) + \varepsilon \sqrt{\frac{\bar{\nu}}{n \varepsilon \kappa K_{\min}} m_b}}\right)$, $\bar{u} \leftarrow \log\left(\frac{\varepsilon}{\kappa} \vee \frac{\bar{\mu}}{m \varepsilon K_{\min}}\right)$;
 12. $\underline{v} \leftarrow \log\left(\varepsilon \kappa \vee \frac{\nu}{\varepsilon(n-n_b) + \varepsilon \sqrt{\frac{\kappa \bar{\mu}}{m \varepsilon K_{\min}} n_b}}\right)$, $\bar{v} \leftarrow \log\left(\varepsilon \kappa \vee \frac{\bar{\nu}}{n \varepsilon K_{\min}}\right)$;
 13. $\bar{\theta} \leftarrow \text{stack}(\bar{u} \mathbf{1}_{n_b}, \bar{v} \mathbf{1}_{m_b})$;
 14. $\underline{\theta} \leftarrow \text{stack}(\underline{u} \mathbf{1}_{n_b}, \underline{v} \mathbf{1}_{m_b})$;
 - step 2:** L-BFGS
 15. $u^{(0)} \leftarrow \log(\varepsilon \kappa^{-1}) \mathbf{1}_{n_b}$;
 16. $v^{(0)} \leftarrow \log(\varepsilon \kappa) \mathbf{1}_{m_b}$;
 17. $\theta^{(0)} \leftarrow \text{stack}[u^{(0)}, v^{(0)}]$;
 18. $\theta \leftarrow \text{L-BFGS}(\theta^{(0)}, \underline{\theta}, \bar{\theta})$;
 19. $\theta_u \leftarrow (\theta_1, \dots, \theta_{n_b})^\top$, $\theta_v \leftarrow (\theta_{n_b+1}, \dots, \theta_{n_b+m_b})^\top$;
 20. $u_i^{sc} \leftarrow (\theta_u)_i$ if $i \in I_{\varepsilon, \kappa}$ and $u_i \leftarrow \log(\varepsilon \kappa^{-1})$ if $i \in I_{\varepsilon, \kappa}^c$;
 21. $v_j^{sc} \leftarrow (\theta_v)_j$ if $j \in J_{\varepsilon, \kappa}$ and $v_j \leftarrow \log(\varepsilon \kappa)$ if $j \in J_{\varepsilon, \kappa}^c$;
 22. **return** $B(u^{sc}, v^{sc})$.
-

146 **Proposition 1.** *One has*

$$\begin{aligned} \|\mu - \mu^{sc}\|_1^2 &\leq 7n_b(\kappa - \log(\kappa) - 1) \max_i \mu_i + 7(n - n_b) \left(\frac{m_b \max_j \nu_j}{n \kappa \min_{i,j} K_{ij}} + (m - m_b) \varepsilon^2 - \min_i \mu_i \right) \\ &\quad + \max_i \mu_i \log \left(\frac{\kappa(n - n_b + 1) \max_i \mu_i}{m_b \min_{i,j} K_{ij} \min_{j \in J_{\varepsilon, \kappa}} \nu_j} + \frac{\kappa^2 \max_i \mu_i}{m m_b \varepsilon^2 (\min_{i,j} K_{ij})^2 \min_{j \in J_{\varepsilon, \kappa}} \nu_j} \right) \end{aligned}$$

147 *and*

$$\begin{aligned} \|\nu - \nu^{sc}\|_1^2 &\leq 7m_b(\kappa - \log(\kappa) - 1) \max_j \nu_j + 7(m - m_b) \left(\frac{n_b \kappa \max_i \mu_i}{n \min_{i,j} K_{ij}} + (n - n_b) \varepsilon^2 - \min_j \nu_j \right) \\ &\quad + \max_j \nu_j \log \left(\frac{\kappa(m - m_b + 1) \max_j \nu_j}{n_b \min_{i,j} K_{ij} \min_{i \in I_{\varepsilon, \kappa}} \mu_i} + \frac{\kappa^2 \max_j \nu_j}{n m_b \varepsilon^2 (\min_{i,j} K_{ij})^2 \min_{i \in I_{\varepsilon, \kappa}} \mu_i} \right). \end{aligned}$$

148 5 Numerical experiments

149 References

- 150 [1] B. K. Abid and R. Gower. Stochastic algorithms for entropy-regularized optimal transport problems. In
- 151 Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference*
- 152 *on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages
- 153 1505–1512, Playa Blanca, Lanzarote, Canary Islands, 2018. PMLR.
- 154 [2] J. Altschuler, F. Bach, A. Rudi, and J. Weed. Massively scalable sinkhorn distances via the nyström method,
- 155 2018.
- 156 [3] J. Altschuler, J. Weed, and P. Rigollet. Near-linear time approximation algorithms for optimal transport via
- 157 sinkhorn iteration. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and

- 158 R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 1964–1974. Curran
159 Associates, Inc., 2017.
- 160 [4] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In Doina Precup and
161 Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70
162 of *Proceedings of Machine Learning Research*, pages 214–223, International Convention Centre, Sydney,
163 Australia, 2017. PMLR.
- 164 [5] J. D. Benamou, G. Carlier, M. Cuturi, L. Nenna, and G. Peyré. Iterative bregman projections for regularized
165 transportation problems. *SIAM J. Scientific Computing*, 37, 2015.
- 166 [6] J. Bigot, R. Gouet, T. Klein, and A. López. Geodesic pca in the wasserstein space by convex pca. *Ann.
167 Inst. H. Poincaré Probab. Statist.*, 53(1):1–26, 2017.
- 168 [7] M. Blondel, V. Seguy, and A. Rolet. Smooth and sparse optimal transport. In Amos Storkey and Fernando
169 Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence
170 and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 880–889, Playa Blanca,
171 Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.
- 172 [8] N. Bonneel, M. van de Panne, S. Paris, and W. Heidrich. Displacement interpolation using lagrangian mass
173 transport. *ACM Trans. Graph.*, 30(6):158:1–158:12, 2011.
- 174 [9] Villani C. *Optimal Transport: Old and New*, volume 338 of *Grundlehren der mathematischen Wis-
175 senschaften*. Springer Berlin Heidelberg, 2009.
- 176 [10] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In C. J. C. Burges, L. Bottou,
177 M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing
178 Systems 26*, pages 2292–2300. Curran Associates, Inc., 2013.
- 179 [11] M. Cuturi and G. Peyré. A smoothed dual approach for variational wasserstein problems. *SIAM Journal
180 on Imaging Sciences*, 9(1):320–343, 2016.
- 181 [12] P. Dvurechensky, A. Gasnikov, and A. Kroshnin. Computational optimal transport: Complexity by
182 accelerated gradient descent is better than by sinkhorn’s algorithm. In Jennifer Dy and Andreas Krause,
183 editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings
184 of Machine Learning Research*, pages 1367–1376, Stockholmsmässan, Stockholm Sweden, 10–15 Jul
185 2018. PMLR.
- 186 [13] J. Ebert, V. Spokoiny, and A. Suvorikova. Construction of non-asymptotic confidence sets in 2-wasserstein
187 space, 2017.
- 188 [14] R. Flamary and N. Courty. Pot python optimal transport library, 2017.
- 189 [15] R. Flamary, M. Cuturi, N. Courty, and A. Rakotomamonjy. Wasserstein discriminant analysis. *Machine
190 Learning*, 107(12):1923–1945, 2018.
- 191 [16] C. Frogner, C. Zhang, H. Mobahi, M. Araya, and T. A. Poggio. Learning with a wasserstein loss. In
192 C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural
193 Information Processing Systems 28*, pages 2053–2061. Curran Associates, Inc., 2015.
- 194 [17] A. Genevay, M. Cuturi, G. Peyré, and F. Bach. Stochastic optimization for large-scale optimal transport. In
195 D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information
196 Processing Systems 29*, pages 3440–3448. Curran Associates, Inc., 2016.
- 197 [18] N. Ho, X. L. Nguyen, M. Yurochkin, H. H. Bui, V. Huynh, and D. Phung. Multilevel clustering via
198 wasserstein means. In *Proceedings of the 34th International Conference on Machine Learning - Volume
199 70, ICML’17*, pages 1501–1509. JMLR.org, 2017.
- 200 [19] B. Kalantari, I. Lari, F. Ricca, and B. Simeone. On the complexity of general matrix scaling and entropy
201 minimization via the ras algorithm. *Mathematical Programming*, 112(2):371–401, Apr 2008.
- 202 [20] L. Kantorovich. On the transfer of masses (in russian). *Doklady Akademii Nauk*, 2:227–229, 1942.
- 203 [21] P. Knight. The sinkhorn–knopp algorithm: Convergence and applications. *SIAM Journal on Matrix
204 Analysis and Applications*, 30(1):261–275, 2008.
- 205 [22] S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing
206 and machine-learning applications. *IEEE Signal Processing Magazine*, 34(4):43–59, July 2017.

- 207 [23] M. Kusner, Y. Sun, N. Kolkin, and K. Weinberger. From word embeddings to document distances. In
 208 Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine*
 209 *Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 957–966, Lille, France, 07–09
 210 Jul 2015. PMLR.
- 211 [24] T. Lin, N. Ho, and M. I. Jordan. On efficient optimal transport: An analysis of greedy and accelerated
 212 mirror descent algorithms. *CoRR*, abs/1901.06482, 2019.
- 213 [25] V. M. Panaretos and Y. Zemel. Amplitude and phase variation of point processes. *Ann. Statist.*, 44(2):771–
 214 812, 04 2016.
- 215 [26] O. Pele and M. Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th International*
 216 *Conference on Computer Vision*, pages 460–467, 2009.
- 217 [27] G. Peyré and M. Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*,
 218 11(5-6):355–607, 2019.
- 219 [28] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval.
 220 *International Journal of Computer Vision*, 40(2):99–121, 2000.
- 221 [29] R. Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *The American*
 222 *Mathematical Monthly*, 74(4):402–405, 1967.
- 223 [30] J. Solomon, F. de Goes, G. Peyré, M. Cuturi, A. Butscher, A. Nguyen, T. Du, and L. Guibas. Convolutional
 224 wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Trans. Graph.*,
 225 34(4):66:1–66:11, 2015.
- 226 [31] J. Solomon, R. Rustamov, L. Guibas, and A. Butscher. Wasserstein propagation for semi-supervised
 227 learning. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on*
 228 *Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 306–314, Beijing,
 229 China, 22–24 Jun 2014. PMLR.
- 230 [32] M. Werman, S. Peleg, and A. Rosenfeld. A distance metric for multidimensional histograms. *Computer*
 231 *Vision, Graphics, and Image Processing*, 32(3):328 – 336, 1985.
- 232 [33] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal. Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale
 233 bound-constrained optimization. *ACM Trans. Math. Softw.*, 23(4):550–560, 1997.