

Onlicorne: optimal transportation distances from sample streams

Anonymous Authors¹

Abstract

Optimal Transport (OT) distances are now routinely used as loss functions in ML tasks. Yet, computing OT distances between arbitrary (i.e. not necessarily discrete) probability distributions remains an open problem. This paper introduces a new online estimator of entropy-regularized OT distances between two such arbitrary distributions. It uses streams of samples from both distributions to iteratively enrich a non-parametric representation of the transportation plan. Compared to the classic Sinkhorn algorithm, our method leverages new samples at each iteration, which enables a consistent estimation of the true regularized OT distance. We cast our algorithm as a block-convex mirror descent in the space of positive distributions, which enables a theoretical analysis of its convergence. We numerically illustrate the performance of our method in comparison with concurrent approaches.

1. Introduction

Optimal transport (OT) distances are fundamental in statistical learning, both as a tool for analyzing the convergence of various algorithms, and as a data-dependant term for estimating data density, e.g. using generative models. OT lifts a given distance over data points living in space \mathcal{X} into a distance between probability distributions over the data space $\mathcal{X} \mathcal{P}(\mathcal{X})$; as such, it allows to compare distributions with disjoint support. To alleviate the computational burden of optimal transport, that is cubic in the number of points, it is common to regularize the linear problem that defines it, using an entropic barrier term. This approach, that has been rediscovered many times in the previous thirty years, allows to approximate OT distances using a matrix balancing algorithm, amenable to GPU computations.

The Sinkhorn algorithm was introduced in a discrete setting,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

i.e. when both distributions to compare are a set of realizations. The so-called Sinkhorn distances between empirical distributions indeed form an estimate of the OT distance between the true distributions from which the samples are drawn. This approach estimates the OT distance in two distinct phases: we draw samples and evaluate a pairwise distance matrix in the first phase; we balance this distance matrix using Sinkhorn-Knopp iterations in the second phase, thereby obtaining a transportation plan and distance.

In this paper, we show how mingling together these two phases can be beneficial to quickly estimate OT distances. Our approach relies on three observations. First, Sinkhorn iterations can be rewritten as a block convex mirror descent on the space of positive distributions. This formulation is valid in the discrete and continuous setting. Second, we can modify these iterations to rely on realizations $\hat{\alpha}_t, \hat{\beta}_t$ of the two distributions α and β , renewed at each iteration t . Finally, we can represent the iterates produced by such approximations in a space of mixtures of simple functions. Those iterates are a simple transformation of the potentials in the Sinkhorn optimization problem.

Contribution. These observations allows us to propose the following material.

- We introduce a new *online Sinkhorn* algorithm. It produces a sequence of estimates $(\hat{w}_t)_t \in \mathbb{R}$ and of transportation plans $\hat{\pi}_t \in \mathcal{P}(\mathcal{X} \times \mathcal{X})$, using two stream of renewed samples $\hat{\alpha}_t = \sum_{i=1}^n \delta_{x_i^t}, \hat{\beta}_t = \sum_{i=1}^n \delta_{y_i^t}$, where x_i^t and y_i^t are sampled from α and β .
- We show that those estimations are consistent, in the sense that $\hat{w}_t \rightarrow \mathcal{W}_{C,\varepsilon}(\alpha, \beta)$, and $\hat{\pi}_t \rightarrow \pi^*$.
- We empirically demonstrate that our algorithm permits a faster estimation of optimal transportation distances for discrete distributions, and a convincing estimation of OT distances between *continuous* distributions.

2. Background: optimal transport distances

We recall the definition of optimal transport distances between arbitrary distributions (i.e. not necessarily discrete), then review how these are estimated using finite samples.

2.1. Optimal transport distances and algorithms

Wasserstein distances. We consider a complete space \mathcal{X} , equipped with a distance function $C : \mathcal{X} \rightarrow \mathbb{R}$. Optimal transport lifts this *ground metric* distance into a distance between probability distributions over the space \mathcal{X} . The Wasserstein distance between α and β is defined as the minimal cost required to move each element of mass of α to each element of mass of β . It rewrites as the solution of a linear problem (LP) over the set of transportation plans

$$\mathcal{W}_C(\alpha, \beta) = \min_{\pi \in \mathbb{P}(\mathcal{X} \times \mathcal{X})} \langle C, \pi \rangle, \quad (1)$$

where $\pi_1 = \int_{y \in \mathcal{X}} d\pi(\cdot, y)$ and $\pi_2 = \int_{x \in \mathcal{X}} d\pi(x, \cdot)$ are the first and second marginals of the transportation plan π . It can shown (ref) that \mathcal{W}_C is a distance over $\mathcal{P}(\mathcal{X})$, that measures the convergence in law.

Entropic regularization and Sinkhorn algorithm. The solution of (2) can be approximated by a simpler optimisation problem, where an entropic term is added to the linear objective to force curvature. The so-called Sinkhorn distance

$$\mathcal{W}_{C,\varepsilon}(\alpha, \beta) = \min_{\pi \in \mathbb{P}(\mathcal{X} \times \mathcal{X})} \langle C, \pi \rangle + \varepsilon \text{KL}(\pi | \alpha \otimes \beta), \quad (2)$$

is indeed an ε -approximation of $\mathcal{W}_C(\alpha, \beta)$. The later problem admits a dual form, which is a maximization problem in the space of continuous *potential* function:

$$\max_{f, g \in \mathcal{C}(\mathcal{X})} \langle f, \alpha \rangle + \langle g, \beta \rangle + \varepsilon \left(\langle \alpha \otimes \beta, \exp\left(\frac{f \oplus g - C}{\varepsilon}\right) \rangle - 1 \right). \quad (3)$$

The dual problem (3), a regularized version of the dual of (2), can be solved by alternated maximization, performing at iteration t

$$f_{t+1}(\cdot) = -T_{C,\varepsilon}(g_t, \beta) \quad g_{t+1}(\cdot) = -T_{C^\top, \varepsilon}(f_{t+1}, \alpha), \quad (4)$$

$$\text{where } T_C(h, \mu) \triangleq \int_{y \in \mathcal{X}} \exp\left(\frac{h(y) - C(\cdot, y)}{\varepsilon}\right) d\mu,$$

is a *soft C-transform*, and the notation $f_t(\cdot)$ emphasizes the belonging of f_t and g_t to the space of continuous functions. $(f_t)_t$ and $(g_t)_t$ converge in $(\mathcal{C}(\mathcal{X}), \|\cdot\|_{\text{var}})$ to a solution (f^*, g^*) of (3), where $\|f\|_{\text{var}} = \max_x f(x) - \min_x f(x)$ is the so-called variation norm. Convergence is due to the strict contraction of the operators $T_C(\cdot, \beta)$ and $T_{C^\top}(\cdot, \alpha)$ in the space $(\mathcal{C}(\mathcal{X}), \|\cdot\|_{\text{var}})$.

2.2. Estimating OT distances with realizations

Regularized optimal transport is elegantly written in functional spaces but the iterations (4) transfers into code

only for discrete distributions $\hat{\alpha} = \sum_{i=1}^n a_i \delta_{x_i}$ and $\hat{\beta} = \sum_{i=1}^n a_i \delta_{x_i}$. In this case, they correspond to the well-known Sinkhorn-Knopp algorithm for balancing the matrix $\exp(\frac{-C}{\varepsilon})$. The algorithm is run in two phases: a first one, during which the cost matrix is computed (with a cost in $\mathcal{O}(n^2 d)$), and a second one, during which it is balanced (each iteration have a cost in $\mathcal{O}(n^2)$).

Sample complexity results. Fortunately, the OT and Sinkhorn distances between two arbitrary distributions α and β can be approximated by the distance between discrete realizations $\hat{\alpha}_n = \frac{1}{n} \sum_i \delta_{x_i}$, $\hat{\beta}_n = \frac{1}{n} \sum_i \delta_{y_i}$, where $(x_i)_i$ and $(y_i)_i$ are i.i.d samples from α and β . Consistency holds, as $\mathcal{W}_{C,(\varepsilon)}(\hat{\alpha}_n, \hat{\beta}_n) \rightarrow \mathcal{W}_{C,(\varepsilon)}(\alpha, \beta)$, with a convergence rate in $\mathcal{O}(n^{-1/2})$ for Sinkhorn distances and $\mathcal{O}(n^{-1/d})$ for Wasserstein distances.

Bias in distance estimation. Although consistency is a reassuring result, the sample complexity of transport in high dimensions with low regularization remains high. For computational reasons, we cannot choose n to be much more than 10^5 , which is not sufficient to ensure that $\mathcal{W}_{C,\varepsilon}(\alpha, \beta)$ is ε -close to \mathcal{W}_C in the typical case where $d = ?$ and $\varepsilon = ?$.

We may wonder whether we can improve the estimation of $\mathcal{W}_C(\alpha, \beta)$ using several sets of samples $(\hat{\alpha}_n^t)_t$ and $(\hat{\beta}_n^t)_t$. Those should be of reasonable size to allow Sinkhorn estimation, and may for example come from a temporal stream. () have proposed to use the Monte-Carlo estimate $\hat{\mathcal{W}}(\alpha, \beta) = \frac{1}{T} \sum_{t=1}^T \mathcal{W}(\hat{\alpha}_n^t, \hat{\beta}_n^t)$. However, this yields a wrong estimation as the distance $\mathcal{W}(\hat{\alpha}_n, \hat{\beta}_n)$ between discrete realizations is a *biased* estimator of $\mathcal{W}(\alpha, \beta)$:

$$\mathcal{W}(\alpha, \beta) \neq \mathbb{E}_{\hat{\alpha}_n \sim \alpha, \hat{\beta}_n \sim \beta} [\mathcal{W}(\hat{\alpha}_n, \hat{\beta}_n)].$$

Bias in gradients. In several applications, the distance $\mathcal{W}(\alpha, \beta)$ is used as a loss function. This is the case in generative modeling, when we parametrize α as the push-forward of some noise distribution μ through a neural network g_θ . We are then interested in computing the displacement gradient $\delta_\alpha \mathcal{W}(\alpha, \beta) \in \mathcal{P}(\mathcal{X})$, in order to train θ by backpropagation. This gradient turns out to be the spatial derivative $\nabla_x f^*$ of the solution of (3). Yet, similarly, estimating this gradient through sampling is biased, as $f^*(\alpha, \beta) \neq \mathbb{E}_{\hat{\alpha}_n \sim \alpha, \hat{\beta}_n \sim \beta} [f^*(\hat{\alpha}_n, \hat{\beta}_n)]$.

3. OT distances from sample streams

We introduce a novel understanding of the Sinkhorn algorithm in this section, whence we derive average and online adaptations. We wish to construct an estimator of $\mathcal{W}(\alpha, \beta)$ from multiple sets of samples $(\hat{\alpha}_n^t)_t$ and $(\hat{\beta}_n^t)_t$. This estimator should successively use these samples to enrich a

representation of the solution of (3), that may be arbitrary complex. $(\hat{\alpha}_n^t)_t$ and $(\hat{\beta}_n^t)_t$ may be seen as mini-batches within a training procedure, or as a temporal stream.

References

3.1. Unbiased Sinkhorn iterations

3.2. Offline mini-batch averaging

3.3. Online estimation of Sinkhorn distances

3.4. Bias-variance trade-offs

4. Analysis

4.1. Offline Sinkhorn averaging

Estimator $\langle \alpha, -\log \mathbb{E}[\exp(-\hat{f})] \rangle + \langle \beta, -\log \mathbb{E}[\exp(-\hat{g})] \rangle$

- estimator properties

4.2. Online Sinkhorn convergence

- slowed down Sinkhorn convergence

- Random iterated functions

- Combining both

4.3. Non-convex mirror descent

5. Experiments

5.1. Offline distance averaging

5.2. Online distance computations

5.3. Training generative models