

Computational Optimal Transport

Professor Xiangfeng Wang

Institute of Trustworthy Artificial Intelligence

SH Research Institute for Intelligent Autonomous Systems

Typed by Yuelin Li

November 2020

1 Introduction to Optimal Transport

Pure math \rightarrow Applied math \rightarrow Computational math \rightarrow Machine learning

Optimal transport (OT) is firstly introduced by Monge two centuries ago: transport or reshape the first pile of sand so that it matches the second pile of sand with least cost. **The goal** is defining geometric tools that are useful to compare probability distributions.

2 Preliminaries of Optimal Transport

For probability vector \mathbf{a} ,

$$\mathbf{a} \in \Sigma_n := \left\{ \mathbf{a} \in \mathbb{R}_+^n \mid \sum_{i=1}^n a_i = 1 \right\}$$

Discrete measures: a discrete measure with weights \mathbf{a} and locations $x_1, \dots, x_n \in \mathcal{X}$ reads $\alpha = \sum_{i=1}^n a_i \delta_{x_i}$, where δ_x is the Dirac at position x , intuitively a unit of mass which is infinitely concentrated at location x . For example, given a cost matrix $(C_{i,j})_{i \in [n], j \in [m]}$ (let $n = m$).

2.1 Optimal Assignment Problem

Suppose an optimal assignment problem: $\min_{\sigma \in \text{perm}(n)} \frac{1}{n} \sum_{i=1}^n C_{i, \sigma(i)}$. This problem may have several optimal solutions.

For discrete measures:

$$\alpha = \sum_{i=1}^n a_i \delta_{x_i}, \quad \beta = \sum_{j=1}^m b_j \delta_{y_j}$$

Monge problem: seeks for a map $T : [n] \rightarrow [m]$

$$\forall j \in [m], \quad b_j = \sum_{i: T(x_i)=y_j} a_i$$

$T\alpha = \beta$. This map should minimize some transportation cost $c(x, y)$ with $(x, y) \in \mathcal{X} \times \mathcal{Y}$

$$\min_T \left\{ \sum_i c(x_i, T(x_i)) \mid T\alpha = \beta \right\}$$

Discrete point case: $\sigma : [n] \rightarrow [m]$ with $j = \sigma(i)$

$$\sum_{i \in \sigma^{-1}(j)} a_i = b_j$$

$n = m$ and $a_i = b_j = \frac{1}{n}$ case $\implies T(x_i) = y_{\sigma(i)}$.

2.2 Kantorovich Relaxation

Background: limitations of the Monge problem are

- (1) Compare two points clouds of the same size
- (2) Combinatorial or non-convex

The Key idea is relaxing the deterministic nature of transportation. The mass at any point x_i be potentially dispatched across several locations.

σ or $T \rightarrow$ a coupling matrix $P \in \mathbb{R}_+^{n \times m}$.

Coupling:

$$U(a, b) := \{P \in \mathbb{R}_+^{n \times m} \mid P\mathbf{1}_m = a, P^T\mathbf{1}_n = b\}$$

Matrix-vector notation:

$$P\mathbf{1}_m = \left(\sum_j P_{i,j} \right) \in \mathbb{R}^n, P^T \mathbf{1}_n = \left(\sum_i P_{i,j} \right) \in \mathbb{R}^m$$

The set is bounded, defined by $n + m$ equality constraints, and therefore a convex polytope.

Kantorovich's optimal transport problem:

$$L_C(a, b) = \min_{P \in U(a, b)} \left\{ \langle C, P \rangle = \sum_{i,j} C_{i,j} P_{i,j} \right\}$$

A linear programming problem Its solutions are not necessarily unique.

This relaxation is better in that 1. relatively easy to calculate and 2. the Kantorovich relaxation is tight when considered on assignment problems.

2.3 Matric Properties of Optimal Transport

Suppose $n = m$, and that for some $p \geq 1, C = D^p$.

$D^p = (D_{i,j}^p)_{i,j} \in \mathbb{R}^{n \times n}$, $D \in \mathbb{R}_+^{n \times n}$ is a distance on $[n]$, note that D is symmetric, $D_{i,j} = 0$ if and only if $i = j$, $\forall (i, j, k) \in [n]^3$, $D_{i,k} \leq D_{i,j} + D_{j,k}$.

Define the p -Wasserstein distance on Σ_n $W_p(a, b) = L_{D^p}(a, b)^{\frac{1}{p}}$.

Wasserstein Barycenter Given a set of distributions $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$, their Wasserstein barycenter is defined as

$$q^*(\mathcal{P}, \lambda) = \arg \min_{q \in \mathcal{Q}} \sum_{k=1}^K \lambda_k W(q, p_k)$$

where W denotes the Wasserstein distance, \mathcal{Q} is in the space of probability distributions, and $\sum_{k=1}^K \lambda_k = 1$.

3 Applications

Applications include Machine Learning, Computer Vision, Robust Optimization, etc.

Simple model of OT: The label distribution learning is an extension of multi-label learning, which cares more about the relative importance of difference. The labels in description of an instance. Suppose we have a training dataset: $\{(x_i, y_i)\}_{i=1}^m$

$$\begin{aligned} \min_{K,h} \quad & \sum_{i=1}^m \langle P_i, M \rangle + \frac{\mu}{2} \|K - K_0\|_F^2 \\ \text{s.t.} \quad & P_i \in \cup(h(x_i), y_i), K \in \mathcal{S}_+ \end{aligned}$$

where $M_{ij} = K_{ii} - 2K_{ij} + K_{jj}$.

Wasserstein Generative Adversarial Networks (WGAN), suppose we have a dataset $\{x_i\}$ and a source of noise data $\{z_j\}$, then we find a parameterized function $g_\theta(\cdot)$ that

$$\begin{aligned} \min_{\theta} \quad & W(\{x_i\}, \{g(z_j, \theta)\}) \\ W(\{x_i\}, \{g_\theta(z_j)\}) = \min_P \quad & \langle C(\theta), P \rangle \\ \text{s.t.} \quad & P \mathbf{1}_n = \frac{1}{n} \mathbf{1}_n, P^T \mathbf{1}_n = \frac{1}{n} \mathbf{1}_n \end{aligned}$$

where g is parameterized by a neural network with θ .

Zero-Shot Learning: The key issue is how to effectively transfer the model (information) learned from the seen classes to the unseen classes. It has a heterogeneous graph with hierarchical structure. The representative nodes are connections between classes, and each instance follows a probability distribution from feature space. Wasserstein barycenter is useful in this case.

4 How to Calculate Optimal Transport

The classic method: **Entropic Regularization of Optimal Transport**, which does not solve the linear programming but adds an entropic regularization on the target function.

$$L_C^\epsilon(a, b) = \min_{P \in \mathcal{U}(a, b)} \langle P, C \rangle - \epsilon H(P) \quad (1)$$

The discrete entropy of a coupling matrix: (ϵ -strongly convex function)

$$H(P) = - \sum_{i,j} P_{i,j} (\log(P_{i,j}) - 1)$$

The above problem has a unique optimal solution. The optimal P_ϵ progressively moves toward an entropic center of the triangle.

The unique solution P_ϵ converges to the optimal solution with maximal entropy within the set of all optimal solutions of the Kantorovich problem, namely

$$P_\epsilon \xrightarrow{\epsilon \rightarrow 0} \arg \min_P \{-H(P) \mid P \in U(a, b), \langle P, C \rangle = L_C(a, b)\}$$

so that in particular $L_C^\epsilon(a, b) \xrightarrow{\epsilon \rightarrow 0} L_C(a, b)$. One has $P_\epsilon \xrightarrow{\epsilon \rightarrow \infty} ab^T = (a_i b_j)_{i,j}$.
Sinkhorn's Algorithm and its Convergence Sinkhorn's algorithm was originally introduced with a proof of convergence by Sinkhorn (1960's). It provides a scalable way to approximate optimal transport by seamless parallelization when solving several OT problems simultaneously.

Kullback-Leibler divergence:

$$KL(P \mid K) = \sum_{i,j} P_{i,j} \log \left(\frac{P_{i,j}}{K_{i,j}} \right) - P_{i,j} + K_{i,j} \quad (2)$$

The unique solution P_ϵ is a projection onto $U(a, b)$ of the Gibbs kernel associated to the cost matrix C

$$K_{i,j} = \exp \left(-\frac{C_{i,j}}{\epsilon} \right) \\ P_\epsilon = \text{Proj}_{U(a,b)}^{KL}(K) = \arg \min_{P \in U(a,b)} KL(P \mid K)$$

The solution P_ϵ is unique and has the form

$$\forall (i, j) \in [n] \times [m], \quad P_{i,j} = u_i K_{i,j} v_j$$

for two scaling variable $(u, v) \in \mathbb{R}_+^n \times \mathbb{R}_+^m$. $P = \text{diag}(u)K\text{diag}(v)$, $\text{diag}(u)K\text{diag}(v)\mathbf{1}_m = a$, and $\text{diag}(v)K^T\text{diag}(u)\mathbf{1}_n = b$, $u \odot (Kv) = a$, and $v \odot (K^T u) = b$.

sinkhorn's algorithm: $v^{(0)} = \mathbf{1}_m$

$$\mathbf{u}^{(\ell+1)} = \frac{\mathbf{a}}{\mathbf{K}_v(\ell)}, \quad \mathbf{v}^{(\ell+1)} = \frac{\mathbf{b}}{\mathbf{K}^T \mathbf{u}^{(\ell+1)}}$$

However, the main computational bottleneck is the vector-matrix multiplication against kernels K and K^T

Compute simultaneously several regularized Wasserstein distances between

pairs of histograms $a_1, \dots, a_N \in \Sigma_n$ and $b_1, \dots, b_N \in \Sigma_m$,
 $L_C^\epsilon(a_1, b_1), \dots, L_C^\epsilon(a_N, b_N)$, $A = [a_1, \dots, a_N] \in \mathbb{R}^{n \times N}$, and $B = [b_1, \dots, b_N] \in \mathbb{R}^{m \times N}$

$$U^{\ell+1} = \frac{A}{KV^{(\ell)}}, \quad V^{(\ell+1)} = \frac{B}{K^T U^{(\ell+1)}}$$

The problem becomes,

$$L_C^\epsilon(a, b) = \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle - \epsilon \left\langle e^{\frac{f}{\epsilon}}, K e^{\frac{g}{\epsilon}} \right\rangle$$

The optimal (f, g) are linked to scaling (u, v) through $(u, v) = \left(e^{\frac{f}{\epsilon}}, e^{\frac{g}{\epsilon}} \right)$, Sinkhorn as a Block Coordinate Ascent on the dual problem, alternating solving f and g , $(f^\ell, g^\ell) = \epsilon \left(\log(u^\ell), \log(v^\ell) \right)$.

Other efficient algorithms for OT include, Greenkhorn, Conjugate gradient method, Primal-dual methods, Stochastic gradient descent methods, ADMM type method... [1–3]

5 Case Study

Inexact Proximal Point Algorithm

$$L_C(a, b) = \min_{P \in U(a, b)} \langle C, P \rangle$$

$$L_C^\epsilon(a, b) = \min_{P \in U(a, b)} \langle P, C \rangle - \epsilon H(P)$$

It approximates original optimal transport (but not). The performance both in terms of numerical stability and computational complexity is very sensitive to the choice of ϵ .

Iterative scheme

$$P^{(\ell+1)} = \arg \min_{P \in U(a, b)} \langle P, C \rangle + \gamma D_h(P, P^{(\ell)})$$

Bregman divergence

$$D_h(P, P^{(\ell)}) = \sum_{i,j} P_{i,j} \left(\log \frac{P_{i,j}}{P_{i,j}^{(\ell)}} \right) - \sum_{i,j} P_{i,j} + \sum_{i,j} P_{i,j}^{(\ell)}$$

$$\implies P^{(\ell+1)} = \arg \min_{P \in U(a, b)} \langle P, C - \gamma \log(P^{(\ell)}) \rangle - \gamma H(P) .$$

The major issue is we have not proved the convergence performance of it when implementing one Sinkhorn iteration.

...

References

- [1] M. Cuturi and G. Peyre, *Semi-dual regularized optimal transport*, SIAM Imaging (2018).
- [2] Marco Cuturi, *Sinkhorn distance: Lightspeed computation of optimal transport*, NeurIPS (2013).
- [3] J. Niles-Weed J. Altschuler and P. Rigollet, *Near-linear time approximation algorithms for optimal transport via sinkhorn iteration*, (NIPS) (2017).