Entropic Optimal Transport On Sinkhorn algorithm and the link with Schröfinger problem

Luca Nenna (LMO, Université Paris-Saclay), 24/01/2024

Monge-Kantorovich problem and some numerics

$$\mathscr{K}_{c}(\mu,\nu) := \inf \left\{ \int_{X \times Y} c(x) \right\}$$

Three main ways to solve numerically this problem:

- Discrete to Discrete: X and Y are finite set and the measures are supported on diracs;
- 2) Discrete to continuous: one of the measure is a.c. With respect to Lebesgue (see Quentin Mérigot's works);
- 3) Continuous to continuous: both the measures are ac with respect to Lebesgue: the celebrated Benamou-Bernier formulation of Optimal transport.

 $(x, y)d\gamma(x, y) \mid \gamma \in \Pi(\mu, \nu)$



Some remarks:

- power p distance);
- Continuous OT works if the Monge-Kantorovich problem admits a dynamic formulation;
- quadratic case (Benamou, Mirebeau, Froese, Oberman, Duval)

Today we will main focus on discrete optimal and entropic regularization to solve it

- Semi-discrete OT works with quadratic cost (or some close variant such as some

- Other approaches exist: e.g. column generation (Friesecke and Penka), moment constraints relaxation (Virginie, Alfonsi et al), solving Monge-Ampére equation for the



Assume X and Y finite sets (with both cardinality N),

$$\mu = \sum_{x \in X} \mu_x \delta_x$$

Then the problem is formulated as follows

$$\inf \left\{ \sum_{x \in X, y \in Y} c(.$$

 N^2 unknonws and 2N constraints to verify $O(N^3)$ complexity via standard linear programming

 $\delta_x \text{ and } \nu = \sum_{y \in Y} \nu_y \delta_y$

 $x(x, y)\gamma_{x,y} \mid \gamma \in \Pi(\mu, \nu)$

Entropic Optimal Transport: the discrete case

<u>Main idea:</u> penalize the non-negativity of $\gamma_{x,v} \ge 0$ by means of an entropy term Ent(γ) := $\sum e(\gamma_{x,y})$ where X,Y $e(r) = \begin{cases} r(\log r - 1) & \text{if } r > 0\\ 0 & \text{if } r = 0\\ +\infty & \text{if } r < 0 \end{cases}$

The regularized problem takes the following form

$$\mathscr{K}_{c}^{\varepsilon}(\mu,\nu) := \inf \left\{ \left\langle c,\gamma \right\rangle + \varepsilon \operatorname{Ent}(\gamma) \mid \sum_{y} \gamma_{x,y} = \mu_{x}, \ \sum_{x} \gamma_{x,y} = \nu_{y} \right\}$$

Where $\left\langle c,\gamma \right\rangle = \sum_{x,y} c(x,y)\gamma_{x,y}$.

1st good news

Thm: problem $\mathcal{K}_c^{\varepsilon}(\mu,\nu)$ has a unique solution γ_{ε} which belongs to $\Pi(\mu,\nu)$.

2nd good news

Thm[convergence in ε]: Consider the sequence of unique solutions γ_{ε} , then it solution of $\mathscr{K}_{c}(\mu,\nu)$ that is

 $\gamma_{\varepsilon} \to \operatorname{argmin}\{\operatorname{Ent}(\gamma) \mid \gamma \in \Pi(\mu, \nu), \langle c, \gamma \rangle = \mathscr{K}_{c}(\mu, \nu)\}.$

converges to the optimal solution with the minimal entropy within the set of all optimal





Sketch

- TAKE \mathcal{E}_{k} S.t. $\mathcal{E}_{h} \rightarrow 9$ Soution With EEn
- TAKE & OPTIMAL FOR K. (4,v)

AND DENOTE YR THE

. T(KI) IS POUNDED AND CLOSE ZYKn-> X EU (MI)

$0 \leq \langle \mathcal{X}_{n}, \mathcal{C} \rangle - \langle \mathcal{X}, \mathcal{C} \rangle \leq \mathcal{E}_{n} \left(\operatorname{Ent}(n) - \operatorname{Ent}(\mathcal{M}) \right)$





The effect of the regularization



Marginals μ and ν



Support of the optimal γ_{ε} as $\varepsilon \to 0$

The matching problem and the regularized counterpart



problem, that is

$$\mathscr{L}(\gamma, \varphi, \psi) := \langle c, \gamma \rangle + \varepsilon \operatorname{Ent}(\gamma) +$$

Where, as in the unregularized case, φ and ψ are the Lagrange multipliers. Then,

$$\mathscr{K}_c^{\varepsilon}(\mu,\nu) = \inf_{\gamma}$$

Deriving the dual problem. First, consider the Lagrangian associated to the entropic



f sup $\mathscr{L}(\gamma, \varphi, \psi)$ ϕ, ψ

Rmk: by KKT (optimality) conditions $\partial_{\gamma} \mathscr{L}$ between primal and dual variables

$$c(x, y) - \varphi(x) - \psi(y) + \varepsilon \log(\gamma_{x, y}) = 0$$

Which gives the following form of the optimal γ

$$\gamma_{x,y} = \exp\left(\frac{\varphi(y) + z}{1-z}\right)$$

Now, by interchanging inf and sup, as we did last week, we obtain...

$$\ell(\gamma, \varphi, \psi) = 0$$
, we have the following relat

 $\frac{\psi(y) - c(x, y)}{\varepsilon}$



 $\mathscr{D}_{c}^{\varepsilon}(\mu,\nu) := \sup \Phi_{\varepsilon}(\varphi,\psi),$ φ, ψ

Where

$\Phi_{\varepsilon}(\varphi, \psi) := \sum \varphi(x)\mu_x + \sum \psi(y)\nu_y$ ${\mathcal X}$ y

Thm: strong duality holds $\mathscr{K}_{c}^{\varepsilon}(\mu,\nu) = \mathscr{D}_{c}^{\varepsilon}(\mu,\nu)$

$$f_y - \varepsilon \sum_{x,y} \exp\left(\frac{\phi(x) + \psi(y) - c(x,y)}{\varepsilon}\right)$$

Rmk 1: the optimal coupling can be written as

Where $D_{\varphi} = \operatorname{diag}(\varphi/\varepsilon)$ and $D_{\psi} = \operatorname{diag}(\psi/\varepsilon)$ are diagonal matrices and $K \in \mathbb{R}^N \times \mathbb{R}^N$ is such that $K_{x,y} = \exp\left(\frac{-c(x,y)}{\varepsilon}\right)$. This actually makes our problem very similar to a matrix scaling problem.

Def (matrix scaling problem): Given a matrix K with positive coefficients find (D_{φ}, D_{ψ}) such that $D_{\varphi}KD_{\psi}$ is doubly stochastic.

Rmk 2: Notice that if (D_{φ}, D_{ψ}) is a solut

 $\gamma = D_{\varphi} K D_{\psi}$

tion then
$$(cD_{\varphi}, \frac{1}{c}D_{\psi})$$
 for any c .

Sinkhorn algorithm

Algorithm 1 Sinkhorn-Knopp algorithm for f
1: function SINKHORN-KNOPP (K)
2: $D^0_{\varphi} \leftarrow 1_N, \ D^0_{\psi} \leftarrow 1_N$
3: for $0 \leq k < k_{\max} do$
4: $D_{\varphi}^{k+1} \leftarrow 1_N./(KD_{\psi}^k)$
5: $D_{\psi}^{k+1} \leftarrow 1_N . / (K^T D_{\varphi}^{k+1})$
6: end for
7: end function

Algorithm 2 Sinkhorn-Knopp algorithm for the regularised optimal transport problem 1: function SINKHORN-KNOPP $(K_{\varepsilon}, \mu, \nu)$ $D^0_{\varphi} \leftarrow \mathbf{1}_X, \ D^0_{\psi} \leftarrow \mathbf{1}_Y$ 2: for $0 \leq k < k_{\max} \operatorname{do}$ 3: $D_{\varphi}^{k+1} \leftarrow \mu./(KD_{\psi}^k)$ 4: $D_{\psi}^{k+1} \leftarrow \nu./(K^T D_{\varphi}^{k+1})$ 5:end for 6: 7: end function

Rmk: $\mathscr{K}_{c}^{\varepsilon}$ can be recasted as a matrix scaling problem by taking

the matrix scaling problem

The importance of being sparse: a multi-scale approach.

In order to reduce the number of grid points used one can apply a multiscale approach and refine the mesh where the solution is supported



Figure: support of the optimal γ_{e} for the Coulomb cost.



log-domain and obtain an iterative methods acting on the dual variable, that is

$$\varphi^k = \varepsilon \log$$

$$\psi^k = \varepsilon \log(k)$$

It turns out that $-\epsilon \log(KD_w^k)$ (resp. $-\epsilon \log(KD_w^k)$) is the soft c-transform of ψ (resp. φ).

In particular the relations above still hold for the optimal dual variables and we have

 $\varphi^* = \varepsilon \log(\mu) - \varepsilon \log(KD_{\psi}^*) \to \min c(x, y) - \psi(y) \text{ as } \varepsilon \to 0.$

Consider the entropic with a kernel K. Then we can re-write the Sinkhorn iterations on

 $g(\mu) - \varepsilon \log(KD_w^k),$

 $g(\nu) - \varepsilon \log(KD_{\omega}^k).$





Convergence of Sinkhorn by using the Hilbert metric

$$\forall (u, v) \in (R^n_{+,\star})^2, d_H(u, v)$$

Where $||x||_V = \max_i x_i - \min_i x_i$. **Theorem A.2** ([2, 12]). Let $K \in \mathbb{R}^{n \times n}_{+,\star}$, then for $d_H(Ku, Kv) \leq \lambda(Ku)$

where

 $\lambda(K) = \frac{\sqrt{\eta(K)} - 1}{\sqrt{n(K)} + 1} < 1$

and

$$\eta(K) = \max_{i,j,kl} \frac{K_{ik}K_{jl}}{K_{jk}K_{il}}.$$

Def (Hilbert projective metric): the Hilbert projective metric on $\mathbb{R}^n_{+,\star}$ is defined as $P) := ||\log(u) - \log(v)||_{V}$

$$r(u,v) \in (\mathbb{R}^n_{+,\star})^2$$
$$K)d_H(u,v),$$

[2] Birkhoff, Extensions of Jentzsch's theorem, Transaction of the American Mathematical Society 85 (1957), no. 1, 219-227.

[12] Samelson et al, On the Perron-Frobenieus theorem, The Michigan Mathematical Journal 4 (1957), no. 1, 57-59



Thm(Franklin and Lorenz '89): One has $(D_{\omega}^{k}, D_{\omega}^{k}) \rightarrow (D_{\omega}^{*}, D_{\omega}^{*})$ and

 $d_H(D_{\omega}^k, D_{\omega}^*) = O(\lambda(K)^{2k}), \quad d_H(D_{w}^k, D_{w}^*) = O(\lambda(K)^{2k}).$

Moreover,

 $d_H(D_{\varphi}^k, D_{\varphi}^*) \leq \frac{d_H(\gamma^k \mathbf{1}_n, \mu)}{1 - \lambda(K)^2},$ $d_H(D_{\psi}^k, D_{\psi}^*) \leq \frac{d_H(\gamma^{k,T}\mathbf{1}_n, \nu)}{1 - \lambda(K)^2},$

Where $\gamma^k = D_{\omega}^k K D_{\omega}^k$. Last, one has $\left\| \log(\gamma^k) - \log(\gamma_{\varepsilon}) \right\|_{\infty} \le d_H(D_{\varphi}^k, D_{\varphi}^*) + d_H(D_{\psi}^k, D_{\psi}^*).$

Sketch of proof:

First, notice that

$$d_H(u,v) = d_H(u/v)$$

Together with the previous theorem, this gives

$$d_H(D_{\varphi}^k, D_{\varphi}^{\star}) = d_H(\frac{\mu}{KD_{\psi}^k}, \frac{\mu}{KD_{\psi}^{\star}}) =$$

 $(\mathbf{1}_n) = d_H(\mathbf{1}_n/u, \mathbf{1}_n/v).$

 $= d_H(KD_{\psi}^k, KD_{\psi}^{\star}) \leqslant \lambda(K)d_H(D_{\psi}^k, D_{\psi}^{\star}).$



Then by using the triangular inequality we have

$$d_{H}(D_{\varphi}^{k}, D_{\varphi}^{\star}) \leq d_{H}(D_{\varphi}^{k+1}, D_{\varphi}^{k}) =$$
$$\leq d_{H}(\frac{\mu}{KD_{\psi}^{k}}, D_{\varphi}^{k})$$
$$= d_{H}(\mu, D_{\varphi}^{k} \odot (K))$$
$$= d_{H}(\mu, \gamma^{k} \mathbf{1}_{n}) + \lambda$$

Where \odot denotes the element-wise multiplication.

$+ d_H(D_{\varphi}^{k+1}, D_{\varphi}^{\star})$ $+ \lambda(K) d_H(D_{\varphi}^k, D_{\varphi}^{\star})$ $T D_{\psi}^k)) + \lambda(K)^2 d_H(D_{\varphi}^k, D_{\varphi}^{\star})$ $\lambda(K)^2 d_H(D_{\varphi}^k, D_{\varphi}^{\star}),$

Rmk (stopping criteria): the bounds d_H $d_{H}(D_{\psi}^{k}, D_{\psi}^{*}) \leq \frac{d_{H}(\gamma^{k, T} \mathbf{1}_{n}, \nu)}{1 - \lambda(K)^{2}} \text{ shows that some error measures on the marginal constraints violation, for instance } \|\gamma^{k} \mathbf{1}_{n} - \mu\|_{1} \text{ and } \|\gamma^{k, T} \mathbf{1}_{n} - \nu\|_{1}, \text{ are useful}$ stopping criteria to monitor the convergence.

becomes exponentially bad as $\varepsilon \to 0$, since it scales like $e^{-1/\varepsilon}$.

$$d_{H}(D_{\varphi}^{k}, D_{\varphi}^{*}) \leq \frac{d_{H}(\gamma^{k} \mathbf{1}_{n}, \mu)}{1 - \lambda(K)^{2}}$$
 and

Rmk: This theorem shows that Sinkhorn algorithm converges linearly, but the rates



Back to the continuous case

$$\mathcal{K}_{\varepsilon}(\mu,\nu) = \inf \left\{ \int_{X \times Y} c(x,y) d\gamma(x,y) + \varepsilon \mathcal{H}(\gamma \mid \mu \otimes \nu) \mid \gamma \in \Pi(\mu,\nu) \right\},\$$

Where

$$\mathcal{H}(\rho \mid \pi) = \begin{cases} \int_{X \times Y} \left(\log \left(\frac{d\rho(x, y)}{d\pi(x, y)} \right) \\ +\infty, \end{cases} \right)$$

And μ , ν are probability measures on the compact sets X and Y.

One can easily recast the regularized OT in the continuous framework as follows

$$-1)d\rho(x,y), \quad \text{if } \rho \ll \pi$$

otherwise,



Linear convergence of Sinkhorn for bounded cost

Consider the following variant of Sinkhorn algorithm

$$\begin{split} \varphi^{k+1}(x) &= -\varepsilon \log \left(\int_X \exp\left(\frac{1}{\varepsilon} (\psi^k(y) - c(x, y)) d\nu(y)\right) + \lambda^k \\ \psi^{k+1}(y) &= -\varepsilon \log \left(\int_X \exp\left(\frac{1}{\varepsilon} (\varphi^{k+1}(x) - c(x, y)) d\mu(x)\right) \right), \\ \text{Where } \lambda^k &= \varepsilon \int_X \log \left(\int_X \exp\left(\frac{1}{\varepsilon} (\psi_k(y) - c(x, y)) d\nu(y)\right) d\mu(x). \end{split}$$

This is equivalent to previous Sinkhorn...I am just fixing a constant.

method

$$\varphi^{k+1} = \operatorname{argmax}_{\varphi, \int \varphi d\mu = 0} \Phi_{\varepsilon}(\varphi, \psi^{k}),$$
$$\psi^{k+1} = \operatorname{argmax}_{\psi} \Phi_{\varepsilon}(\varphi^{k+1}, \psi).$$

$$\varphi^{k+1} = \operatorname{argmax}_{\varphi, \int \varphi d\mu = 0} \Phi_{\varepsilon}(\varphi, \psi^{k}),$$
$$\psi^{k+1} = \operatorname{argmax}_{\psi} \Phi_{\varepsilon}(\varphi^{k+1}, \psi).$$

Where

$$\begin{split} \Phi_{\varepsilon}(\phi,\psi) &:= \int_{X} \varphi(x) d\mu(x) + \int_{Y} \psi(y) d\nu(y) \\ &- \varepsilon \int_{X \times Y} \exp\left(\frac{\phi(x) + \psi(y) - c(x,y)}{\varepsilon}\right) d\mu \otimes d\nu(x,y) \,. \end{split}$$

The algorithm in the previous slide is equivalent to the following coordinate ascent

Rmk [alternative dual formulation]: one can sho that the dual problem can also rewritten in the following way

 ϕ,ψ

Where

$$\begin{split} \tilde{\Phi}_{\varepsilon}(\phi,\psi) &:= \int_{X} \varphi(x) \mathrm{d}\mu(x) + \int_{Y} \psi(y) \mathrm{d}\nu(y) \\ &-\varepsilon \mathrm{log}\bigg(\int_{X \times Y} \exp\bigg(\frac{\phi(x) + \psi(y) - c(x,y)}{\varepsilon} \bigg) \mathrm{d}\mu \otimes \mathrm{d}\nu(x,y) \bigg) \,. \end{split}$$

To see it just use the variational representation of the relative entropy, that is

$$\mathscr{H}(\rho \mid \pi) = \sup_{\phi} \left(\int \phi d\rho - \log \left(\int e^{\phi} d\pi \right) \right).$$

 $\sup \tilde{\Phi}_{\varepsilon}(\varphi, \psi),$



Lemma: For every $k \ge 0$ we have $\left\| \left| \varphi^k \right| \right\|_{\infty} \le 2 \left\| \left| c \right| \right\|_{\infty} \quad \epsilon$ **Proof** Just compute $\varphi^k(x_1) - \varphi^k(x_2)$.

- **Thm:** Let (φ^*, ψ^*) the unique solution of the dual entropic problem with $\int_{V} \varphi^*(x) d\mu(x) = 0$. The iterates of Sinkhorn satisfy
- $\Phi_{\varepsilon}(\varphi^*, \psi^*) \Phi_{\varepsilon}(\varphi^k, \psi^k) \leq \beta^k$ $\| \varphi^* - \varphi^k \|_{I^2}^2 + \| \psi^* - \psi^k \|_{I^2}^2 \le \eta \beta$ Where $\beta := 1 - e^{-24||c||_{\infty}/\varepsilon}$ and $\eta = 2e$

and
$$||\psi^k||_{\infty} \leq 3||c||_{\infty}$$
.

$$k(\Phi_{\varepsilon}(\varphi^{*},\psi^{*}) - \Phi_{\varepsilon}(\varphi^{0},\psi^{0})),$$

$$\beta^{k}(\Phi_{\varepsilon}(\varphi^{*},\psi^{*}) - \Phi_{\varepsilon}(\varphi^{0},\psi^{0})),$$

$$\beta^{6}||c||_{\infty}/\varepsilon$$

Sketch of proof:

The basic idea is to use strong convexity of the exponential function on an interval $[-\alpha, +\infty)$, that is

$$e^{b} - e^{a} \ge (b - a)e^{a} + \frac{e^{-\alpha}}{2} |b - a|^{2}, \quad \text{for} \quad a, b \in [-\alpha, +\infty).$$

Step 1: Given $\sigma = e^{-6||c||_{\infty}/\varepsilon}$ we have

$$\Phi_{\varepsilon}(\varphi^{k+1}, \psi^{k+1}) - \Phi_{\varepsilon}(\varphi^{k}, \psi^{k}) \ge \frac{\sigma}{2}(||\varphi^{k+1} - \varphi^{k}||_{L^{2}}^{2} + \psi^{k+1} - \psi^{k}||_{L^{2}}^{2}).$$

for
$$a, b \in [-\alpha, +\infty)$$
.



Where

 $\partial_1 \Phi_{\mathcal{L}}(\varphi, \psi) = 1 - 1$

 $+ \int_{V} \partial_2 \Phi_{\varepsilon}(\varphi^k, \psi^k)(y) [\psi^k(y) - \psi^*(y)] d\nu(y)$

 $+\frac{\sigma}{2}(||\varphi^{k}-\varphi^{*}||_{L^{2}}^{2}+||\psi^{k}-\psi^{*}||_{L^{2}}^{2})$

$$\varepsilon \int_{Y} e^{\frac{\varphi + \psi - c}{\varepsilon}} d\nu(y)$$

Step 3: By exploring the zero mean iterates and Young's inequality we get
$$\Phi_{\varepsilon}(\varphi^*, \psi^*) - \Phi_{\varepsilon}(\varphi^k, \psi^k) \leq \frac{1}{2\sigma} ||\partial_1 \Phi_{\varepsilon}(\varphi^k, \psi^k) - \partial_1 \Phi_{\varepsilon}(\varphi^{k+1}, \psi^k)||_{L^2}^2.$$
Now by Lipschitz continuity of the exponential and **step 1** we have
$$\Phi_{\varepsilon}(\varphi^*, \psi^*) - \Phi_{\varepsilon}(\varphi^k, \psi^k) \leq \frac{1}{\sigma^4} (\Phi_{\varepsilon}(\varphi^{k+1}, \psi^{k+1}) - \Phi_{\varepsilon}(\varphi^k, \psi^k))$$
Taking $\Delta^k = \Phi_{\varepsilon}(\varphi^*, \psi^*) - \Phi_{\varepsilon}(\varphi^k, \psi^k)$ the above inequality can be expressed
$$\Delta^{k+1} \leq (1 - \sigma^4) \Delta^k$$

Iterating we get the result.

et

ressed as

