

Uniform Quantitative Stability of Sinkhorn

George Deligiannidis

Joint with V. de Bortoli and A. Doucet

Department of Statistics, Oxford University

Athens Probability Colloquium, March 2023

Optimal Transport

- **Coupling:** A coupling of two probability measures μ, ν , on spaces $X, Y \subset \mathbb{R}^d$ respectively, a **coupling** is a **joint distribution** π , with μ, ν as its marginals,

$$\int f(x)\pi(dx, dy) = \int f(x)\mu(dx), \quad \int g(y)\pi(dx, dy) = \int g(y)\nu(dy).$$

- We write $\mathcal{C}(\mu, \nu)$ for the collection of couplings of μ, ν .
- **Optimal Transport:** the basic problem of Optimal Transport is to find a coupling $\pi \in \mathcal{C}(\mu, \nu)$ that minimizes

$$\inf_{\pi \in \mathcal{C}(\mu, \nu)} \int \pi(dx, dy)c(x, y); \quad \text{OT}(\mu, \nu)$$

where $c : X \times Y \rightarrow \mathbb{R}$ is any cost function.

Optimal Transport

$$\inf_{\pi \in \mathcal{C}(\mu, \nu)} \int \pi(dx, dy) c(x, y);$$

- this formulation goes back to Kantorovich;
- the earlier formulation by Monge considers mappings $y = T(x)$, that is $\pi(dx, dy) = \mu(dx) \delta_{T(x)}(dy)$ and results in non-convex optimisation.
- **Brenier 1987**: Assuming $c(x, y) = |x - y|^2$, μ, ν have finite second moments and μ gives measure 0 to all sets of Hausdorff dimension $\leq d - 1$, there exists a convex mapping $\varphi : \mathcal{X} \mapsto \mathcal{Y}$ such that optimal coupling takes the form of an optimal transport map $\mu(dx) \delta_{\nabla \varphi(x)}(dy)$
- When μ is discrete, π is a proper coupling.
- For discrete measures with N atoms, worst case computational cost is $O(N^{5/2})$ for assignment problem, using Hungarian (auction) algorithm, see e.g. **Méridot and Oudet 2016**.
- This can be improved by considering the entropy regularised version.

Entropy Regularised Optimal Transport

- **Cuturi 2013** realised that the entropy regularised problem

$$\Pi_{\epsilon}^{\mu, \nu} = \arg \min_{\pi \in \mathcal{C}(\mu, \nu)} \int \pi(dx, dy) \|x - y\|^2 + \epsilon \text{KL}(\pi | \mu \otimes \nu) \quad \text{OT}_{\epsilon}(\mu, \nu)$$

can be solved efficiently using the **Iterative Proportional Fitting Procedure (IPFP)**, aka **Sinkhorn's algorithm**.

- the computational cost is roughly $O(N^2)$, see **Altschuler, Weed, and Rigollet 2017**, ignoring dimensionality.
- Equivalent to (static) Schrödinger bridge problem

$$\inf_{\pi \in \mathcal{C}(\mu, \nu)} \text{KL}(\pi | \Gamma_{\epsilon}), \quad \Gamma_{\epsilon}(dx, dy) = \exp[-c(x, y)/\epsilon] \mu(dx) \nu(dy).$$

- Solution always in the form

$$\frac{d\Pi_{\epsilon}^{\mu, \nu}}{d\mu \otimes \nu} = \exp[\varphi_{\epsilon}(x) + \psi_{\epsilon}(y) - c(x, y)/\epsilon],$$

where $(\varphi_{\epsilon}, \psi_{\epsilon})$ are known as the **potentials**.

Schrödinger bridge problem

- **Question:** what is the most likely evolution of a stochastic process $\{X_t : t \in [0, 1]\}$, $X_0 \sim \mu$ and $X_1 \sim \nu$?
- one formulation is

$$\operatorname{argmin}_{\mathbb{P}} \operatorname{KL}(\mathbb{P} | \mathbb{W}_\epsilon) \quad (1)$$

- minimization over $\mathcal{P}(\mathcal{C}([0, 1]; \mathbb{R}^n))$ with marginals μ and ν at times 0 and 1 respectively;
- \mathbb{W}_ϵ is the law of Brownian motion multiplied by $\sqrt{\epsilon/2}$.
- The solution to (1) can then be expressed in terms of the solution to $\operatorname{OT}_\epsilon(\mu, \nu)$ and the bridges of the measure \mathbb{W}_ϵ .
- The quadratic cost arises from the log-transition density of Brownian motion—different reference measures can give rise to different costs.

Recent applications of Schrödinger bridges

How is the static Schrödinger bridge useful in practice?

- Idea is to use solution π_ϵ^* of OT_ϵ and treat it as an approximation to the solution π^* of OT.
- A lot of recent progress quantifying $\pi_\epsilon^* \rightarrow \pi^*$.
- In recent applications, the Schrödinger bridge is used for its own benefits rather than as a computationally feasible approximation to the standard optimal transport problem.

Differentiable Particle Filtering

- **Corenflos et al. 2021** used the solution to Schrödinger bridge to build a particle filtering scheme.
- Suppose $\hat{\pi}_t = N^{-1} \sum_{j=1}^N \delta_{X_j^t}$ is particle approximation of π_t ;
- Let

$$\tilde{X}_j^{t+1} \sim q_{t+1}(\cdot | X_j^t), j \in [N]; \quad \omega_j^{t+1} = \frac{f_\theta(\tilde{X}_j^{t+1} | X_j^t) g_\theta(Y^{t+1} | \tilde{X}_j^{t+1})}{q_{t+1}(\tilde{X}_j^{t+1} | X_j^t)}$$

$$\mathbb{P}(I_j^{t+1} = k) = \frac{\omega_k^{t+1}}{\sum_{j=1}^N \omega_j^{t+1}}$$

$$X_j^{t+1} \stackrel{\text{iid}}{\sim} \sum_{k=1}^N \frac{\omega_k^{t+1}}{\sum_{j=1}^N \omega_j^{t+1}} \delta_{\tilde{X}_k^{t+1}}$$

- Resampling: estimators not differentiable wrt θ .
- There are approximate ways to bypass this.

Differentiable Particle Filtering

- **Corenflos et al. 2021** solves the ϵ -entropy regularised OT problem between

$$\alpha_N^{(t)} = \frac{1}{N} \sum_{i=1}^N \delta_{\tilde{X}_t^i}, \quad \beta_N^{(t)} = \sum_i w_t^i \delta_{\tilde{X}_t^i}$$

to obtain a matrix $\mathbf{P}_\epsilon = \mathbf{P}_\epsilon(\theta)$;

- resampling replaced by **ensemble transform** $\mathbf{X}_t = \mathbf{P}_\epsilon \tilde{\mathbf{X}}_t$
- for $\epsilon > 0$ $\theta \mapsto \mathbf{P}_\epsilon(\theta)$ differentiable \implies end to end differentiable estimators.
- Proof of consistency: compare with the solution of the unregularised OT problem (needs $\epsilon = \epsilon_N \rightarrow 0$)
- OR we can compare the solution of $\text{OT}_\epsilon(\alpha_N, \beta_N)$ with that of $\text{OT}_\epsilon(\alpha_N, \beta_N)$ where $\alpha = \lim_{N \rightarrow \infty} \alpha_N, \beta = \lim_N \beta_N$.

Stability of OT

This brings us to the following question

- **Question:** suppose $\alpha_n \rightarrow \alpha$ and $\beta_n \rightarrow \beta$ (say weakly);
- Does the solution of $OT(\alpha_n, \beta_n)$ converges in some sense to that of $OT(\alpha, \beta)$?
- In the realm of Brenier's theorem, we could talk about convergence of the transport maps.
- In classical Optimal Transport there is a classical, qualitative solution.

Stability of OT

Theorem (5.20 in Villani 2009)

Let \mathcal{X}, \mathcal{Y} be Polish spaces and let $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a continuous, lower bounded cost function. Let $(\mu_k), (\nu_k)$ be sequences of probability measures on \mathcal{X}, \mathcal{Y} respectively, such that $\mu_k \rightarrow \mu$ and $\nu_k \rightarrow \nu$ weakly. For each k , let π_k be an optimal transportation plan between μ_k and ν_k .

If for all $k \in \mathbb{N}$, $\int c d\pi_k < +\infty$, then up to extraction of a subsequence, π_k converges weakly to some c -cyclically monotone transportation plan $\pi \in \mathcal{C}(\mu, \nu)$.

If moreover $\liminf_k \int c d\pi_k < \infty$, then the optimal transport cost between μ, ν is finite and π is an optimal transport plan.

Stability of OT

- If the optimal transport plan between μ and ν is unique then there is no need to extract a subsequence.
- If plan is given by map, maps converge in probability (Corollary 5.23 Villani 2009).
- Proof is based on compactness and characterization of optimal transport plans in terms of c -cyclical monotonicity: a set $\Gamma \subset X \times Y$ is called c -cyclically monotone if for all collections $(x_1, y_1), \dots, (x_N, y_N) \subset \Gamma$ we have

$$\sum_{i=1}^N c(x_i, y_i) \leq \sum_{i=1}^N c(x_i, y_{i+1}), \quad y_{N+1} := y_1.$$

a transport plan is called c -cyclically monotone if its support is.

Quantitative Stability of OT

Li and Nochetto 2021

- $\mu_h \approx \mu$, $\nu_h \approx \nu$, γ_h solves $\text{OT}(\mu_h, \nu_h)$, and $T = \nabla\phi$ is the unique λ -Lipschitz optimal transport map between μ, ν . Then

$$\left(\int_{X \times Y} |T(x) - y|^2 d\gamma_h(x, y) \right)^{1/2} \leq 2\lambda^{1/2} \Delta_h^{1/2} [\mathbf{W}_2(\mu, \nu) + \Delta_h]^{1/2} + (\lambda \vee 1) \Delta_h, \quad (2)$$

$$\Delta_h := \mathbf{W}_2(\mu, \mu_h) + \mathbf{W}_2(\nu, \nu_h) \quad (3)$$

Méridot, Delalande, and Chazal 2020

- ρ, μ, ν p.m.s, $T_\mu^\# \rho = \mu$, $T_\nu^\# \rho = \nu$ Brenier maps.
- Then

$$\mathbf{W}_2(\mu, \nu) \leq \|T_\mu - T_\nu\|_{L^2(\rho)} \leq C \mathbf{W}_\rho(\mu, \nu)^{2/15},$$

for any $\rho \geq 1$, where C depends on the dimension and the sets X, Y .

- If T_μ is K -Lipschitz then the exponent can be improved to $1/2$.
- They also show that the Monge embedding $\mu \mapsto T_\mu$ is in general not better than $1/2$ -Holder.

Stability of Regularised OT

- Until 2-3 years ago much less was known on stability of regularised OT.
- Analogue of cyclical monotonicity recently done in Ghosal, Nutz, and Bernton 2021 who proved the qualitative stability of Schrödinger bridges (no compactness required).
- For compact state spaces Luise et al. 2019 prove that

$$\|\varphi_{\epsilon}^{\mu, \nu} - \varphi_{\epsilon}^{\mu', \nu'}\|_{\infty} \leq C(d, \epsilon, X, Y) \{\|\mu - \mu'\|_{\text{TV}} + \|\nu - \nu'\|_{\text{TV}}\}.$$

NOTE: TV is too strong to capture convergence of empirical measures.

- For smooth costs, $c \in C^{s+1}$, $s > d/2$, and ν_n an empirical version of ν , using MMD metrics, w.prob $> 1 - \tau$

$$\|\varphi_{\epsilon}^{\mu, \nu} - \varphi_{\epsilon}^{\mu, \nu_n}\|_{\infty} \leq C(d, \epsilon, X, Y) \log(3/\tau) n^{-1/2},$$

Sample complexity of cost

Stability is related to **sample complexity**.

- Suppose we want to estimate $\mathbf{W}_\rho(\mu, \nu)$ using samples from μ, ν .
- **How many samples do we need to be within ϵ from the truth?**
- Let $S_\epsilon(\mu, \nu)$ be the cost of the OT_ϵ , write $S(\mu, \nu)$ for the unregularised.
- If μ_n, ν_n empirical approximations to μ, ν in general $S(\mu, \nu) = S(\mu_n, \nu_n) + O(n^{-1/d})$;
- BUT miraculously for any $\epsilon > 0$, $S_\epsilon(\mu, \nu) = S_\epsilon(\mu_n, \nu_n) + O(n^{-1/2})$, see **Genevay et al. 2019; Mena and Niles-Weed 2019**.

Sample complexity of Coupling

- Interested in full coupling $\Pi_{\epsilon}^{\mu, \nu}$: when measured say in Wasserstein distance, sample complexity depends on dimension;
- Say μ_n, ν_n are empirical versions of μ, ν . Then trivially we have that

$$\mathbf{W}_1(\Pi_{\epsilon}^{\mu_n, \nu_n}, \Pi_{\epsilon}^{\mu, \nu}) \geq \mathbf{W}_1(\mu_n, \mu) = \Omega(n^{-1/d}),$$

see e.g. [Fournier and Guillin 2015](#).

- But \mathbf{W} is a uniform metric. What if we only care about expectations of just a few functionals?
- [Rigollet and Stromme 2022](#) obtain $O(n^{-1/2})$ rates in [non-uniform metrics](#), e.g. measuring expectation of a single function.

Iterative Proportional Fitting (aka Sinkhorn) I

- Given two probability measures μ, ν and $\epsilon > 0$ let $\Pi^{(0)} = e^{-\epsilon} \mu \otimes \nu$;
- the IPFP iteratively solves the following one-sided Bridge problems

$$\Pi^{(2j+1)} = \arg \min_{\Pi: \Pi_0 = \mu} \text{KL} \left(\Pi \mid \Pi^{(2j)} \right)$$

$$\Pi^{(2j+2)} = \arg \min_{\Pi: \Pi_1 = \nu} \text{KL} \left(\Pi \mid \Pi^{(2j+1)} \right).$$

- it alternates between projecting on measures with first marginal μ and measures with second marginal ν .
- It's easy to see that we can write

$$\frac{d\Pi^{(j)}}{d\Pi^{(j-1)}}(x, y) := \frac{d\nu}{d\Sigma_2^\# \Pi^{(j-1)}}(y), \quad j \text{ odd,}$$

$$\frac{d\Pi^{(j)}}{d\Pi^{(j-1)}}(x, y) := \frac{d\mu}{d\Sigma_1^\# \Pi^{(j-1)}}(x), \quad j \text{ even,}$$

where $\Sigma_1(x, y) = x$, $\Sigma_2(x, y) = y$.

Iterative Proportional Fitting (aka Sinkhorn) II

- The **dual formulation** in terms of potentials is

$$\varphi^{(t+1)}(x) := -\log \int \exp\{\psi^{(t)}(y) - c(x,y)/\epsilon\} \nu(dy)$$

$$\psi^{(t+1)}(y) := -\log \int \exp\{\varphi^{(t+1)}(x) - c(x,y)/\epsilon\} \mu(dx)$$

- For any $c \in \mathbb{R}$, $(\varphi^{(t)} - c, \psi^{(t)} + c)$ defines the same measure; we fix this choice following **Carlier 2021**, so that $\mu[\varphi^{(t+1)}] = 0$.

Main results

Theorem 1

Suppose that X, Y are compact metric spaces and $c \in \text{Lip}(X \times Y)$. For any $\pi_0, \hat{\pi}_0 \in \mathcal{P}(X)$, $\pi_1, \hat{\pi}_1 \in \mathcal{P}(Y)$ let $(\mathbb{P}^n)_{n \in \mathbb{N}}$ and $(\hat{\mathbb{P}}^n)_{n \in \mathbb{N}}$ the IPFP sequence with marginals (π_0, π_1) respectively $(\hat{\pi}_0, \hat{\pi}_1)$. Then any $n \in \mathbb{N}$ we have

$$\mathbf{W}_1(\mathbb{P}^n, \hat{\mathbb{P}}^n) \leq C \{ \mathbf{W}_1(\pi_0, \hat{\pi}_0) + \mathbf{W}_1(\pi_1, \hat{\pi}_1) \}, \quad (4)$$

with

$$C = e^{10\|c\|_\infty} \{ 1 + (2\text{Lip}(c) + 10)(\text{diam}(X) + \text{diam}(Y)) \}. \quad (5)$$

- **Good:** The result is **uniform** in the iterations of Sinkhorn.
- **Good:** Sample complexity is sharp.
- **The bad:** Constant is terrible.

Main results (ctd)

As an immediate consequence of Theorem 1 and the fact that the IPFP sequence converges, we obtain the quantitative stability of Schrödinger bridge.

Corollary 2

For any $\pi_0, \hat{\pi}_0 \in \mathcal{P}(X)$, $\pi_1, \hat{\pi}_1 \in \mathcal{P}(Y)$ let \mathbb{P}^ , respectively $\hat{\mathbb{P}}^*$, be the Schrödinger bridge with marginals (π_0, π_1) , respectively $(\hat{\pi}_0, \hat{\pi}_1)$. Then we have*

$$\mathbf{W}_1(\mathbb{P}^*, \hat{\mathbb{P}}^*) \leq C \{ \mathbf{W}_1(\pi_0, \hat{\pi}_0) + \mathbf{W}_1(\pi_1, \hat{\pi}_1) \}, \quad (6)$$

with C as in Theorem 1.

Background on Hilbert projective metric I

Will now sketch the main idea for the Schrödinger bridge, rather than IPFP as it is a little clearer.

- In compact spaces we can employ the machinery of the **Birkhoff-Hopf contraction theorem**;
- contraction w.r.t. the **Hilbert metric**.
- It is a projective metric in the sense that it measures distances between rays $\{\lambda x : \lambda \geq 0\}$ rather than points.
- For function spaces: **oscillation** in the log-scale, ie

$$d_H(f, g) = \|\log(f/g)\|_{\text{osc}} := \sup_x \log \frac{f(x)}{g(x)} - \inf_y \log \frac{f(y)}{g(y)}.$$

We are now ready to state the Birkhoff contraction theorem.

- $(V, \|\cdot\|), (V', \|\cdot\|)$ normed real vector spaces;
- $K \subset V, K' \subset V'$ cones;
- $C \subset K, C' \subset K'$ convex parts;

Background on Hilbert projective metric II

- d_H, d'_H the Hilbert metric on C, C' respectively;
- $T : V \rightarrow V'$ is a linear mapping such that $T(C) \subset C'$.

Theorem (Birhoff Contraction Theorem)

Then

$$\kappa(T) := \sup_{x,y \in C} \frac{d'_H(T(x), T(y))}{d_H(x,y)} \leq \tanh(\Delta(T)/4), \quad (7)$$

where the projective diameter $\Delta(T)$ of T is defined by

$$\Delta(T) := \sup \{d'_H(T(x), T(y)) : x, y \in C, \|x\| = \|y\| = 1\}.$$

Since $\Delta(T)$ is finite $\kappa(T) < 1$.

Sketch of proof I

- In our context

$$C = \mathcal{C}(X; (0, \infty)), \quad C' = \mathcal{C}(Y; (0, \infty)).$$

- Letting $\mu \in \mathcal{P}(X), \nu \in \mathcal{P}(Y)$ define two linear maps of interest through

$$(\mathcal{E}_\mu f)(y) := \int f(x)K(x,y)\mu(dx); \quad \mathcal{E}_\mu : C \rightarrow C' \quad (8)$$

$$(\mathcal{E}_\nu g)(x) := \int g(y)K(x,y)\nu(dy); \quad \mathcal{E}_\nu : C' \rightarrow C. \quad (9)$$

- Let also $\mathcal{J} : f \mapsto f^{-1}$, where we overload the operator to act on both C and C' .
- As pointed out in [Chen, Georgiou, and Pavon 2016](#) it is easy to show that \mathcal{J} is an isometry w.r.t. the Hilbert metric;
- Also we can bound the projective diameter of $\mathcal{E}_\mu, \mathcal{E}_\nu$.

Sketch of proof II

- In this notation the Sinkhorn iteration that takes $\mathbf{e}^{\varphi^{(t)}} \rightarrow \mathbf{e}^{\varphi^{(t+1)}}$ can be written as

$$\mathcal{S}_{\mu,\nu} \exp[\varphi^{(t)}] := [\mathcal{I} \circ \mathcal{E}_\nu \circ \mathcal{I} \circ \mathcal{E}_\mu] [\exp \varphi^{(t)}], \quad (10)$$

$$\mathcal{S}_{\mu,\nu}^\dagger \exp[\psi^{(t)}] := [\mathcal{I} \circ \mathcal{E}_\mu \circ \mathcal{I} \circ \mathcal{E}_\nu] \exp[\psi^{(t)}]. \quad (11)$$

- The Birkhoff contraction theorem show that both of these maps are contractions in the Hilbert metric.
- The **potentials** $(\varphi_\epsilon, \psi_\epsilon)$, $(\widehat{\varphi}_\epsilon, \widehat{\psi}_\epsilon)$ corresponding to (μ, ν) , $(\widehat{\mu}, \widehat{\nu})$ are fixed points of $(\mathcal{S}_{\mu,\nu}, \mathcal{S}_{\mu,\nu}^\dagger)$ and $(\mathcal{S}_{\widehat{\mu},\widehat{\nu}}, \mathcal{S}_{\widehat{\mu},\widehat{\nu}}^\dagger)$ resp.

Sketch of proof: main ideas I

For $h \in \text{Lip}(X \times Y)$, writing

$$F_\epsilon(x, y) = \exp\{\varphi_\epsilon(x) + \psi_\epsilon(y)\}, \quad \widehat{F}_\epsilon(x, y) = \exp\{\widehat{\varphi}_\epsilon(x) + \widehat{\psi}_\epsilon(y)\}$$

we want to write

$$\begin{aligned} & \int h(x, y) \left[\Pi_\epsilon^{\mu, \nu}(dx, dy) - \Pi_\epsilon^{\widehat{\mu}, \widehat{\nu}}(dx, dy) \right] \\ &= \int h(x, y) F_\epsilon(x, y) K_\epsilon(x, y) \mu(dx) \nu(dy) - \int h(x, y) \widehat{F}_\epsilon(x, y) K_\epsilon(x, y) \widehat{\mu}(dx) \widehat{\nu}(dy) \\ &= \int h \left[F_\epsilon - \widehat{F}_\epsilon \right] K_\epsilon(x, y) \mu(dx) \nu(dy) + \int h \widehat{F}_\epsilon K_\epsilon(x, y) [\mu \otimes \nu - \widehat{\mu} \otimes \widehat{\nu}]. \end{aligned}$$

- The second term looks like it can be controlled by $\mathbf{W}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu})$;
- But the potentials \widehat{F}_ϵ are not even necessarily smooth, or even well-defined everywhere....

Sketch of proof: main ideas II

- **Canonical extensions:** although $\hat{\varphi}_\epsilon, \hat{\psi}_\epsilon$ are only defined on the supports of $\hat{\mu}, \hat{\nu}$ they are fixed points of the Sinkhorn iteration, that is

$$\hat{\varphi}_\epsilon(y) = -\log \int \exp\{\hat{\psi}_\epsilon(y) - c(x,y)/\epsilon\} \hat{\nu}(dy) \quad (12)$$

$$\hat{\psi}_\epsilon(y) = -\log \int \exp\{\hat{\varphi}_\epsilon(x) - c(x,y)/\epsilon\} \hat{\mu}(dx). \quad (13)$$

- We can use this to **extend** $\hat{\varphi}_\epsilon, \hat{\psi}_\epsilon$ X, Y ;
- they inherit smoothness from cost, see also [Luise et al. 2019](#),
- So indeed

$$\int h(x,y) \left[\Pi_\epsilon^{\mu,\nu}(dx, dy) - \Pi_\epsilon^{\hat{\mu},\hat{\nu}}(dx, dy) \right] \leq \int h \left[F_\epsilon - \hat{F}_\epsilon \right] K_\epsilon(x,y) \mu(dx) \nu(dy) + \mathbf{W}_1(\mu \otimes \nu, \hat{\mu} \otimes \hat{\nu}).$$

- If we can control $\|F_\epsilon - \hat{F}_\epsilon\|_\infty$ we should be done.

Sketch of proof: main ideas I

- Next we use the fact that $F_\epsilon, \hat{F}_\epsilon$ are **fixed points** of **Sinkhorn iterations** to control $F_\epsilon - \hat{F}_\epsilon$.
- Idea here is

$$\begin{aligned}d_H(F_\epsilon, \hat{F}_\epsilon) &= d_H(\mathcal{S}_{\mu, \nu} F_\epsilon, \mathcal{S}_{\hat{\mu}, \hat{\nu}} \hat{F}_\epsilon) \\&= \underbrace{d_H(\mathcal{S}_{\mu, \nu} F_\epsilon, \mathcal{S}_{\mu, \nu} \hat{F}_\epsilon)}_{\text{contraction of Sinkhorn}} + \underbrace{d_H(\mathcal{S}_{\mu, \nu} \hat{F}_\epsilon, \mathcal{S}_{\hat{\mu}, \hat{\nu}} \hat{F}_\epsilon)}_{\leq \mathbf{W}_1(\mu \otimes \nu, \hat{\mu} \otimes \hat{\nu})} \\&= \kappa d_H(F_\epsilon, \hat{F}_\epsilon) + C \mathbf{W}_1(\mu \otimes \nu, \hat{\mu} \otimes \hat{\nu}) \\d_H(F_\epsilon, \hat{F}_\epsilon) &\leq \frac{C}{1 - \kappa} \mathbf{W}_1(\mu \otimes \nu, \hat{\mu} \otimes \hat{\nu}).\end{aligned}$$

- Final issue is that $d_H(F_\epsilon, \hat{F}_\epsilon)$ only controls the oscillations $\|\log F_\epsilon - \log \hat{F}_\epsilon\|_{\text{osc}}$ rather than the supremum;

Sketch of proof: main ideas II

- To bypass this issue notice that \widehat{F}_ϵ defines a **probability density** w.r.t. reference measure $K_\epsilon d\widehat{\mu} \otimes \widehat{\nu}$.

$$\int \frac{\widehat{F}_\epsilon}{F_\epsilon} F_\epsilon K_\epsilon d\mu \otimes \nu = \int \widehat{F}_\epsilon K_\epsilon d\widehat{\mu} \otimes \widehat{\nu} + \mathbf{CW}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu}) \quad (14)$$

$$= \int \widehat{F}_\epsilon K_\epsilon \widehat{\mu} \otimes \widehat{\nu} + \mathbf{CW}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu}) \quad (15)$$

$$= 1 + \mathbf{CW}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu}) \quad (16)$$

$$(17)$$

- Recall $F_\epsilon K_\epsilon \mu \otimes \nu$ is a probability measure;
- thus the random variable $\widehat{F}_\epsilon / F_\epsilon(X, Y)$ with $(X, Y) \sim F_\epsilon K_\epsilon \mu \otimes \nu$, must either be a.s. equal to $1 + \mathbf{CW}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu})$, or must take values both above and below $1 + \mathbf{CW}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu})$.
- But recall that we have controlled the oscillation $\|\log F_\epsilon - \log \widehat{F}_\epsilon\|_{\text{osc}}$;

Sketch of proof: main ideas III

- therefore we can anchor $\log F_\epsilon / \widehat{F}_\epsilon$ near 1 and therefore

$$\log \frac{F_\epsilon}{\widehat{F}_\epsilon} \leq C \mathbf{W}_1(\mu \otimes \nu, \widehat{\mu} \otimes \widehat{\nu}) + \|\log \widehat{F}_\epsilon(x, y) - \log F_\epsilon(x, y)\|_{\text{osc}}.$$

- We can similarly control $\widehat{F}_\epsilon / F_\epsilon$.
- A little more work obtains stability **uniform** in the iterations of Sinkhorn. □

Recent progress (an incomplete summary)

Eckstein and Nutz 2021

- **non-compact** space; consider the Schrödinger bridge, ie the limit of Sinkhorn and prove Hölder continuity, rather than Lipschitz, in the marginals; rate for Sinkhorn convergence.

Rigollet and Stromme 2022

- **Dimension free rates** for barycentric approximation, density in L_2 and expectations of single functionals.

Carlier, Chizat, and Laborde 2022

- Mapping sending marginals to potentials Lipschitz continuous wrt $(W_2(\mu, \mu')^2 + W_2(\nu, \nu')^2)^{1/2}$.

Chiarini et al. 2022

- gradient estimates for potentials and stability in terms of KL and a negative order Sobolev norm.

Bayraktar, Eckstein, and Zhang 2022

- Quantitative stability estimates for general f -divergences.

Thanks for your attention

An incomplete list of references I

- [1] Jason Altschuler, Jonathan Weed, and Philippe Rigollet. "Near-linear time approximation algorithms for optimal transport via Sinkhorn iteration". In: *Advances in Neural Information Processing Systems*. 2017.
- [2] Erhan Bayraktar, Stephan Eckstein, and Xin Zhang. "Stability and Sample Complexity of Divergence Regularized Optimal Transport". In: *arXiv preprint arXiv:2212.00367* (2022).
- [3] Yann Brenier. "Polar decomposition and increasing rearrangement of vector-fields". In: *COMPTES RENDUS DE L ACADEMIE DES SCIENCES SERIE I-MATHEMATIQUE* 305.19 (1987), pp. 805–808.
- [4] Guillaume Carlier. "On the linear convergence of the multi-marginal Sinkhorn algorithm". In: *HAL preprint hal-03176512* (2021).
- [5] Guillaume Carlier, Lénaïc Chizat, and Maxime Laborde. "Lipschitz Continuity of the Schrödinger Map in Entropic Optimal Transport". In: *arXiv preprint arXiv:2210.00225* (2022).

An incomplete list of references II

- [6] Yongxin Chen, Tryphon Georgiou, and Michele Pavon. “Entropic and displacement interpolation: a computational approach using the Hilbert metric”. In: *SIAM Journal on Applied Mathematics* 76.6 (2016), pp. 2375–2396.
- [7] Alberto Chiarini, Giovanni Conforti, Giacomo Greco, and Luca Tamanini. “Gradient estimates for the Schrödinger potentials: convergence to the Brenier map and quantitative stability”. In: *arXiv preprint arXiv:2207.14262* (2022).
- [8] Adrien Corenflos, James Thornton, George Deligiannidis, and Arnaud Doucet. “Differentiable particle filtering via entropy-regularized optimal transport”. In: *International Conference on Machine Learning*. PMLR, 2021, pp. 2100–2111.
- [9] Marco Cuturi. “Sinkhorn distances: Lightspeed computation of optimal transport”. In: *Advances in Neural Information Processing Systems*. 2013.
- [10] Stephan Eckstein and Marcel Nutz. “Quantitative stability of regularized optimal transport and convergence of sinkhorn’s algorithm”. In: *arXiv preprint arXiv:2110.06798* (2021).

An incomplete list of references III

- [11] Nicolas Fournier and Arnaud Guillin. “On the rate of convergence in Wasserstein distance of the empirical measure”. In: *Probability Theory and Related Fields* 162.3 (2015), pp. 707–738.
- [12] Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. “Sample complexity of sinkhorn divergences”. In: *The 22nd international conference on artificial intelligence and statistics*. PMLR. 2019, pp. 1574–1583.
- [13] Promit Ghosal, Marcel Nutz, and Espen Bernton. “Stability of Entropic Optimal Transport and Schrödinger Bridges”. In: *arXiv preprint arXiv:2106.03670* (2021).
- [14] Wenbo Li and Ricardo H Nochetto. “Quantitative stability and error estimates for optimal transport plans”. In: *IMA Journal of Numerical Analysis* 41 (3 2021), pp. 1941–1965.

An incomplete list of references IV

- [15] Giulia Luise, Saverio Salzo, Massimiliano Pontil, and Carlo Ciliberto. "Sinkhorn Barycenters with Free Support via Frank-Wolfe Algorithm". In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc., 2019.
- [16] Gonzalo Mena and Jonathan Niles-Weed. "Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem". In: *Advances in Neural Information Processing Systems* 32 (2019).
- [17] Quentin Mérigot, Alex Delalande, and Frédéric Chazal. "Quantitative stability of optimal transport maps and linearization of the 2-Wasserstein space". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2020, pp. 3186–3196.
- [18] Quentin Mérigot and Edouard Oudet. "Discrete optimal transport: complexity, geometry and applications". In: *Discrete & Computational Geometry* 55.2 (2016), pp. 263–283.

An incomplete list of references V

- [19] Philippe Rigollet and Austin J Stromme. "On the sample complexity of entropic optimal transport". In: *arXiv preprint arXiv:2206.13472* (2022).
- [20] Cédric Villani. *Optimal transport: old and new*. Vol. 338. Springer, 2009.