

# Overview on optimal transport

Bernhard Schmitzer

Chamonix / Les Houches, March 2023

## 1 Introduction: Monge and Kantorovich

### 1.1 Monge formulation of OT

**Problem statement.** Given a pile of sand and a hole to fill; what is the most efficient way to fill the hole with the sand?

**Mathematical modelling of problem.**

- space  $X$  in which problem lives
- cost function  $c : X \times X \rightarrow \mathbb{R}$ ,  $c(x, y)$  is amount of work to move one unit of sand from  $x$  to  $y$
- $T : X \rightarrow X$  indicates for each position  $x$  that sand is to be sent to  $T(x)$
- how to describe pile and hole? as probability measures  $\mu, \nu \in \mathcal{P}(X)$ . mass in region  $A \subset X$ ?  $\mu(A) = \int_A d\mu$

**Examples for measures.**

- density. let  $f(x)$  be height of sand pile at  $x$ . then volume in region  $A$  given by

$$\mu(A) = \int_A f(x) dx$$

where  $dx$  denotes integration against the ‘usual’ Lebesgue measure

- Dirac measure  $\delta_x$ , unit mass at  $x$ ,

$$\delta_x(A) = \begin{cases} 1 & \text{if } x \in A, \\ 0 & \text{otherwise.} \end{cases}$$

- combination of Diracs,  $\mu := \sum_{i=1}^N m_i \delta_{x_i}$

**Push-forward.**

- if we take mass from  $x$  to  $T(x)$ , how does this transform pile  $\mu$ ?
- denote transformed measure by  $T_{\#}\mu$ : ‘push-forward of  $\mu$  under  $T$ ’
- which mass particles get mapped into  $A$ ?

$$T^{-1}(A) := \{x \in X \mid T(x) \in A\}$$

therefore,  $T_{\#}\mu(A) = \mu(T^{-1}(A))$

### Monge problem.

- look for map  $T$  that fills hole and causes least amount of work
- to fill the hole, need  $\nu = T_{\#}\mu$
- total work associated with map  $T$ :

$$\int_X c(x, T(x)) d\mu(x)$$

- so arrive at:

$$\inf \left\{ \int_X c(x, T(x)) d\mu(x) \mid T : X \rightarrow X, T_{\#}\mu = \nu \right\}$$

- main issue: feasible set is highly non-trivial, may even be empty

## 1.2 Kantorovich formulation

### Transport plans.

- new concept: transport plan  $\pi \in \mathcal{P}(X \times X)$ 
  - intuition:  $\pi(x, y)$  gives (infinitesimal) amount of mass that goes from  $x$  to  $y$
  - or:  $\pi(A \times B)$  gives amount of mass that goes from  $A \subset X$  to  $B \subset X$
- need that  $\pi$  transports  $\mu$  onto  $\nu$ , so need

$$\pi(A \times X) = \mu(A) \forall A \subset X, \quad \pi(X \times A) = \nu(A) \forall A \subset X.$$

- can write this with push-forward: let

$$p_i : X \times X \rightarrow X, \quad (x_1, x_2) \mapsto x_i$$

then find

$$p_1^{-1}(A) = \{(x, y) \in X \times X \mid p_1(x, y) \in A\} = A \times X$$

so need for all  $A \subset X$  that

$$\mu(A) = \pi(A \times X) = \pi(p_1^{-1}(A)) = p_{1\#}\pi(A)$$

and therefore  $p_{1\#}\pi = \mu$ .

- so the set of admissible transport plans can be written as

$$\Pi(\mu, \nu) := \{\pi \in \mathcal{P}(X \times X) \mid p_{1\#}\pi = \mu, p_{2\#}\pi = \nu\}$$

- $\Pi(\mu, \nu) \neq \emptyset$  since  $\mu \otimes \nu \in \Pi(\mu, \nu)$  where

$$(\mu \otimes \nu)(A \times B) := \mu(A) \cdot \nu(B).$$

### Kantorovich problem.

- cost associated with plan:

$$\int_{X \times X} c(x, y) d\pi(x, y)$$

- Kantorovich problem:

$$\inf \left\{ \int_{X \times X} c(x, y) d\pi(x, y) \mid \pi \in \Pi(\mu, \nu) \right\}$$

objective is linear, feasible set is non-empty, convex, ‘polyhedral’

**Example: discrete setting.**

- let  $X = \{x_1, \dots, x_N\}$ , discrete set of points
- identify  $\mathcal{P}(X)$  with probability simplex

$$\sigma_N := \left\{ \mu \in \mathbb{R}_+^N \mid \sum_{i=1}^N \mu_i = 1 \right\}.$$

measure belonging to vector:  $\sum_i \mu_i \delta_{x_i}$

- identify  $\Pi(\mu, \nu)$  with

$$\left\{ \pi \in \mathbb{R}_+^{N \times N} \mid \sum_j \pi_{i,j} = \mu_i \forall i, \sum_i \pi_{i,j} = \nu_j \forall j \right\}$$

for convenience introduce row- and column sum operators:

$$(P_1 \pi)_i = \sum_j \pi_{i,j}, \quad (P_2 \pi)_j = \sum_i \pi_{i,j}$$

$P_i : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^N$ , linear, can be represented as discrete matrices.

- $c : X \times X \rightarrow \mathbb{R}$  becomes matrix in  $\mathbb{R}^{N \times N}$ ,

$$\int_{X \times X} c \, d\pi = \sum_{i,j} c_{i,j} \pi_{i,j} =: \langle c, \pi \rangle$$

- arrive at finite-dimensional linear program; can in principle be solved with standard solvers; in fact: even a min cost flow problem, for which there are special, more efficient variants of the simplex algorithm

### 1.3 Kantorovich duality

**Dual problem.** We now give a formal derivation of the Kantorovich dual problem. For simplicity consider the discrete case.

- primal problem:

$$\inf_{\pi \in \mathbb{R}_+^{N \times N}} \langle c, \pi \rangle \quad \text{s.t. } P_1 \pi = \mu, P_2 \pi = \nu$$

- add Lagrange multipliers  $\phi, \psi \in \mathbb{R}^N$  for constraints:

$$= \inf_{\pi \in \mathbb{R}_+^{N \times N}} \sup_{\phi, \psi \in \mathbb{R}^N} \langle c, \pi \rangle + \langle \phi, \mu - P_1 \pi \rangle + \langle \psi, \nu - P_2 \pi \rangle$$

- duality theorem for (finite-dimensional) linear programs: can swap order of inf and sup; also: re-arrange terms

$$= \sup_{\phi, \psi \in \mathbb{R}^N} \langle \phi, \mu \rangle + \langle \psi, \nu \rangle + \inf_{\pi \in \mathbb{R}_+^{N \times N}} \langle c - P_1^\top \phi - P_2^\top \psi, \pi \rangle$$

- here use transpose (or adjoint) of  $P_i$ . find for  $P_1$ :

$$\begin{aligned} \langle P_1^\top \phi, \pi \rangle_{\mathbb{R}^{N \times N}} &:= \langle \phi, P_1 \pi \rangle_{\mathbb{R}^N} = \sum_i \phi_i (P_1 \pi)_i = \\ &= \sum_i \phi_i \sum_j \pi_{i,j} = \sum_{i,j} \phi_i \pi_{i,j} = \sum_{i,j} (P_1^\top \phi)_{i,j} \pi_{i,j} \end{aligned}$$

- now explicitly evaluate infimum over  $\pi$ , can do this entry-wise for each  $i, j$ :

$$\inf_{\pi \in \mathbb{R}_+^{N \times N}} \sum_{i,j} \hat{c}_{i,j} \pi_{i,j} = \sum_{i,j} \inf_{\pi \in \mathbb{R}_+} \hat{c}_{i,j} \cdot \pi = \begin{cases} 0 & \text{if } \hat{c}_{i,j} \geq 0, \\ -\infty & \text{else} \end{cases} = \begin{cases} 0 & \text{if } \hat{c} \geq 0, \\ -\infty & \text{else} \end{cases}$$

- so arrive at dual problem:

$$\sup_{\phi, \psi \in \mathbb{R}^N} \langle \phi, \mu \rangle + \langle \psi, \nu \rangle \quad \text{s.t.} \quad \underbrace{(\phi_i + \psi_j)}_{=:(\phi \oplus \psi)_{i,j}} \leq c_{i,j} \quad \forall i, j$$

- continuous version:

$$\sup_{\phi, \psi: X \rightarrow \mathbb{R}} \int_X \phi \, d\mu + \int_X \psi \, d\nu \quad \text{s.t.} \quad \phi(x) + \psi(y) \leq c(x, y) \quad \forall x, y$$

- A particular property of the dual problem is that it is invariant under constant shifts of the dual variables,  $(\phi, \psi) \rightarrow (\phi + \lambda, \psi - \lambda)$  for  $\lambda \in \mathbb{R}$ . The shifted dual variables are still dual feasible and have the same objective value.

### Primal-dual optimality condition.

- from previous derivation know for  $\pi \in \Pi(\mu, \nu)$  and  $\phi, \psi : X \rightarrow \mathbb{R}$  with  $\phi \oplus \psi \leq c$  that

$$\int c \, d\pi \geq \int \phi \, d\mu + \int \psi \, d\nu$$

with equality if and only if  $\pi$  is optimal (primal) plan and  $(\phi, \psi)$  are dual optimal. so we have:

$$0 \leq \int c \, d\pi - \int \phi \underbrace{d\mu}_{=dP_1\pi} - \int \psi \underbrace{d\nu}_{=dP_2\pi} = \int_{X \times X} \underbrace{[c(x, y) - \phi(x) - \psi(y)]}_{\geq 0} \underbrace{d\pi(x, y)}_{\geq 0}$$

- so equality if and only if  $c(x, y) = \phi(x) + \psi(y)$   $\pi(x, y)$ -almost everywhere; in discrete case:

$$[\pi_{i,j} > 0] \quad \Leftrightarrow \quad [c_{i,j} = \phi_i + \psi_j]$$

### $c$ -concave functions.

- In dual problem, for fixed  $\psi$ , find best  $\phi$ . Since  $\mu \geq 0$ , try to make  $\phi$  at each point as large as possible without violating the constraint  $c \geq \phi \oplus \psi$ . So set

$$\phi(x) = \inf_{y \in X} c(x, y) - \psi(y) =: \psi^c(x).$$

This is called  $c$ -transform (assume here for simplicity that  $c$  is symmetric, otherwise need to define ‘forward and backwards transform’ separately).

- A function  $\phi$  that can be written as  $\phi = \psi^c$  is called  $c$ -concave.  $c$ -concavity is a strong structural property and various ‘big’ theorems in OT are derived from special properties of  $c$ -concave functions. For  $c(x, y) = \|x - y\|^2$ ,  $\phi$  is  $c$ -concave, iff  $x \mapsto \|x\|^2 - \phi(x)$  is convex.
- Alternating maximization of the dual potentials  $\phi$  and  $\psi$  is easy, but not a good optimization algorithm. In fact,  $\phi^{ccc} = \phi^c$ , so the iterations become stationary after three iterations, without guarantee of optimality. The auction and Sinkhorn algorithms can be interpreted as fixes of this issue, such that alternating dual maximization becomes an efficient (approximate) algorithm.

## Differentiability of Kantorovich cost.

- Consider function  $C : \mathcal{P}(X) \times \mathcal{P}(X) \rightarrow \mathbb{R}$ ,  $(\mu, \nu) \mapsto \inf_{\pi \in \Pi(\mu, \nu)} \int c, d\pi$ . Can also be written via dual:

$$C(\mu, \nu) = \sup_{\substack{\phi, \psi: X \rightarrow \mathbb{R}, \\ \phi \oplus \psi \leq c}} \int \phi d\mu + \int \psi d\nu$$

- This is a supremum over a collection of linear functions, parametrized by slopes  $(\phi, \psi)$ . Convex analysis:  $C$  is lower semi-continuous and sub-gradients are given by optimal  $(\phi, \psi)$ .
- For  $(\mu, \nu)$  let  $(\phi, \psi)$  be optimal duals. Then

$$C(\mu + \delta\mu, \nu) \geq C(\mu, \nu) + \int \phi d\delta\mu.$$

This can in principle be used to predict small changes in transport cost; but careful: this holds for any dual optimal  $\phi$ .

- So if  $\delta\mu$  is not mean zero, can just add big  $\lambda$  to  $\phi$  and make this as large as we want. Consistent:  $C(\mu + \delta\mu, \nu) = \infty$  if  $\mu + \delta\mu \notin \mathcal{P}(X)$ : there is no feasible transport plan.
- If  $\delta\mu$  has zero mean and  $\phi$  is unique up to constant shifts (e.g. in setting of Benier theorem) then this yields a first order estimate. But no simple regularity theory as for finite-dimensional differentiable functions.

## 1.4 Relation between Monge and Kantorovich problem; Brenier's theorem

### Kantorovich as relaxation of Monge

- for  $\mu, \nu \in \mathcal{P}(X)$ , assume  $T : X \rightarrow X$  is such that  $T_{\#}\mu = \nu$ , then define measure by

$$\pi := (\text{id}, T)_{\#}\mu \quad \text{where} \quad (\text{id}, T) : X \rightarrow X \times X, \quad x \mapsto (x, T(x))$$

- find:  $p_{1\#}\pi = p_{1\#}(\text{id}, T)_{\#}\mu = [p_1 \circ (\text{id}, T)]_{\#}\pi = \text{id}_{\#}\mu = \mu$   
likewise:  $p_{2\#}\pi = p_{2\#}(\text{id}, T)_{\#}\mu = [p_2 \circ (\text{id}, T)]_{\#}\pi = T_{\#}\mu = \nu$   
so  $\pi \in \Pi(\mu, \nu)$ .

- compare transport costs:

$$\int_X c(x, T(x)) d\mu(x) = \int_X [c \circ (\text{id}, T)](x) d\mu(x) = \int_{X \times X} c d[(\text{id}, T)_{\#}\mu](x, y) = \int_{X \times X} c d\pi$$

here we used the change of variables formula for push-forward measures

- so each Monge transport map (not necessarily optimal) induces a Kantorovich transport plan with the same cost. we find:

$$\inf_{\substack{T: X \rightarrow \mathbb{R}, \\ T_{\#}\mu = \nu}} \int_X c(x, T(x)) d\mu(x) \geq \inf_{\pi \in \Pi(\mu, \nu)} \int_{X \times X} c(x, y) d\pi(x, y)$$

this inequality can be strict, in particular if there are no feasible transport maps.

### Brenier's theorem.

- In special cases it can be shown that the above inequality is in fact an equality and that the optimal Kantorovich plan is indeed induced by an optimal Monge map. For the special case of  $X \subset \mathbb{R}^d$  with  $c(x, y) = \|x - y\|^2$  this is called Brenier's theorem.
- Assume that  $\mu$  has a Lebesgue density (a slightly weaker condition would be sufficient) and that  $\mu$  and  $\nu$  'decay fast enough to zero as  $\|x\| \rightarrow +\infty$ '. Then the minimizing  $\pi$  in the Kantorovich problem is unique. It has the form  $\pi = (\text{id}, T)_{\#}\mu$  for a map  $T : X \rightarrow X$  with  $T_{\#}\mu = \nu$ .

- In addition,  $T$  is the gradient of a convex potential  $\phi : X \rightarrow \mathbb{R}$ , i.e.  $T = \nabla\phi$ . (At least Lebesgue-almost everywhere, and  $\phi$  may be non-differentiable on a small set.)
- This can be proved via the above primal-dual relation and indeed the convex potential  $\phi$  is very closely related to the dual Kantorovich potential  $\hat{\phi}$  in the previous section via  $\phi(x) = \frac{1}{2}(\|x\|^2 - \hat{\phi}(x))$ .

### Semi-discrete transport.

- As above, let  $X \subset \mathbb{R}^d$ ,  $c(x, y) = \|x - y\|^2$ , let  $\mu = f \cdot \mathcal{L}$  (where  $\mathcal{L}$  denotes the Lebesgue measure) and  $\nu = \sum_i m_i \delta_{y_i}$ .
- By Brenier's theorem, the optimal  $\pi$  is unique and induced by a map  $T$ ,  $\pi = (\text{id}, T)_\# \mu$ . Clear:  $T$  must be piecewise constant, mapping regions  $C_i \subset X$  to the mass locations  $y_i$ .
- This might model assignments of students (a lot of them, essentially 'continuously' distributed over the country or city) to schools (discrete centers), or customers to supermarkets and similar situations with continuous demand and discrete supply.
- Locally, at first, customers simply want to go to the nearest market, where  $c(x, y_i) = \min_j c(x, y_j)$ . So customers in

$$V_i := \{x \in X | c(x, y_i) = \min_j c(x, y_j)\}$$

would all go to market  $i$ . These are called Voronoi cells. But then the demand  $\mu(V_i)$  might not match the supply  $m_i$ . Economically, we need to introduce prices to act as additional incentive, in addition to travel distance, to influence customers decisions. If  $\mu(V_i) > m_i$ , then the market  $i$  should raise its prices to 'repulse' some customers, if  $\mu(V_i) < m_i$ , then lower the prices to attract near customers from other nearby markets. OT provides a natural description of this phenomenon.

- By the primal-dual optimality relation, one has  $c = \phi \oplus \psi$   $\pi$ -almost everywhere. The optimal  $\phi$  can be written as  $c$ -transform of  $\psi$ , so get

$$\phi(x) = \psi^c(x) = \min_i c(x, y_i) - \psi_i.$$

For a given  $x$ , mass can only be transported to some  $y_i$  that is minimizing in this expression. I.e. location  $i$  can only receive mass from within the region

$$C_i(\psi) := \{x \in X | c(x, y_i) - \psi_i = \min_j c(x, y_j) - \psi_j\}.$$

Now we see that  $-\psi$  acts as the aforementioned price.

- For the squared distance, the above condition can be written as

$$x \in C_i(\psi) \quad \text{iff} \quad \langle x, y_i \rangle - \frac{1}{2}(\|y_i\|^2 - \psi_i) = \max_j \langle x, y_j \rangle - \frac{1}{2}(\|y_j\|^2 - \psi_j)$$

. From this we deduce that the boundaries between adjacent cells are straight lines and that in fact each cell is a convex polytope.

- Eliminating  $\phi$  via the  $c$ -transform, one can obtain a finite-dimensional unconstrained dual problem in terms of  $\psi$ :

$$\sup_{\psi \in \mathbb{R}^N} J(\psi) \quad \text{where} \quad J(\psi) := \int \psi^c d\mu + \sum_i \psi_i m_i$$

The integral can be written as

$$\sum_i \int_{C_i(\psi)} [\|x - y_i\|^2 - \psi_i] d\mu(x).$$

One can show that  $J$  is differentiable and

$$\partial_{\psi_i} J(\psi) = - \int_{C_i(\psi)} d\mu + m_i.$$

This corresponds nicely with the economic interpretation: if the demand  $\mu(C_i)$  is lower than the supply  $m_i$ , the partial derivative for  $\psi_i$  is positive.

- Note that if we add the same constant  $\lambda$  to all prices, then all  $C_i$  remain the same. Economically, this invariance is not so realistic. Eventually, customers will stop going to the market entirely and try to find alternatives. This setting might be modelled more appropriately with unbalanced transport.

## 2 Wasserstein distance

### 2.1 Definition and metric axioms

**Definition.**

- Let  $(X, d)$  be a metric space, e.g.  $X \subset \mathbb{R}^d$ ,  $d(x, y) = \|x - y\|^2$  or a curved surface.
- Let  $p \in [1, \infty)$ ,  $\mu, \nu \in \mathcal{P}(X)$ . Set:

$$W_p(\mu, \nu) := \inf \left\{ \int_{X \times X} d^p(x, y) d\pi(x, y) \mid \pi \in \Pi(\mu, \nu) \right\}^{1/p}$$

- claim:  $W_p$  is a metric on  $\mathcal{P}(X)$ , called Wasserstein distance (technically: a metric on probability measures with finite  $p$ -th moment)

**Symmetry, definiteness.**

- Since  $d(x, y) \geq 0$  and  $\pi \geq 0$ , one has  $W_p(\mu, \nu) \geq 0$ .
- Symmetry,  $W_p(\mu, \nu) = W_p(\nu, \mu)$ : let  $\pi \in \Pi(\mu, \nu)$ , set  $\hat{\pi} := (p_2, p_1)_\# \pi$  and find  $\hat{\pi} \in \Pi(\nu, \mu)$ . Further:

$$\int d d\hat{\pi} = \int d \circ (p_2, p_1) d\pi = \int d d\pi$$

where we used the symmetry of  $d$ . Consequently, for any plan in  $\Pi(\mu, \nu)$ , the ‘reversed plan’ is in  $\Pi(\nu, \mu)$ , has the same transport cost, and vice versa. Hence, both problem have the same infimal value.

- Definiteness, part 1:  $W_p(\mu, \mu) = 0$ . In this case,  $\pi = (\text{id}, \text{id})_\# \mu$  is a feasible transport plan. We get:

$$0 \leq W_p(\mu, \mu) \leq \int_{X \times X} d d\pi = \int_X d \circ (\text{id}, \text{id}) d\mu = \int_X d(x, x) d\mu(x) = 0$$

- Definiteness, part 2:  $[W_p(\mu, \nu) = 0] \Rightarrow [\mu = \nu]$ . Let  $\pi$  be an optimal plan for  $W_p(\mu, \nu)$  (do not address existence here). Then:

$$0 = \int_{X \times X} d(x, y)^p d\pi(x, y)$$

Since  $d \geq 0$  and  $\pi \geq 0$ , this is only possible if  $d(x, y) = 0$   $\pi(x, y)$ -almost everywhere. So  $\pi$  must live exclusively on the ‘diagonal’ and thus have the form  $\pi = (\text{id}, \text{id})_\# \rho$  for some  $\rho \in \mathcal{P}(X)$ . Now use  $\pi \in \Pi(\mu, \nu)$ :  $\mu = p_{1\#} \pi = \rho = p_{2\#} \pi = \nu$ .

### 2.2 Triangle inequality

The metric  $d$  must satisfy

$$d(x, y) + d(y, z) \geq d(x, z) \quad \text{for all } x, y, z \in X.$$

The same must hold for  $W_p$  on triplets of probability measures  $\mu, \nu, \rho \in \mathcal{P}(X)$ .

**Simple proof for transport maps.**

- We sketch the proof for the simple case where the optimal plans for  $W_p(\mu, \nu)$  and  $W_p(\nu, \rho)$  are induced by maps  $S, T : X \rightarrow X$ ,  $S_\# \mu = \nu$ ,  $T_\# \nu = \rho$ .



- In this case,  $(T \circ S)_{\#}\mu = \rho$  and so  $T \circ S$  is a feasible (not necessarily optimal) transport map from  $\mu$  to  $\rho$  for  $W_p(\mu, \rho)$ . Therefore we get

$$\begin{aligned}
W_p(\mu, \rho) &\leq \left[ \int d(x, T(S(x)))^p d\mu(x) \right]^{1/p} \\
&\leq \left[ \int [d(x, S(x)) + d(S(x), T(S(x)))]^p d\mu(x) \right]^{1/p} \\
&\leq \left[ \int d(x, S(x))^p d\mu(x) \right]^{1/p} + \left[ \int d(S(x), T(S(x)))^p d\mu(x) \right]^{1/p} \\
&= W_p(\mu, \nu) + \left[ \int d(x, T(x))^p d(S_{\#}\mu)(x) \right]^{1/p} = W_p(\mu, \nu) + W_p(\nu, \rho)
\end{aligned}$$

### Gluing lemma.

- The proof strategy also works when the optimal plans are not map-like. Will not go through all steps, but show how two plans  $\pi \in \Pi(\mu, \nu)$ ,  $\lambda \in \Pi(\nu, \rho)$  can be concatenated.
- We will construct a measure  $\eta \in \mathcal{P}(X \times X \times X)$  with

$$(p_1, p_2)_{\#}\eta = \pi, \quad (p_2, p_3)_{\#}\eta = \lambda.$$

$\eta(x, y, z)$  can intuitively be interpreted as (infinitesimal) mass going from  $x$  to  $z$  via  $y$ .

- For simplicity, do this in discrete setting, but same idea works in continuum with slightly more complex notation. Denote mass at individual points by  $\mu(x)$ ,  $\pi(x, y)$ , et cetera. Then  $\pi \in \Pi(\mu, \nu)$  implies  $\sum_y \pi(x, y) = \mu(x)$  and so forth.
- Consider mass arriving in  $y$  after being transported according to  $\pi$ . How should we send the mass along? Use ‘conditional probability’

$$\lambda_y(z) := \frac{\lambda(y, z)}{\nu(y)}$$

to distribute mass. (In the continuous setting, this conditional probability is provided by the disintegration.)

- Then we set

$$\eta(x, y, z) := \pi(x, y) \frac{\lambda(y, z)}{\nu(y)} = \frac{\pi(x, y)\lambda(y, z)}{\nu(y)}$$

(Be a little careful when  $\eta(y) = 0$ . In this case no mass will arrive at  $y$  and no mass will leave. So simply set  $\eta(x, y, z) = 0$ .)

- Then find

$$\sum_z \eta(x, y, z) = \sum_z \frac{\pi(x, y)\lambda(y, z)}{\nu(y)} = \frac{\pi(x, y)}{\nu(y)} \sum_z \lambda(y, z) = \pi(x, y).$$

Likewise, find  $\sum_x \eta(x, y, z) = \lambda(y, z)$  and thus  $\sum_y \eta(\cdot, y, \cdot) \in \Pi(\mu, \rho)$ .

- Some things to take away:

- this construction also works in continuum
- evaluating Wasserstein cost of  $\eta$  yields triangle inequality, similar to above
- perfectly reasonable to have ‘transport plans’ between more than two marginals

### 2.3 Shortest paths

For  $x, y \in X$ , a curve  $\gamma : [0, 1] \rightarrow X$  is called (constant speed) geodesic from  $x$  to  $y$  if

$$\gamma(0) = x, \quad \gamma(1) = y, \quad d(\gamma(s), \gamma(t)) = |s - t| \cdot d(x, y) \quad \forall s, t \in [0, 1].$$

We will call  $(X, d)$  geodesic, if such curves exist for all pairs  $(x, y)$ . We will now show: If  $(X, d)$  is geodesic, then so is  $(\mathcal{P}(X), W_p)$ .

### Constructing geodesics.

- For simplicity, assume  $X \subset \mathbb{R}^d$ ,  $d(x, y) = \|x - y\|$ , and let the optimal plan  $\pi$  be induced by a map  $T$ ,  $\pi = (\text{id}, T)_{\#}\mu$ .
- Intuition for geodesic in  $\mathcal{P}(X)$  from  $\mu$  to  $\nu$ : each mass particle moves along geodesic in  $X$  from initial to final position. Particle at  $x$  moves to  $T(x)$ , so it will move on curve

$$\gamma(\cdot, x) : [0, 1] \rightarrow X, \quad \gamma(t, x) := (1 - t) \cdot x + t \cdot T(x)$$

Taking each particle from its initial position to  $\gamma(t, \cdot)$  means we apply the push-forward of  $\gamma(t, \cdot)$  to  $\mu$ . So our conjectured geodesic takes the form

$$\rho : [0, 1] \rightarrow \mathcal{P}(X), \quad \rho(t) := \gamma(t, \cdot)_{\#}\mu.$$

It is easy to verify:  $\rho(0) = \mu$ ,  $\rho(1) = \nu$ .

### Estimating distances along geodesic.

- How do we estimate the distance between  $\rho(0) = \mu$  and  $\rho(s)$ ? Need to guess a transport map. Clear: particle from  $x$  moves to  $\gamma(t, x)$ , so use this as transport map candidate:  $T_{0,s} := \gamma(t, \cdot)$ . We find:

$$W_p(\mu, \rho(s))^p \leq \int \underbrace{\|T_{0,s}(x) - x\|^p}_{=s\|T(x)-x\|^p} d\mu(x) = s^p \int \|T(x) - x\|^p d\mu(x) = s^p W_p(\mu, \nu)^p$$

- How about optimal map from  $\rho(s)$  to  $\rho(t)$ ? Mass from  $\gamma(s, x)$  moves to  $\gamma(t, x)$ . For  $s < 1$  can show:  $\gamma(s, \cdot)$  is invertible. So use as candidate:  $T_{s,t} := \gamma(t, \cdot) \circ \gamma(s, \cdot)^{-1}$ . Find:

$$\begin{aligned} W_p(\rho(s), \rho(t))^p &\leq \int \|T_{s,t}(x) - x\|^p d\rho(s) = \int \|\gamma(t, \cdot) \circ \gamma(s, \cdot)^{-1} - \text{id}\|^p d(\gamma(s, \cdot)_{\#}\mu) \\ &= \int \|\gamma(t, \cdot) - \gamma(s, \cdot)\|^p d\mu = |t - s|^p \int \|T - \text{id}\|^p d\mu = |t - s|^p W_p(\mu, \nu)^p. \end{aligned}$$

- If any of these inequalities were strict, we would violate the triangle inequality. The above inequalities provide

$$W_p(\mu, \rho(s)) + W_p(\rho(s), \rho(t)) + W_p(\rho(t), \nu) \leq W_p(\mu, \nu)$$

whereas the triangle inequality gives the opposite inequality. Thus, all inequalities above must be equalities.

- The proof strategy generalizes to non-map transport plans and more general metric spaces.

## 2.4 Continuity equation and Benamou–Brenier formula

### Continuity equation.

- Let  $X \subset \mathbb{R}^d$ ,  $d(x, y) = \|x - y\|$ , set  $p = 2$  for simplicity. Consider a curve of measures  $\rho(t) := \gamma(t, \cdot)_{\#}\mu$ , assume  $\gamma$  differentiable in time,  $\gamma(t, \cdot)$  invertible.
- A particle starting at  $t = 0$  at  $x$  will move with (Lagrangian) velocity  $\dot{\gamma}(t, x)$ . But at time  $t$  is at  $\gamma(t, x)$ , so an external observer (who cannot distinguish the mass particles) will see the (Eulerian) velocity field

$$v(t, \cdot) := \dot{\gamma}(t, \cdot) \circ \gamma(t, \cdot)^{-1}.$$

- If all mass particles of  $\rho$  follow a (Eulerian) velocity field  $v$ ,  $(\rho, v)$  will satisfy the continuity equation

$$\partial_t \rho + \text{div}(v \cdot \rho) = 0,$$

in our case with the boundary conditions  $\rho(0) = \mu$  and  $\rho(1) = \nu = \gamma(1, \cdot)_{\#}\mu$ .

- This equation is to be understood in a distributional sense. This means that for every differentiable test function  $\phi \in C^1([0, 1] \times X)$  one imposes that

$$\int_0^1 \int_X \partial_t \phi(t, \cdot) d\rho(t) dt + \int_0^1 \int_X \nabla \phi(t, \cdot) \cdot v(t, \cdot) d\rho(t) dt = \int_X \phi(1, \cdot) d\nu - \int_X \phi(0, \cdot) d\mu$$

- Now we show that  $\rho(t) := \gamma(t, \cdot) \# \mu$  and  $v(t, \cdot) := \dot{\gamma}(t, \cdot) \circ \gamma(t, \cdot)^{-1}$  solve the continuity equation. Set  $F(t) := \int_X \phi(t, \cdot) d\rho(t)$ . Then

$$\int_0^1 \left[ \frac{d}{dt} F(t) \right] dt = F(1) - F(0)$$

and

$$F(1) = \int \phi(1, \cdot) d\nu, \quad F(0) = \int \phi(0, \cdot) d\mu.$$

For  $F(t)$  one has

$$F(t) = \int_X \phi(t, \cdot) d(\gamma(t, \cdot) \# \mu) = \int_X \phi(t, \gamma(t, \cdot)) d\mu$$

and so

$$\begin{aligned} \frac{d}{dt} F(t) &= \int_X [\partial_t \phi(t, \gamma(t, \cdot)) + \nabla \phi(t, \gamma(t, \cdot)) \cdot \dot{\gamma}(t, \cdot)] d\mu \\ &= \int_X \left[ \partial_t \phi(t, \cdot) + \nabla \phi(t, \cdot) \cdot \underbrace{\dot{\gamma}(t, \gamma(t, \cdot)^{-1})}_{=v(t, \cdot)} \right] d\rho(t) \end{aligned}$$

### Benamou–Brenier formula.

- Geodesics in  $(\mathcal{P}(\mathbb{R}^d), W_2)$  are of the above form and therefore solve the continuity equation. Question: which of the solutions of the continuity equation yields the shortest path? Answer: the one with the lowest kinetic energy (more precisely: action).
- For a pair  $(\rho, v)$  that solves the continuity equation, set

$$BB(\rho, v) := \int_0^1 \int_X \|v(t, \cdot)\|^2 d\rho(t) dt.$$

Then one has

$$W_2(\mu, \nu)^2 = \inf \{ BB(\rho, v) \mid (\rho, v) \text{ solve CE between } \mu \text{ and } \nu \}$$

- Sketch of inequality  $BB \leq W$ : take  $(\rho, v)$  generated by shortest path. Recall:  $\gamma(t, \cdot) = (1-t) \text{id} + t \cdot T$ . Lagrangian velocity:  $\dot{\gamma}(t, \cdot) = T - \text{id}$ . Have already shown that  $(\rho, v)$  solve the continuity equation. Now plug into  $BB$  to get:

$$\begin{aligned} \inf BB &\leq BB(\rho, v) = \int_0^1 \int_X \|v(t, \cdot)\|^2 d\rho(t) dt = \int_0^1 \int_X \|\dot{\gamma}(t, \cdot) \circ \gamma(t, \cdot)^{-1}\|^2 d\gamma(t, \cdot) \# \mu dt \\ &= \int_0^1 \int_X \|\dot{\gamma}(t, \cdot)\|^2 d\mu dt = \int_0^1 \int_X \|T - \text{id}\|^2 d\mu dt = W_2(\mu, \nu)^2 \end{aligned}$$

- For converse inequality, for given  $(\rho, v)$  follow the individual mass particles. Let  $\gamma(t, x)$  be the solution of the following ODE:

$$\gamma(0, x) = x, \quad \dot{\gamma}(t, x) = v(t, \gamma(t, x))$$

This means, we recover the Lagrangian coordinate from the Eulerian velocity field. (Mathematically this only works if  $v$  is sufficiently regular.) Can then show with similar calculations as above:  $\rho(t) = \gamma(t, \cdot) \# \rho(0) = \gamma(t, \cdot) \# \mu$ . Now look at BB functional:

$$\begin{aligned}
BB(\rho, v) &= \int_0^1 \int_X \|v(t, \cdot)\|^2 d\rho(t) dt = \int_0^1 \int_X \|v(t, \gamma(t, \cdot))\|^2 d\mu dt \\
&= \int_X \int_0^1 \|v(t, \gamma(t, \cdot))\|^2 dt d\mu = \int_X \int_0^1 \|\dot{\gamma}(t, \cdot)\|^2 dt d\mu \geq \int_X \left\| \int_0^1 \dot{\gamma}(t, \cdot) dt \right\|^2 d\mu \\
&\geq \int_X \|\gamma(1, \cdot) - \gamma(0, \cdot)\|^2 d\mu \geq W_2(\mu, \nu)^2
\end{aligned}$$

where we used Jensen's inequality when pulling out  $\|\cdot\|^2$  from the integral, and the fact that  $\gamma(0, \cdot) = \text{id}$  and that  $\gamma(1, \cdot)$  is a feasible transport map from  $\mu$  to  $\nu$ .

### 3 Entropic regularization and Sinkhorn algorithm

#### 3.1 Entropic regularization

**KL-divergence / relative entropy.**

- For  $\rho, \sigma \in \mathcal{P}(X)$  set

$$\text{KL}(\rho|\sigma) := \begin{cases} \int \varphi\left(\frac{d\rho}{d\sigma}\right) d\sigma & \text{if } \rho \ll \sigma, \\ +\infty & \text{else.} \end{cases},$$

where

$$\varphi(s) := \begin{cases} s \log(s) - s + 1 & \text{if } s > 0, \\ 1 & \text{if } s = 0, \\ +\infty & \text{else.} \end{cases}$$

and similarly on other spaces.

- One has KL is jointly convex in both of its arguments,  $\text{KL}(\cdot|\sigma)$  is strictly convex,  $\text{KL}(\rho|\sigma) \geq 0$  with equality if and only if  $\rho = \sigma$ .

**Regularized primal problem.**

- Now, for some (small) parameter  $\varepsilon \geq 0$ , add  $\varepsilon \cdot \text{KL}(\pi|\mu \otimes \nu)$  to primal Kantorovich objective. The regularized problem is given by:

$$\inf \left\{ \int c \, d\pi + \varepsilon \text{KL}(\pi|\mu \otimes \nu) \mid \pi \in \Pi(\mu, \nu) \right\}$$

In principle other choices instead of  $\mu \otimes \nu$  are possible for the reference measure, but they can all be reduced to this case.

- Minimizers exist under reasonable assumptions. The minimizer is unique for  $\varepsilon > 0$ . In reasonable settings, one has  $\pi_\varepsilon \rightarrow \pi$  as  $\varepsilon \rightarrow 0$  where  $\pi$  is (some) minimizer of the unregularized problem.

**Motivation.** Considering the regularized problem can be useful for various reasons:

- Relation to modelling stochastic processes: optimal  $\pi$  becomes ‘blurry’ and non-deterministic, even in cases where optimal unregularized solution would be given by a Monge map.
- Statistical reasons: if true  $\mu$  and  $\nu$  are unknown, but only finite samples are available, then can show that optimal value of entropic problem converges faster to true value, compared to the unregularized problem; in particular if  $X$  is high-dimensional.
- Uniqueness of optimal  $\pi$  and differentiability of optimal value: when OT is used in a bigger data processing pipeline, it is reassuring to know that optimal transport plan is unique and that gradient-based optimization can be used.
- Simple numerical method: for  $\varepsilon > 0$  the problem can be solved with Sinkhorn’s algorithm, which is particularly simple and also apt for parallelization on GPUs.

**Derivation of dual problem.**

- for simplicity, we will do this again to the discrete setting
- as before, argue via Lagrangian: primal problem is equivalent to

$$\inf_{\pi \in \mathbb{R}_+^{N \times N}} \sup_{\phi, \psi \in \mathbb{R}^N} \langle c, \pi \rangle + \varepsilon \text{KL}(\pi|\mu \otimes \nu) + \langle \phi, \mu - P_1 \pi \rangle + \langle \psi, \nu - P_2 \pi \rangle$$

- swap order of minimization, reorder terms

$$\sup_{\phi, \psi \in \mathbb{R}^N} \langle \phi, \mu \rangle + \langle \psi, \nu \rangle + \inf_{\pi \in \mathbb{R}_+^{N \times N}} \left[ \langle c - P_1^\top \phi - P_2^\top \psi, \pi \rangle + \varepsilon \text{KL}(\pi | \mu \otimes \nu) \right]$$

- since KL is acting ‘entry-wise’ on  $\pi$ , the min can be performed for each entry of  $\pi$  separately. let us solve the following sub-problem:

$$\inf_{s \geq 0} \hat{c} \cdot s + \varepsilon \varphi(s/t) \cdot t$$

where  $t$  takes the role of  $\mu_i \cdot \nu_j$ .

- try first order optimality condition:

$$0 = \partial_s [\hat{c} \cdot s + \varepsilon \varphi(s/t) \cdot t] = \hat{c} + \varepsilon \log(s/t) \quad \Rightarrow \quad s = \exp(-\hat{c}/\varepsilon) \cdot t > 0$$

this value lies in the region where  $\varphi'$  is defined. by strict convexity of  $\varphi$  this must be the unique minimizer. we get:

$$\begin{aligned} \inf_{s \geq 0} \hat{c} \cdot s + \varepsilon \underbrace{\varphi(s/t) \cdot t}_{=s \log(s/t) - s + t} &= \hat{c} \exp(-\hat{c}/\varepsilon) \cdot t + \varepsilon [\exp(-\hat{c}/\varepsilon) \cdot (-\hat{c}/\varepsilon) - \exp(-\hat{c}/\varepsilon) + 1] \cdot t \\ &= -\varepsilon \cdot [\exp(-\hat{c}/\varepsilon) - 1] \cdot t \end{aligned}$$

- back to full problem, arrive at regularized dual:

$$(\dots) = \sup_{\phi, \psi \in \mathbb{R}^N} \langle \phi, \mu \rangle + \langle \psi, \nu \rangle - \varepsilon \sum_{i,j} \left[ \exp\left(-\frac{c_{i,j} - \phi_i - \psi_j}{\varepsilon}\right) - 1 \right] \mu_i \nu_j$$

- discussion: the term  $-\varepsilon \exp\left(-\frac{c_{i,j} - \phi_i - \psi_j}{\varepsilon}\right)$  acts like a smooth approximation of the constraint  $c_{i,j} - \phi_i - \psi_j \geq 0$ . if the constraint is violated, the term tends to  $+\infty$  as  $\varepsilon \rightarrow 0$ . if the constraint is satisfied, the penalty tends to 0.
- the last term,  $-\varepsilon \cdot (-1)$  is constant and tends to 0 as  $\varepsilon \rightarrow 0$ .
- so we have a smooth, unconstrained approximation of the original dual problem
- observe: still have the invariance under constant shifts of  $\phi$  and  $\psi$ , but by convexity of exp the objective is strictly concave up to these constant shifts, i.e. dual maximizers are unique up to these shifts

### Primal-dual optimality condition.

- For the unregularized problem we obtained the primal-dual optimality condition

$$c - \phi \oplus \psi = 0 \quad \pi\text{-almost everywhere.}$$

Now study the equivalent relation for the regularized case. We will show  $\pi \in \Pi(\mu, \nu)$  and  $\phi, \psi : X \rightarrow \mathbb{R}$  are primal-dual optimal if and only if

$$\pi = \exp\left(\frac{-c + \phi \oplus \psi}{\varepsilon}\right) \cdot \mu \otimes \nu.$$

- For simplicity, again we consider only the discrete setting. Recall the derivation of the regularized dual problem, when we explicitly minimized over the entries of  $\pi$ . We obtained:

$$s \cdot \hat{c} + \varepsilon \varphi(s/t) \cdot t \geq -\varepsilon [\exp(-\hat{c}/\varepsilon) - 1] \cdot t$$

where  $s = \pi_{i,j}$ ,  $t = \mu_i \nu_j$ , with equality if and only if  $s = \exp(-\hat{c}/\varepsilon) \cdot t$ .

- Now apply this to the primal dual gap:

$$[\langle c, \pi \rangle + \varepsilon \text{KL}(\pi | \mu \otimes \nu)] - \left[ \langle \phi, \mu \rangle + \langle \psi, \nu \rangle - \varepsilon \sum_{i,j} \left[ \exp\left(\frac{-c_{i,j} + \phi_i + \psi_j}{\varepsilon}\right) - 1 \right] \mu_i \nu_j \right]$$

(now use  $\pi \in \Pi(\mu, \nu)$ , i.e.  $P_1 \pi = \mu, \dots$ )

$$= \sum_{i,j} \left[ [(c_{i,j} - \phi_i - \psi_j) \cdot \pi_{i,j} + \varepsilon \varphi(\pi_{i,j} / \mu_i \nu_j) \cdot \mu_i \nu_j] + \varepsilon \left[ \exp\left(\frac{-c_{i,j} + \phi_i + \psi_j}{\varepsilon}\right) - 1 \right] \mu_i \nu_j \right]$$

(now for each  $i, j$  the corresponding term is of the above form, so we get:)

$$\geq 0$$

with equality if and only if  $\pi_{i,j} = \exp\left(\frac{-c_{i,j} + \phi_i + \psi_j}{\varepsilon}\right)$  for all  $i, j$ .

### 3.2 Sinkhorn algorithm

**Derivation as alternating dual maximization.**

- Consider now the dual objective for fixed  $\psi$  and optimize over  $\phi$ . The objective can be written as:

$$\sum_i \left[ \mu_i \cdot \phi_i - \varepsilon \mu_i \exp(\phi_i / \varepsilon) \cdot \sum_j \exp\left(-\frac{c_{i,j} - \psi_j}{\varepsilon}\right) \nu_j \right] + \langle \psi, \nu \rangle + \varepsilon \sum_{i,j} \mu_i \nu_j$$

- So we can optimize over each  $\phi_i$  individually. Take partial derivative and set to zero:

$$0 = \mu_i \left[ 1 - \exp(\phi_i / \varepsilon) \cdot \sum_j \exp\left(-\frac{c_{i,j} - \psi_j}{\varepsilon}\right) \nu_j \right]$$

- Resolve for  $\phi$ , analogous formula for optimization over  $\psi$ :

$$\phi_i = -\varepsilon \log \left[ \sum_j \exp\left(-\frac{c_{i,j} - \psi_j}{\varepsilon}\right) \nu_j \right], \quad \psi_j = -\varepsilon \log \left[ \sum_i \exp\left(-\frac{c_{i,j} - \phi_i}{\varepsilon}\right) \mu_i \right]$$

- Now, if we start with some initial  $\phi^{(0)}, \psi^{(0)}$ , then generate  $\phi^{(1)}$  by optimizing over  $\phi$ , then  $\psi^{(1)}$  by optimizing over  $\psi$  and keep on iterating, the update rule is given by:

$$\phi_i^{(\ell)} := -\varepsilon \log \left[ \sum_j \exp\left(-\frac{c_{i,j} - \psi_j^{(\ell-1)}}{\varepsilon}\right) \nu_j \right], \quad \psi_j^{(\ell)} := -\varepsilon \log \left[ \sum_i \exp\left(-\frac{c_{i,j} - \phi_i^{(\ell)}}{\varepsilon}\right) \mu_i \right]$$

for  $\ell \geq 1$ .

**Reformulation with scaling factors.**

- Define the matrix  $k \in \mathbb{R}_{++}^{N \times N}$  via  $k_{i,j} := \exp(-c_{i,j}/\varepsilon)$ . Introduce the scaling factors  $u^{(\ell)}, v^{(\ell)} \in \mathbb{R}_+^N$  via

$$u_i^{(\ell)} := \exp(\phi_i^{(\ell)} / \varepsilon), \quad v_j^{(\ell)} := \exp(\psi_j^{(\ell)} / \varepsilon).$$

- Then the above iterations for  $\phi^{(\ell)}$  and  $\psi^{(\ell)}$  can be equivalently rewritten as

$$u_i^{(\ell)} := \frac{1}{\sum_j k_{i,j} v_j^{(\ell-1)} \nu_j}, \quad v_j^{(\ell)} := \frac{1}{\sum_i k_{i,j} u_i^{(\ell)} \mu_i}.$$

- This can be compactly written as

$$u^{(\ell)} := \frac{1}{k \cdot \text{diag}(\nu) \cdot v^{(\ell-1)}}, \quad v^{(\ell)} := \frac{1}{k^\top \cdot \text{diag}(\mu) \cdot u^{(\ell)}}.$$

where the  $\cdot$  denotes matrix-vector multiplication and the fraction of two vectors is to be understood entry-wise. This is the famous Sinkhorn algorithm and its main loop can be written in two lines in most scientific computing environments.

- Note that since  $\phi, \psi \in \mathbb{R}^N$ , one has  $u, v = \exp(\phi/\varepsilon), \exp(\psi/\varepsilon) \in \mathbb{R}_{++}^N$  and also  $k = \exp(-c/\varepsilon) \in \mathbb{R}_{++}^{M \times N}$ , the division is always well-defined. However, numerically this may become an issue.

### Corresponding primal sequence.

- Recall the primal-dual optimality condition:

$$\pi_{i,j} = \exp\left(-\frac{c_{i,j} - \phi_i - \psi_j}{\varepsilon}\right) \mu_i \nu_j = u_i \cdot k_{i,j} \cdot v_j \cdot \mu_i \nu_j$$

So we can associate a sequence of primal iterates with the dual iterates. Note the following:

$$\begin{aligned} \sum_j u_i^{(\ell)} k_{i,j} v_j^{(\ell-1)} \mu_i \nu_j &= \sum_j \frac{1}{\sum_{j'} k_{i,j'} v_{j'}^{(\ell-1)} \nu_{j'}} k_{i,j} v_j^{(\ell-1)} \mu_i \nu_j = \mu_i \\ \sum_i u_i^{(\ell)} k_{i,j} v_j^{(\ell)} \mu_i \nu_j &= \sum_i \frac{1}{\sum_{i'} k_{i',j} u_{i'}^{(\ell)} \mu_{i'}} k_{i,j} u_i^{(\ell)} \mu_i \nu_j = \nu_j \end{aligned}$$

- So after a  $u$ -update, the primal iterate satisfies the row-constraints, after a  $v$ -update it satisfies the column-constraints.
- The updates can be interpreted as re-scaling each row or column such that those constraints are satisfied.
- Since the map from dual to primal iterates is continuous, convergence of dual iterates implies convergence of primal iterates.
- By stationarity of the optimal duals (under further iterations), the limit of the primal iterates satisfies row and column constraints and therefore, by the primal-dual optimality condition, must be the unique primal minimizer.

### Some comments on the Sinkhorn algorithm.

- The algorithm becomes numerically unstable and slow as  $\varepsilon \rightarrow 0$ . Good tricks for numerical stabilization exist.
- Some results on speed of convergence:
  - Franklin, Lorenz, Linear Algebra and its Applications, (1989): linear convergence of dual iterates to maximizer in Hilbert’s projective metric. But: contraction ratio approaches 1 like  $1 - \exp(-\|c\|_\infty/\varepsilon)$  as  $\varepsilon \rightarrow 0$ .
  - Schmitzer, SIAM J. Sci. Comput. (2019): convergence of an asymmetric (‘auction-like’) Sinkhorn algorithm in  $O(1/\varepsilon)$  iterations (measured in  $L^1$ -error of primal iterate marginal constraints)
  - Berman, Numerische Mathematik (2020): convergence of the Sinkhorn algorithm for the  $W_2$  distance on the Torus in  $O(1/\varepsilon)$  iterations, by showing that the iterates asymptotically follow a non-linear PDE
  - $\varepsilon$ -scaling very efficient in practice (at least on ‘normal problems’) but no proof for its efficiency yet (as far as I am aware).
  - There are several variants of Sinkhorn, intended to be faster, such as the ‘Greenhorn’ algorithm.



- One of big strengths of Sinkhorn is, that it can easily be adapted to related problems, such as optimal transport barycenters, multi-marginal transport problems (only efficient, if there is some trick to reduce the problem dimensionality), and unbalanced transport problems.