

Practical Optimal Transport for Machine Learning

Nicolas Courty¹, Laetitia Chapel¹ and Rémi Flamary²

¹ IRISA, Université Bretagne Sud, France

² CMAP – Ecole Polytechnique

Ecole de recherche (part. 2) Rennes – Juin 2022

Divergences statistiques et géométriques pour l'apprentissage machine

Planning of the course:

- (Morning) 1h30 of introductory course to *computational optimal transport* (Nicolas)
- (Afternoon) 1h30 of introduction to Unbalanced OT and Transport between metric spaces (Laetitia)
- (Afternoon) 1h30 of practical sessions in Python (POT)
- (Tomorrow Morning) 1h30 of applications of OT to Machine Learning (Rémi)

Table of content

Optimal transport : introduction

- Introduction to OT

- Simple applications

Wasserstein distances

- Definition

- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Regularized optimal transport

- Dual formulation

- Minimizing the Wasserstein distance

Gradient Flows in Wasserstein Space

Optimal transport : introduction

What is optimal transport ?

The natural geometry of probability measures



Monge

Kantorovich

Koopmans

Dantzig

Brenier

Otto

McCann

Villani

Nobel '75

Fields '10

666. MÉMOIRES DE L'ACADEMIE ROYALE

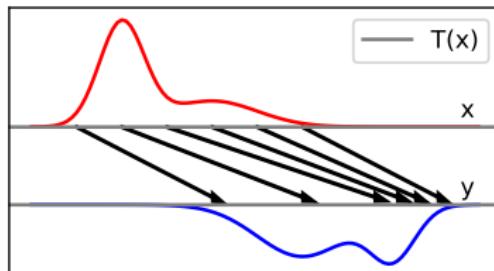
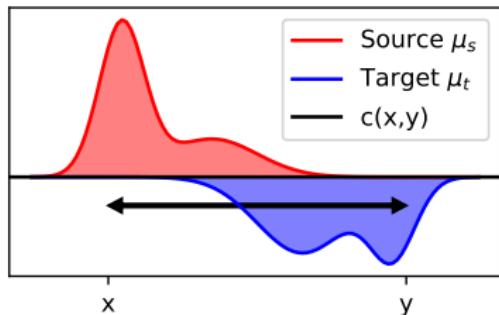
MÉMOIRE
SUR LA
THÉORIE DES DÉBLAIS
ET DES REMBLAIS.
Par M. MONGE.



Problem [Monge, 1781]

- How to move dirt from one place (déblais) to another (remblais) while minimizing the effort ?
- Find a mapping T between the two distributions of mass (transport).
- Optimize with respect to a displacement cost $c(x, y)$ (optimal).

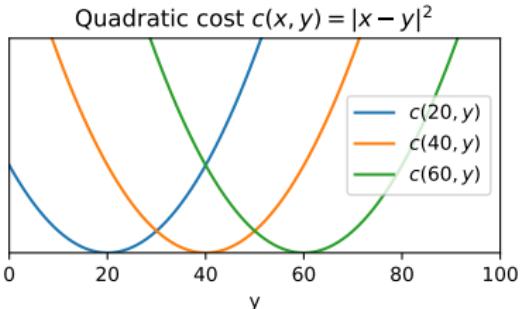
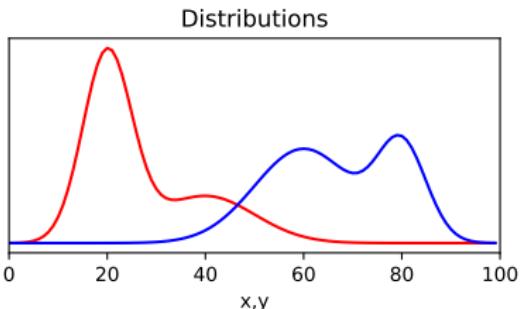
The origins of optimal transport



Problem [Monge, 1781]

- How to move dirt from one place (déblais) to another (remblais) while minimizing the effort ?
- Find a mapping T between the two distributions of mass (transport).
- Optimize with respect to a displacement cost $c(x, y)$ (optimal).

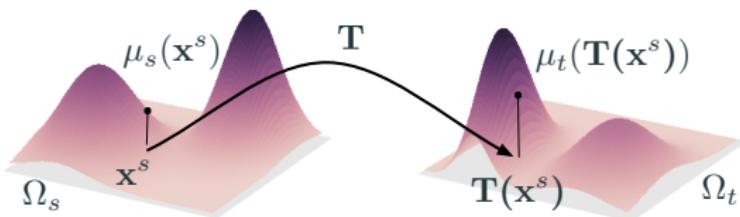
Optimal transport (Monge formulation)



- Probability measures μ_s and μ_t on and a cost function $c : \Omega_s \times \Omega_t \rightarrow \mathbb{R}^+$.
- The Monge formulation [Monge, 1781] aim at finding a mapping $T : \Omega_s \rightarrow \Omega_t$

$$\inf_{T \# \mu_s = \mu_t} \int_{\Omega_s} c(\mathbf{x}, T(\mathbf{x})) \mu_s(\mathbf{x}) d\mathbf{x} \quad (1)$$

What is $T\#\mu_s = \mu_t$?



- $T\#$ is the so called push forward operator
- it transfers measures from one space Ω_s to another space Ω_t
- it is equivalent to:

$$\mu_t(A) = \mu_s(T^{-1}(A))$$

$$\int_{\Omega_t} g(y) d\mu_t(y) = \int_{\Omega_s} g(T(x)) d\mu_s(x)$$

- for smooth measures $\mu_s = \rho(x)dx$ and $\mu_t = \eta(x)dx$

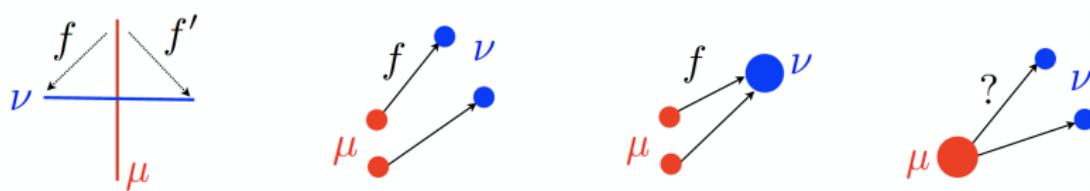
$$T\#\mu_s = \mu_t \equiv \rho(T(x)) |\det(\partial T(x))| = \eta(x)$$

- a.k.a. change of variable formula

Non-existence / Non-uniqueness

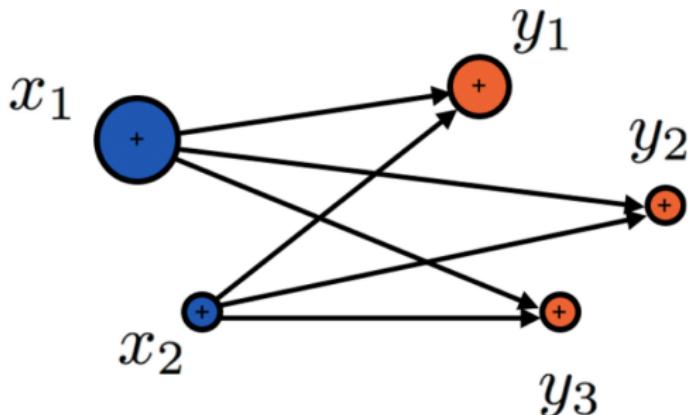
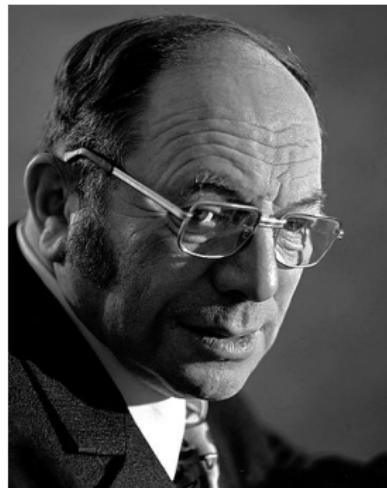
Solving for this push-forward operator is a non-convex optimization problem,

- for which existence is not guaranteed,
- nor unicity



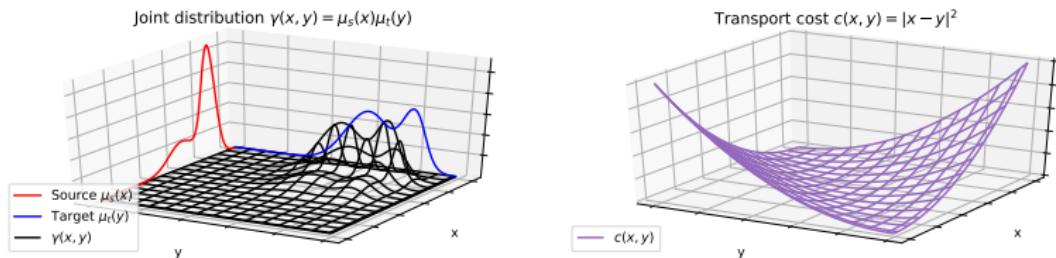
Note: [Brenier, 1991] proved existence and unicity of the Monge map for $c(x, y) = \|x - y\|^2$ and distributions with densities (i.e. continuous).

Kantorovich relaxation



- Leonid Kantorovich (1912–1986), Economy nobelist in 1975, proposed a different formulation of the problem
- with applications mainly for ressource allocation problems

Optimal transport (Kantorovich formulation)



- The Kantorovich formulation [Kantorovich, 1942] seeks for a probabilistic coupling $\gamma \in \mathcal{P}(\Omega_s \times \Omega_t)$ between Ω_s and Ω_t :

$$\gamma_0 = \operatorname{argmin}_{\gamma} \int_{\Omega_s \times \Omega_t} c(\mathbf{x}, \mathbf{y}) \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}, \quad (2)$$

$$\text{s.t. } \gamma \in \mathcal{P} = \left\{ \gamma \geq \mathbf{0}, \int_{\Omega_t} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mu_s, \int_{\Omega_s} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \mu_t \right\}$$

- γ is a joint probability measure with marginals μ_s and μ_t .
- Linear Program that always have a solution.

The 3 ways of optimal transport

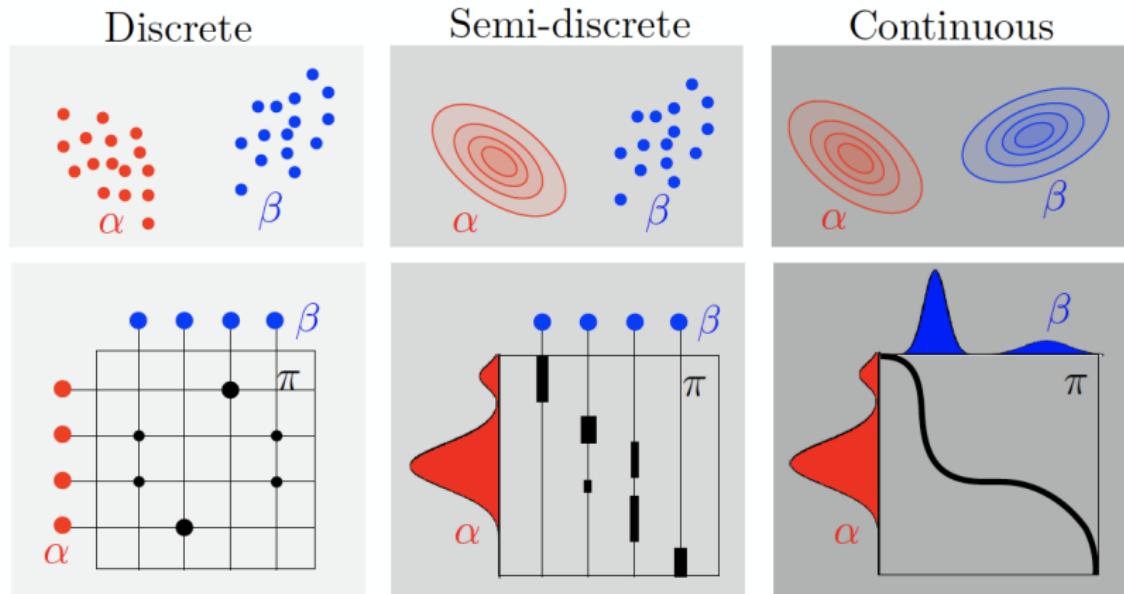
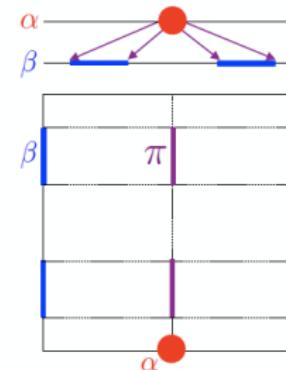
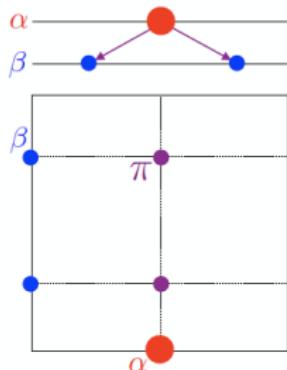
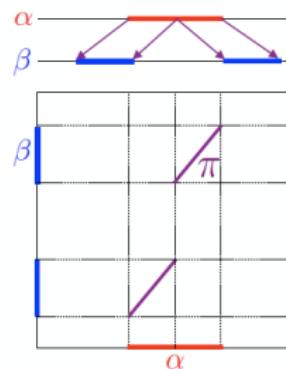
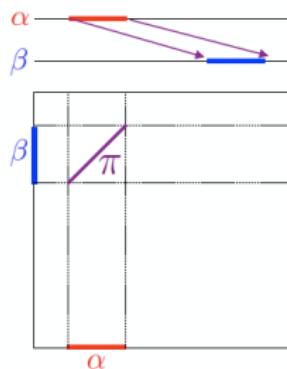


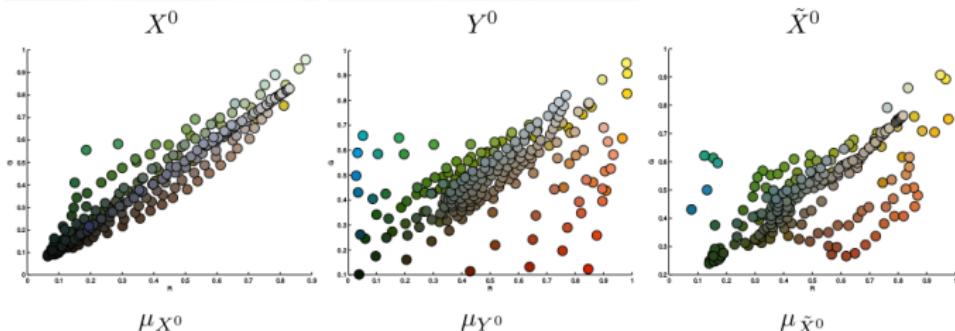
Image from Gabriel Peyré

Couplings



Histogram matching in images : color grading

Pixels as empirical distribution [Ferradans et al., 2014]



Histogram matching in images : color grading

Image colorization [Ferradans et al., 2014]



Matching words embedding



- Words are embedded in a high-dimensional space with neural networks
- Matching two documents is an OT problem, with the cost being the l_2 distance in the embedded space

Outline

Optimal transport : introduction

Introduction to OT

Simple applications

Wasserstein distances

Definition

Barycenters and geometry of optimal transport

Computational aspects of optimal transport

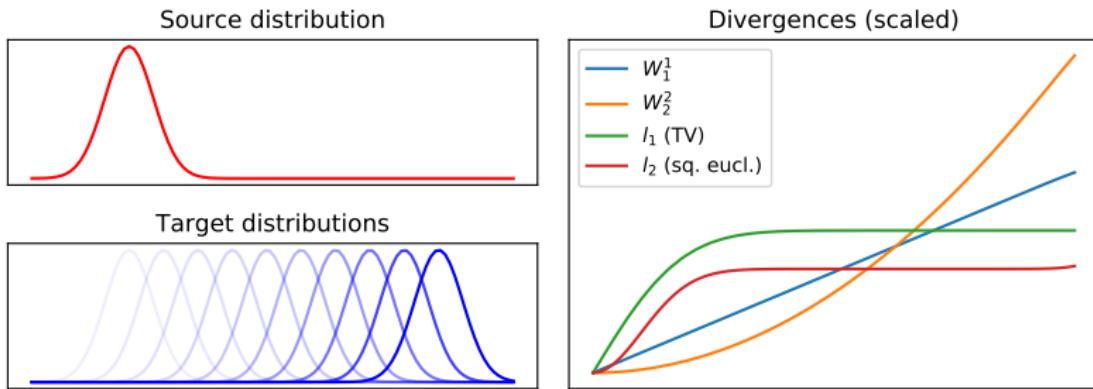
Regularized optimal transport

Dual formulation

Minimizing the Wasserstein distance

Gradient Flows in Wasserstein Space

Wasserstein distance



Wasserstein distance

$$W_p^p(\mu_s, \mu_t) = \min_{\gamma \in \mathcal{P}} \int_{\Omega_s \times \Omega_t} c(\mathbf{x}, \mathbf{y}) \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} [c(\mathbf{x}, \mathbf{y})] \quad (3)$$

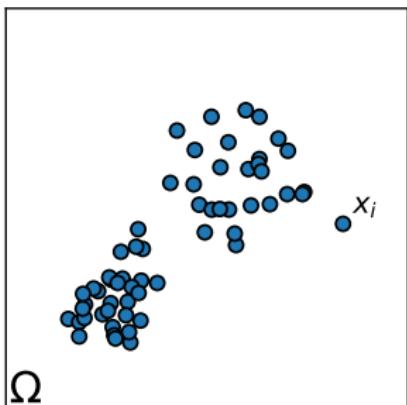
where $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^p$

- A.K.A. Earth Mover's Distance (W_1^1) [Rubner et al., 2000].
- Do not need the distribution to have overlapping support.
- Works for continuous and discrete distributions (histograms, empirical).

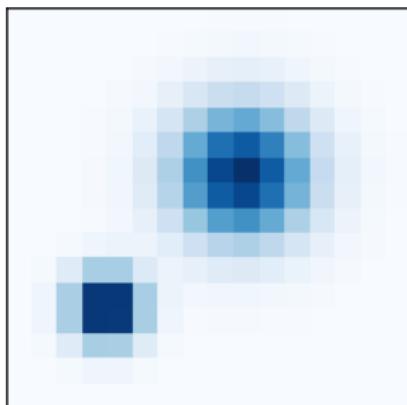
Discrete distributions: Empirical vs Histogram

Discrete measure: $\mu = \sum_{i=1}^n \mu_i \delta_{\mathbf{x}_i}, \quad \mathbf{x}_i \in \Omega, \quad \sum_{i=1}^n \mu_i = 1$

Lagrangian (point clouds)



Eulerian (histograms)



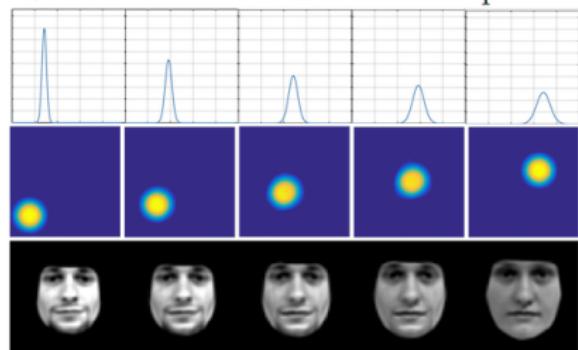
- Constant weight: $\mu_i = \frac{1}{n}$
- Quotient space: Ω^n, Σ_n
- Fixed positions \mathbf{x}_i e.g. grid
- Convex polytope Σ_n (simplex):
$$\{(\mu_i)_i \geq 0; \sum_i \mu_i = 1\}$$

Wasserstein space

The space of probability distribution equipped with the Wasserstein metric ($\mathcal{P}_p(X)$, $W_2^2(X)$) defines a geodesic space with a Riemannian structure [Santambrogio, 2014].

- Geodesics are shortest curves on $\mathcal{P}_p(X)$ that link two distributions

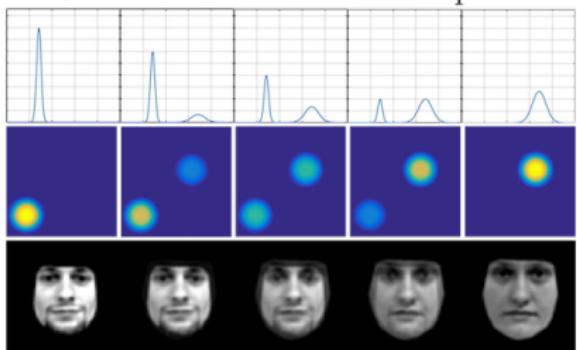
Geodesic in the 2-Wasserstein space



$$t = 0 \quad t = 0.25 \quad t = 0.5 \quad t = 0.75 \quad t = 1$$

$$\begin{aligned}\rho^*(., t) &= ((1-t)id + tf^*)\#\mu \\ d\rho^*(x, t) &= I^*(x, t)dx\end{aligned}$$

Geodesic in the Euclidean space

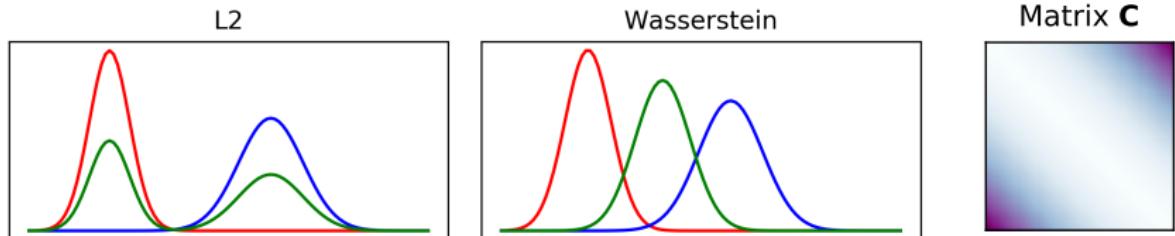


$$t = 0 \quad t = 0.25 \quad t = 0.5 \quad t = 0.75 \quad t = 1$$

$$I(x, t) = (1-t)I_0(x) + tI_1(x)$$

Illustration by S. Kolouri

Wasserstein barycenter

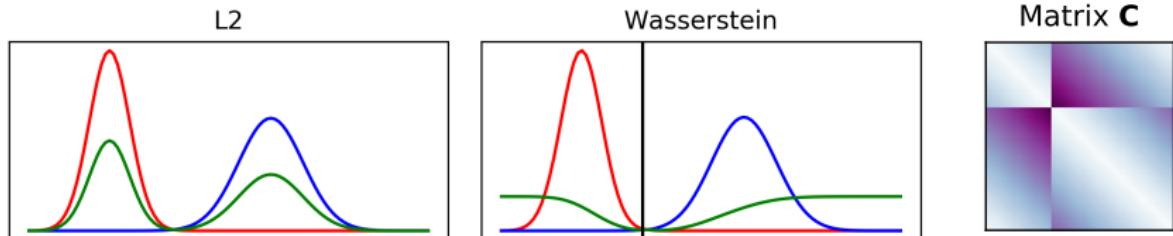


Barycenters [Aguech and Carlier, 2011]

$$\bar{\mu} = \arg \min_{\mu} \sum_i^n \lambda_i W_p^p(\mu^i, \mu)$$

- $\lambda_i > 0$ and $\sum_i^n \lambda_i = 1$.
- Uniform barycenter has $\lambda_i = \frac{1}{n}, \forall i$.
- Interpolation with $n=2$ and $\lambda = [1-t, t]$ with $0 \leq t \leq 1$ [McCann, 1997].
- Regularized barycenters using Bregman projections [Benamou et al., 2015].
- The cost and regularization impacts the interpolation trajectory.

Wasserstein barycenter



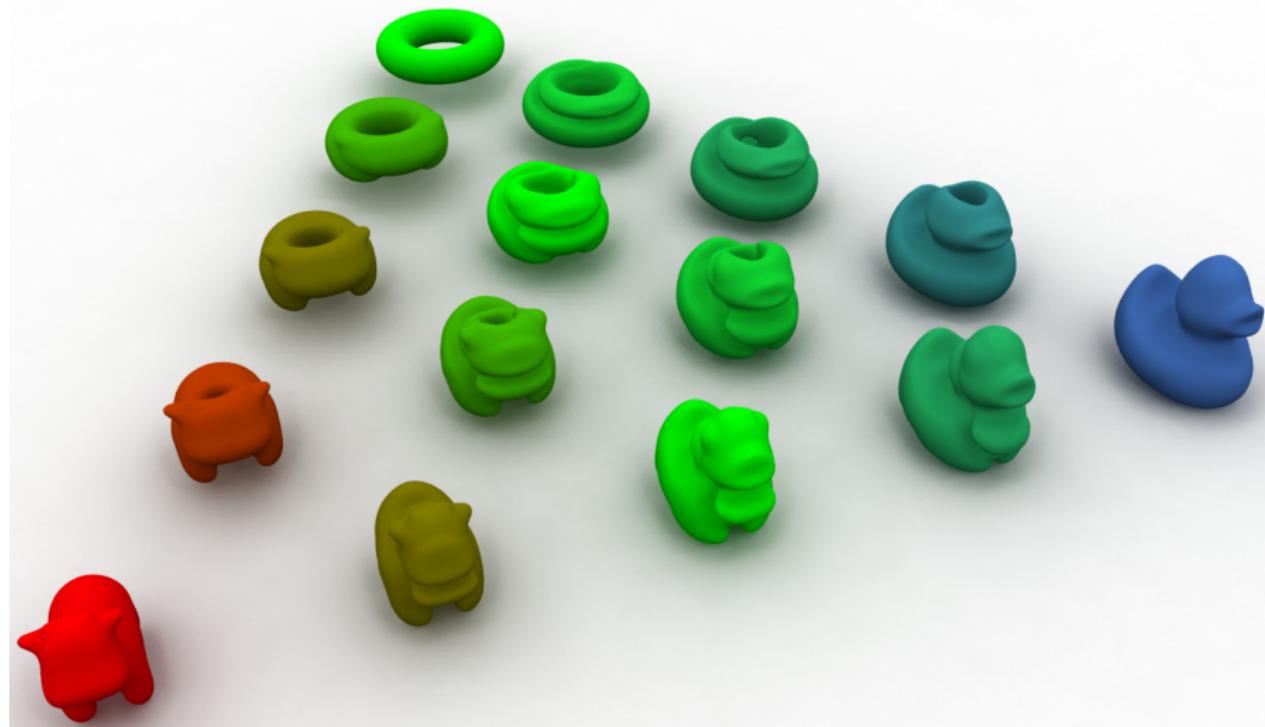
Barycenters [Aguech and Carlier, 2011]

$$\bar{\mu} = \arg \min_{\mu} \sum_i^n \lambda_i W_p^p(\mu^i, \mu)$$

- $\lambda_i > 0$ and $\sum_i^n \lambda_i = 1$.
- Uniform barycenter has $\lambda_i = \frac{1}{n}, \forall i$.
- Interpolation with $n=2$ and $\lambda = [1-t, t]$ with $0 \leq t \leq 1$ [McCann, 1997].
- Regularized barycenters using Bregman projections [Benamou et al., 2015].
- The cost and regularization impacts the interpolation trajectory.

3D Wasserstein barycenter

Shape interpolation [Solomon et al., 2015]



Principal Geodesics Analysis

Class 0			Class 1			Class 4											
PCA			PGA			PCA			PGA			PCA			PGA		
1	2	3	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3
0	0	0	0	0	0	1	X	X	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	X	X	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	X	X	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	1	1	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	1	1	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	1	1	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	1	1	1	1	1	4	4	4	4	4	4
0	0	0	0	0	0	1	1	1	1	1	1	4	4	4	4	4	4

Geodesic PCA in the Wasserstein space [Bigot et al., 2017]

- Generalization of Principal Component Analysis to the Wasserstein manifold.
- Regularized OT [Seguy and Cuturi, 2015].
- Approximation using Wasserstein embedding [Courty et al., 2017].
- Also note recent Wasserstein Dictionary Learning approaches [Schmitz et al., 2017].

Outline

Optimal transport : introduction

- Introduction to OT

- Simple applications

Wasserstein distances

- Definition

- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Regularized optimal transport

- Dual formulation

- Minimizing the Wasserstein distance

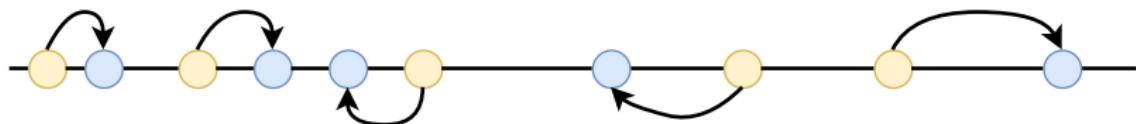
Gradient Flows in Wasserstein Space

Special case: 1D distribution

We consider the case where $c(x, y)$ is a strictly convex and increasing function of $|x - y|$.

- if $x_1 < x_2$ and $y_1 < y_2$, it is easy to check that
 $c(x_1, y_1) + c(x_2, y_2) < c(x_1, y_2) + c(x_2, y_1)$
- As such, any optimal transport plan respects the ordering of the elements, and the solution is given by the monotone rearrangement of μ_1 onto μ_2

This gives very simple algorithm to compute the transport in $O(N \log N)$, by sorting both x_i and y_i and summing the absolute values of differences.



Special case: 1D distribution

Consider the cumulative distribution functions F_μ associated to the μ distribution.

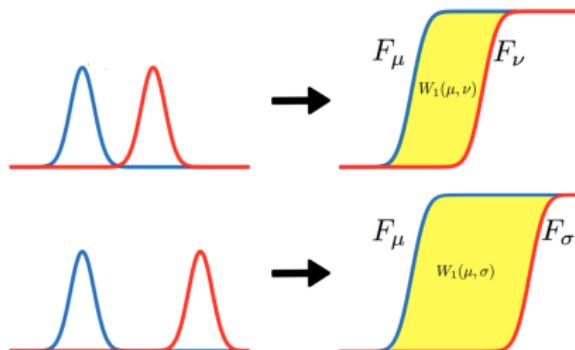
- It is defined such that $F_\mu(t) = \mu(-\infty, t]$.

We will note $F_\mu^{-1}(q)$, $q \in [0, 1]$ the corresponding generalized inverse distribution (or quantile function)

- defined as $F_\mu^{-1}(q) = \inf\{x \in \mathbb{R} : F_\mu(x) \geq q\}$.

Then,

$$W_1(\mu_s, \mu_t) = \int_0^1 c(F_{\mu_s}^{-1}(q), F_{\mu_t}^{-1}(q)) dq$$



Sliced-Wasserstein on \mathbb{R}^d

Wasserstein on \mathbb{R} :

$$\forall p \geq 1, \forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}), W_p^p(\mu, \nu) = \int_0^1 |F_\mu^{-1}(u) - F_\nu^{-1}(u)|^p \, du \quad (4)$$

This property gives a method for computing Wasserstein in higher dimensions ($n > 1$).

The principle is simple. Slice the distribution along lines, project the measures onto it and compute 1D Wasserstein along those projections.

Sliced-Wasserstein [Rabin et al., 2011b]

Let $\mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$,

$$SW_p^p(\mu, \nu) = \int_{S^{d-1}} W_p^p(P_\#^\theta \mu, P_\#^\theta \nu) \, d\lambda(\theta), \quad (5)$$

where $P^\theta(x) = \langle x, \theta \rangle$, λ uniform measure on S^{d-1} .

Properties:

- Distance
- Topologically equivalent to the Wasserstein distance
- Monte-Carlo approximation in $O(Ln \log n)$

Sliced Radon Wasserstein

This has been extended to a more general setting. Consider the Radon transform \mathcal{R} :

$$\mathcal{R}(\mu, \theta) = \int_{\mathbb{S}^{d-1}} \mu(\mathbf{x}) \delta(t - \theta \cdot \mathbf{x}) dx$$

where $t \in \mathbb{R}$ parametrizes the support and $\forall \theta \in \mathbb{S}^{d-1}$ (unit sphere in \mathbb{R}^d). Then, the p-sliced Wasserstein distance is given by:

p-sliced Wasserstein distance pSW [Bonneel et al., 2015]

$$pSW_p^p(\mu_s, \mu_t) = \int_{\mathbb{S}^{d-1}} W_p^p(\mathcal{R}(\mu_s, \theta), \mathcal{R}(\mu_t, \theta)) d\theta$$

Special cases: Wasserstein on the Circle

Let $\mu, \nu \in \mathcal{P}(S^1)$ where $S^1 = \mathbb{R}/\mathbb{Z}$.

- Parametrize S^1 by $[0, 1[$
- $\forall x, y \in [0, 1[, d_{S^1}(x, y) = \min(|x - y|, 1 - |x - y|)$
- For a cost function $c(x, y) = h(d_{S^1}(x, y))$ with $h : \mathbb{R} \rightarrow \mathbb{R}^+$ increasing and convex
- $\forall \mu, \nu \in \mathcal{P}(S^1)$, [Rabin et al., 2011a]

$$W_c(\mu, \nu) = \inf_{\alpha \in \mathbb{R}} \int_0^1 h(|F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t)|) dt. \quad (6)$$

- To find α : binary search [Delon et al., 2010]

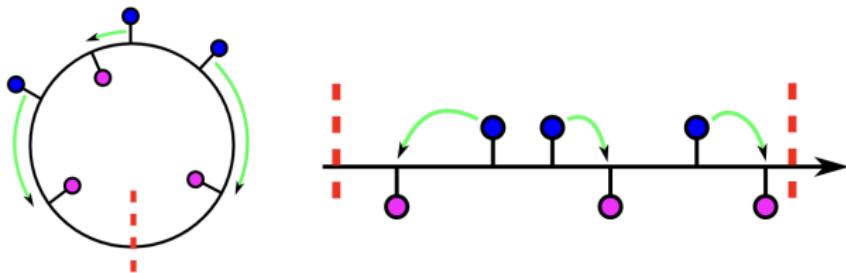


Image from [Rabin et al., 2011a]

Particular Cases

- For $h = \text{Id}$, [Hundrieser et al., 2021]

$$W_1(\mu, \nu) = \int_0^1 |F_\mu(t) - F_\nu(t) - \text{LevMed}(F_\mu - F_\nu)| dt, \quad (7)$$

where

$$\text{LevMed}(f) = \inf \left\{ t \in \mathbb{R}, \beta(\{x \in [0, 1], f(x) \leq t\}) \geq \frac{1}{2} \right\}. \quad (8)$$

- For $h(x) = x^2$ and $\nu = \text{Unif}(S^1)$, [Bonet et al., 2022]

$$W_2^2(\mu, \nu) = \int_0^1 |F_\mu^{-1}(t) - t - \hat{\alpha}|^2 dt \quad \text{with} \quad \hat{\alpha} = \int x d\mu(x) - \frac{1}{2}. \quad (9)$$

In particular, if $x_1 < \dots < x_n$ and $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$, then

$$W_2^2(\mu_n, \nu) = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n^2} \sum_{i=1}^n (n+1-2i)x_i + \frac{1}{12}. \quad (10)$$

Sliced-Wasserstein on the Sphere

The slicing strategy can be extended to manifolds ! In this case slices are geodesics of the considered manifold. Example on the Sphere [Bonet et al., 2022]:

- Great circle: Intersection between 2-plane and S^{d-1}
- Parametrize 2-plane by the Stiefel manifold

$$\mathbb{V}_{d,2} = \{U \in \mathbb{R}^{d \times 2}, U^T U = I_2\}$$

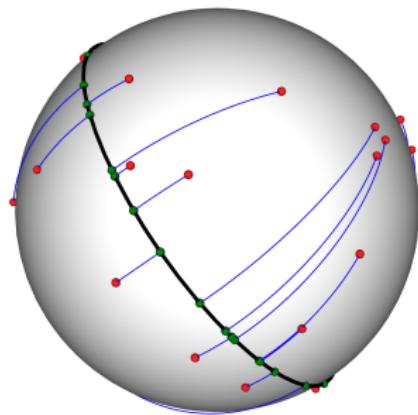
- Projection on great circle C : For a.e. $x \in S^{d-1}$,

$$P^C(x) = \operatorname{argmin}_{y \in C} d_{S^{d-1}}(x, y),$$

where $d_{S^{d-1}}(x, y) = \arccos(\langle x, y \rangle)$.

- For $U \in \mathbb{V}_{d,2}$, $C = \operatorname{span}(UU^T) \cap S^{d-1}$,

$$\begin{aligned} P^U(x) &= U^T \operatorname{argmin}_{y \in C} d_{S^{d-1}}(x, y) \\ &= \frac{U^T x}{\|U^T x\|_2}. \end{aligned}$$



Special case: transport between Gaussians

In the case where $\mu_s \sim \mathcal{N}(\mathbf{m}_1, \Sigma_1)$ and $\mu_t \sim \mathcal{N}(\mathbf{m}_2, \Sigma_2)$ the Wasserstein distance with $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ reduces to:

W_2^2 between Gaussians

$$W_2^2(\mu_s, \mu_t) = \|\mathbf{m}_1 - \mathbf{m}_2\|_2^2 + \mathcal{B}(\Sigma_1, \Sigma_2)^2$$

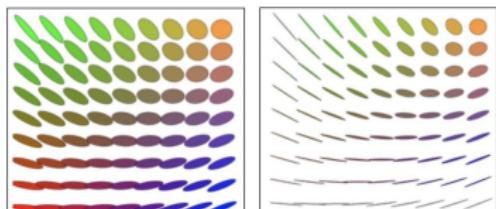
where $\mathbb{B}(,)$ is the so-called Bures metric:

$$\mathcal{B}(\Sigma_1, \Sigma_2)^2 = \text{trace}(\Sigma_1 + \Sigma_2 - 2(\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2}).$$

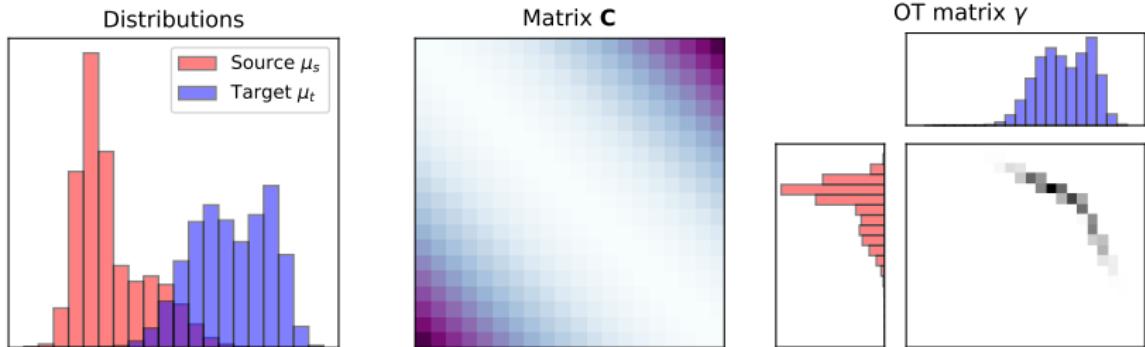
The optimal map T is given by

$$T(\mathbf{x}) = \mathbf{m}_2 + A(\mathbf{x} - \mathbf{m}_1)$$

$$\text{with } A = \Sigma_1^{-1/2} (\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \Sigma_1^{-1/2}$$



Optimal transport with discrete distributions



OT Linear Program

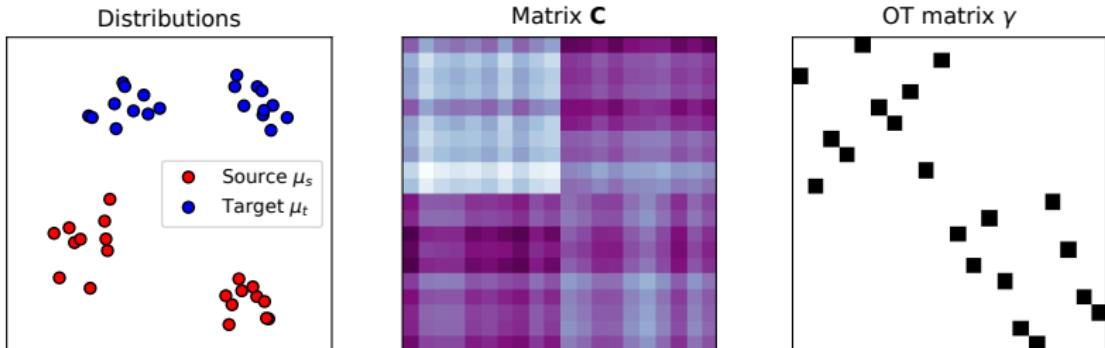
$$\gamma_0 = \operatorname{argmin}_{\gamma \in \mathcal{P}} \left\{ \langle \gamma, \mathbf{C} \rangle_F = \sum_{i,j} \gamma_{i,j} c_{i,j} \right\}$$

where \mathbf{C} is a cost matrix with $c_{i,j} = c(\mathbf{x}_i^s, \mathbf{x}_j^t)$ and the marginal constraints are

$$\mathcal{P} = \left\{ \gamma \in (\mathbb{R}^+)^{n_s \times n_t} \mid \gamma \mathbf{1}_{n_t} = \mu_s, \gamma^T \mathbf{1}_{n_s} = \mu_t \right\}$$

Solved with Network Flow solver of complexity $O(n^3 \log(n))$.

Optimal transport with discrete distributions



OT Linear Program

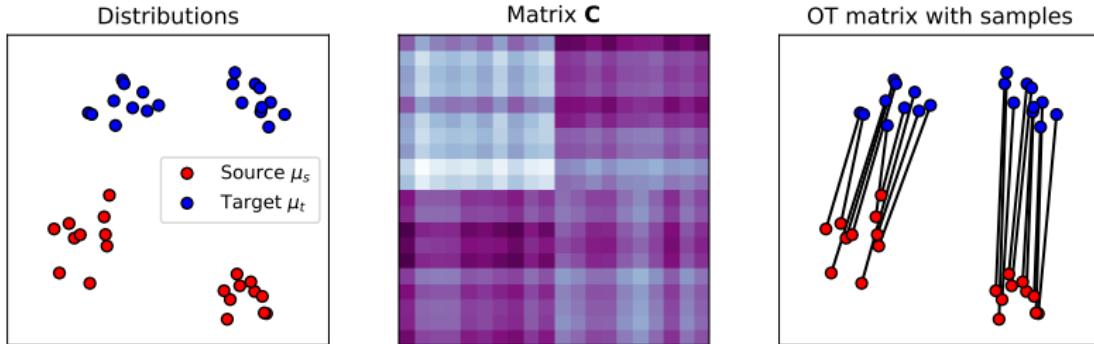
$$\gamma_0 = \operatorname{argmin}_{\gamma \in \mathcal{P}} \left\{ \langle \gamma, \mathbf{C} \rangle_F = \sum_{i,j} \gamma_{i,j} c_{i,j} \right\}$$

where \mathbf{C} is a cost matrix with $c_{i,j} = c(\mathbf{x}_i^s, \mathbf{x}_j^t)$ and the marginal constraints are

$$\mathcal{P} = \left\{ \gamma \in (\mathbb{R}^+)^{n_s \times n_t} \mid \gamma \mathbf{1}_{n_t} = \mu_s, \gamma^T \mathbf{1}_{n_s} = \mu_t \right\}$$

Solved with Network Flow solver of complexity $O(n^3 \log(n))$.

Optimal transport with discrete distributions



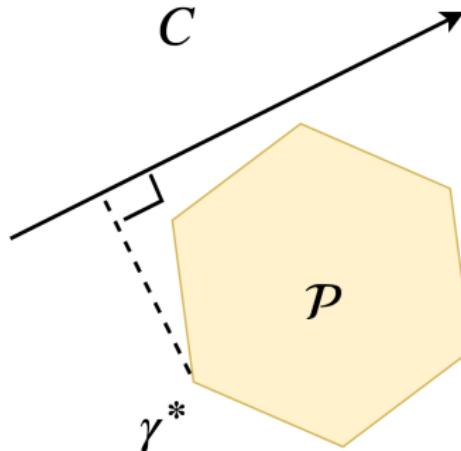
OT Linear Program

$$\gamma_0 = \operatorname{argmin}_{\gamma \in \mathcal{P}} \left\{ \langle \gamma, \mathbf{C} \rangle_F = \sum_{i,j} \gamma_{i,j} c_{i,j} \right\}$$

where \mathbf{C} is a cost matrix with $c_{i,j} = c(\mathbf{x}_i^s, \mathbf{x}_j^t)$ and the marginal constraints are

$$\mathcal{P} = \left\{ \gamma \in (\mathbb{R}^+)^{n_s \times n_t} \mid \gamma \mathbf{1}_{n_t} = \mu_s, \gamma^T \mathbf{1}_{n_s} = \mu_t \right\}$$

Solved with Network Flow solver of complexity $O(n^3 \log(n))$.



- \mathcal{P} is the Birkhoff polytope
- No unique solution in some cases, numerical instabilities
- Not differentiable !

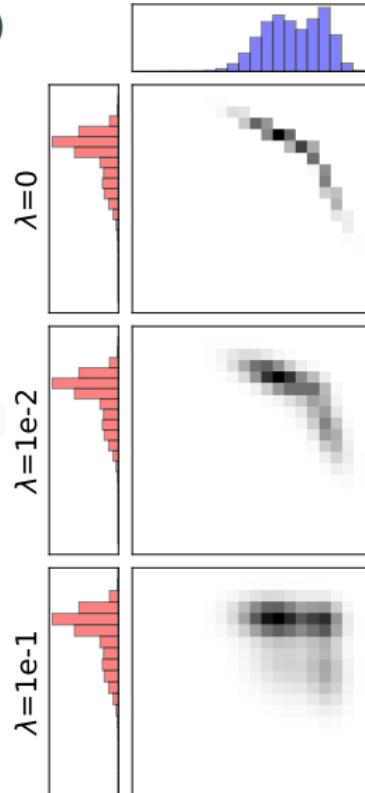
$$\gamma_0^\lambda = \operatorname{argmin}_{\gamma \in \mathcal{P}} \langle \gamma, \mathbf{C} \rangle_F + \lambda \Omega(\gamma), \quad (11)$$

Regularization term $\Omega(\gamma)$

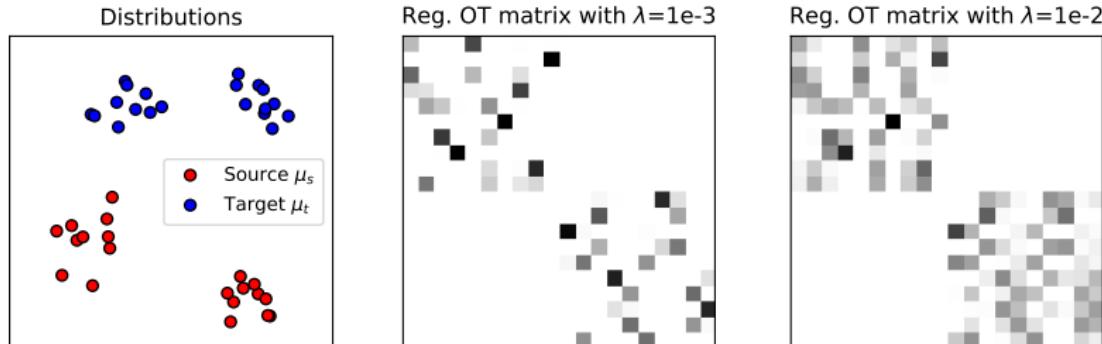
- Entropic regularization [Cuturi, 2013].
- Group Lasso [Courty et al., 2016].
- KL, Itakura Saito, β -divergences, [Dessein et al., 2016].

Why regularize?

- Smooth the “distance” estimation:
$$W_\lambda(\mu_s, \mu_t) = \langle \gamma_0^\lambda, \mathbf{C} \rangle_F$$
- Encode prior knowledge on the data.
- Better posed problem (convex, stability).
- Fast algorithms to solve the OT problem.



Entropic regularized optimal transport

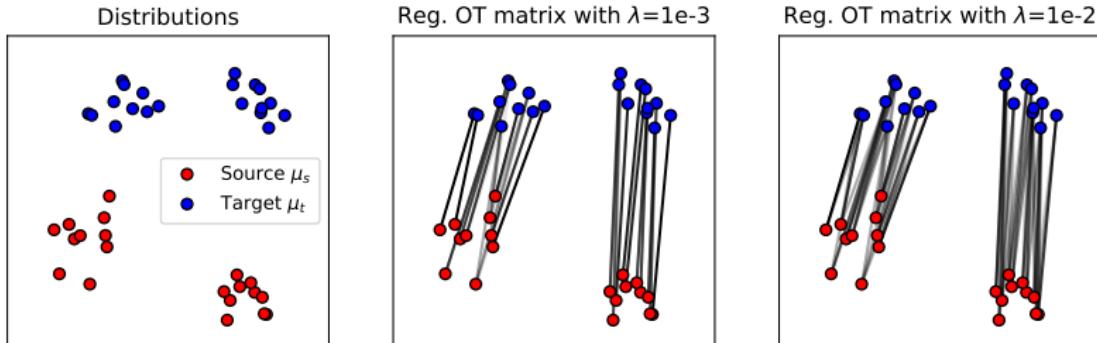


Entropic regularization [Cuturi, 2013]

$$\Omega(\gamma) = \sum_{i,j} \gamma(i,j)(\log \gamma(i,j) - 1)$$

- Regularization with the negative entropy of γ .

Entropic regularized optimal transport



Entropic regularization [Cuturi, 2013]

$$\Omega(\gamma) = \sum_{i,j} \gamma(i,j)(\log \gamma(i,j) - 1)$$

- Regularization with the negative entropy of γ .

Resolving the entropy regularized problem

Entropy-regularized transport

The solution of entropy regularized optimal transport problem is of the form

$$\gamma_0^\lambda = \text{diag}(\mathbf{u}) \exp(-\mathbf{C}/\lambda) \text{diag}(\mathbf{v})$$

Why ? Consider the Lagrangian of the optimization problem:

$$\mathcal{L}(\boldsymbol{\gamma}, \alpha, \beta) = \sum_{ij} \boldsymbol{\gamma}_{ij} \mathbf{C}_{ij} + \lambda \boldsymbol{\gamma}_{ij} (\log \boldsymbol{\gamma}_{ij} - 1) + \alpha^T (\boldsymbol{\gamma} \mathbf{1}_{n_t} - \mu_s) + \beta^T (\boldsymbol{\gamma}^T \mathbf{1}_{n_s} - \mu_t)$$

$$\frac{\partial \mathcal{L}(\boldsymbol{\gamma}, \alpha, \beta)}{\partial \boldsymbol{\gamma}_{ij}} = \mathbf{C}_{ij} + \lambda \log \boldsymbol{\gamma}_{ij} + \alpha_i + \beta_j$$

$$\frac{\partial \mathcal{L}(\boldsymbol{\gamma}, \alpha, \beta)}{\partial \boldsymbol{\gamma}_{ij}} = 0 \implies \boldsymbol{\gamma}_{ij} = \exp\left(\frac{\alpha_i}{\lambda}\right) \exp\left(-\frac{\mathbf{C}_{ij}}{\lambda}\right) \exp\left(\frac{\beta_j}{\lambda}\right)$$

- Through the **Sinkhorn theorem** $\text{diag}(\mathbf{u})$ and $\text{diag}(\mathbf{v})$ exist and are unique.
- Can be solved by the **Sinkhorn-Knopp algorithm** (implementation in parallel, GPU).

Sinkhorn-Knopp algorithm

The Sinkhorn-Knopp algorithm performs alternatively a scaling along the rows and columns of $\mathbf{K} = \exp(-\frac{\mathbf{C}}{\lambda})$ to match the desired marginals.

Algorithm 1 Sinkhorn-Knopp Algorithm (SK).

Require: $\mathbf{a}, \mathbf{b}, \mathbf{C}, \lambda$

$$\mathbf{u}^{(0)} = \mathbf{1}, \mathbf{K} = \exp(-\mathbf{C}/\lambda)$$

for i in $1, \dots, n_{it}$ **do**

$$\mathbf{v}^{(i)} = \mathbf{b} \oslash \mathbf{K}^\top \mathbf{u}^{(i-1)} \text{ // Update right scaling}$$

$$\mathbf{u}^{(i)} = \mathbf{a} \oslash \mathbf{K} \mathbf{v}^{(i)} \text{ // Update left scaling}$$

end for

$$\mathbf{return} \quad \mathbf{T} = \text{diag}(\mathbf{u}^{(n_{it})}) \mathbf{K} \text{diag}(\mathbf{v}^{(n_{it})})$$

- Complexity $O(kn^2)$, where k iterations are required to reach convergence
- Fast implementation in parallel, GPU friendly
- Convolutional/Heat structure for \mathbf{K} [Solomon et al., 2015]

Recalling that the Kullback Leibler (KL) divergence between two distribution is

$$\text{KL}(\gamma, \rho) = \sum_{ij} \gamma_{ij} \log \frac{\gamma_{ij}}{\rho_{ij}} = \langle \gamma, \log \frac{\gamma}{\rho} \rangle_F,$$

Benamou et al. [Benamou et al., 2015] showed that solving for the OT problem is actually a Bregman projection

OT as a Bregman projection

γ^* is the solution of the following Bregman projection

$$\gamma^* = \underset{\gamma \in \mathcal{P}}{\operatorname{argmin}} \text{KL}(\gamma, \zeta), \quad (12)$$

where $\zeta = \exp(-\frac{C}{\lambda})$.

- Sinkhorn in this case is an iterative projection scheme, with alternative projections on marginal constraints.
- Generalizes well for barycenters computation

Dual formulation of optimal transport

- Yet, solving for γ is impractical to intractable when dealing with high-dimensional distributions
- especially if one is interested in computing the gradients of the Wasserstein distance
- Other solving strategies should be taken into consideration
- Recalling that any LP problem can be turned into its dual form:

<p>primal form :</p> <p>minimize $z = \mathbf{c}^T \mathbf{x}$, so that $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$</p>	<p>dual form :</p> <p>maximize $\tilde{z} = \mathbf{b}^T \mathbf{y}$, so that $\mathbf{A}^T \mathbf{y} \leq \mathbf{c}$</p>
---	--

- **Weak duality**: \tilde{z} is a lower bound of z , **Strong duality** $\tilde{z} = z$
- **Strong duality** is usually achieved via Farkas Theorem

Duality: general case with continuous distributions

We now introduce two functions scalar functions ϕ and ψ (also known as Kantorovich potentials) that will act as our dual variables. Then, we consider the optimal problem is equivalent (by the Rockafellar-Fenchel theorem) to:

$$\max_{\phi, \psi} \left\{ \int \phi d\mu_s + \int \psi d\mu_t \mid \phi(x) + \psi(y) \leq c(x, y) \right\} \quad (13)$$

Note that the marginal constraint has been turned into an equality constraint on ϕ and ψ

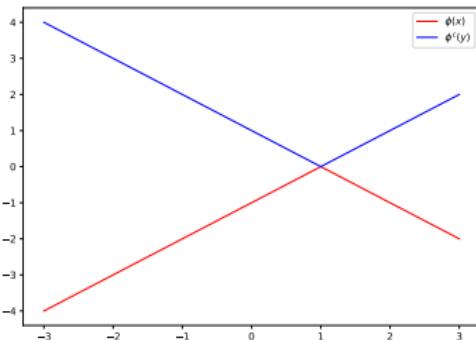
Introducing the *c-transform* (or *c-conjugate*) H^c which is in spirit close to a Legendre transform:

$$\phi^c \stackrel{\text{def}}{=} H^c(\phi) = \inf_x c(x, y) - \phi(x) \quad (14)$$

then the following problem is equivalent:

$$\max_{\phi} \left\{ \int \phi d\mu_s + \int \phi^c d\mu_t \mid \phi(x) + \phi^c(y) \leq c(x, y) \right\} \quad (15)$$

Case $c(x, y) = |x - y|$ (a.k.a W_1^1)



Whenever $c(x, y) = |x - y|$, then:

- existence of a solution but not unique
- For any $\phi \in \text{Lip}^1$ (set of 1-Lipschitz functions), we have $\phi^c(x) = -\phi(x)$

The optimal transport problem then amounts to find $\phi \in \text{Lip}^1$ as

$$\sup_{\phi \in \text{Lip}^1} \int \phi d(\mu_s - \mu_t) = \sup_{\phi \in \text{Lip}^1} \mathbb{E}_{x \sim \mu_s} [\phi(x)] - \mathbb{E}_{y \sim \mu_t} [\phi(y)] \quad (16)$$

- also known as **Kantorovich-Rubinstein duality**
- ϕ can be learnt as a neural network constrained to the set Lip^1 , see next section on GAN

Case $c(x, y) = |x - y|^2/2$ (a.k.a W_2^2)

Whenever the cost is quadratic, $c(x, y) = |x - y|^2/2$, then:

- $T(x)$ the transport mapping exists and is unique
- More remarkably, it is a gradient of a convex functions $\Phi(x)$

$$T(x) = x - \nabla \phi(x) = \nabla\left(\frac{x^2}{2} - \phi(x)\right) = \nabla(\Phi(x)) \quad (17)$$

Brenier's Theorem

Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, μ absolutely continuous with respect to the Lebesgue measure. Then, the optimal coupling γ^* is unique and of the form $\gamma^* = (Id, \nabla \Phi)_\# \mu$ with Φ a convex function.

- can be optimized for instance with Input Convex Neural Networks (ICNN)
[\[Amos et al., 2017\]](#) to model the convex functions

In the case when we have access to discrete distributions, μ_s (resp. μ_t) is characterized by a set of locations \mathbf{X}^s and masses $\mathbf{a} \in \mathbb{R}^{n^s}$ (resp. \mathbf{X}^t and $\mathbf{b} \in \mathbb{R}^{n^t}$)

Discrete dual version of OT

$$W(\mu_s, \mu_t) = \max_{\alpha \in \mathbb{R}^{n^s}, \beta \in \mathbb{R}^{n^t}, \alpha_i + \beta_j \leq c(\mathbf{x}_i^s, \mathbf{x}_j^t)} \alpha^T \mathbf{a} + \beta^T \mathbf{b} \quad (18)$$

i.e. find a scalar values per sample

Regularized case

Adding regularization to the original problem turns the dual computation to an **unconstrained problem** !

In the case of entropy regularization, i.e.

$$W_\lambda(\mu_s, \mu_t) = \min_{\gamma \in \mathcal{P}} \quad \langle \gamma, \mathbf{C} \rangle_F + \lambda \Omega(\gamma) \text{ with } \Omega(\gamma) = \sum_{i,j} \gamma(i,j) \log \gamma(i,j),$$

the dual now reads (in a discrete settings, measures are collections of Diracs):

$$\max_{\alpha, \beta} \alpha^T \mu_s + \beta^T \mu_t - \frac{1}{\lambda} \exp\left(\frac{\alpha}{\lambda}\right)^T \mathbf{K} \exp\left(\frac{\beta}{\lambda}\right) \quad (19)$$

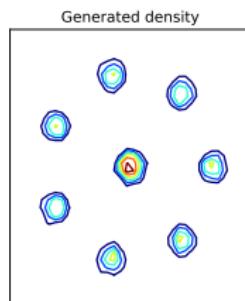
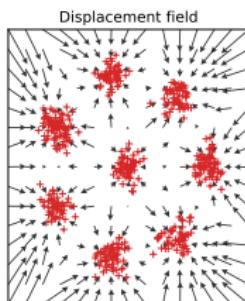
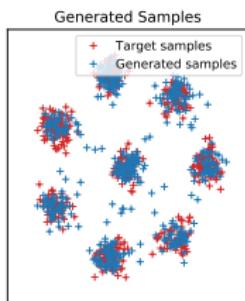
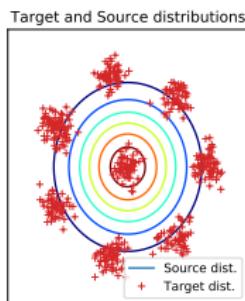
with $\mathbf{K} = \exp\left(-\frac{\mathbf{C}}{\lambda}\right)$.

Remark: The Sinkhorn algorithm is a gradient ascent on the dual variables !

Regularized case

With this unconstrained problem, incremental gradients techniques (SGD, SAG) can be used to solve the problem !

- [Genevay et al., 2016] used the semi-dual formulation (one variable is removed by replacing it with its c-transform) int the first stochastic version of Optimal Transport problem
- [Seguy et al., 2017] used the full dual version with entropic and L2 regularizations, together with neural networks to parameterize the problem.



2 ways of minimizing the Wasserstein distance

In machine learning applications, one can be interested in finding distributions that minimize the Wasserstein distance wrt. a reference measure. There are two ways of understanding this:

- case 1: **for a fixed support \mathbf{X}** , find the corresponding probability masses \mathbf{m}
- case 2: **for a fixed vector of probability masses \mathbf{m}** , e.g. uniform distribution, find the corresponding support \mathbf{X}

Case 1: fixed support

Recalling the form of the dual

$$W(\mu, \mu_t) = \max_{\alpha \in \mathbb{R}^{n^s}, \beta \in \mathbb{R}^{n^t}, \alpha_i + \beta_j \leq c(\mathbf{x}, \mathbf{x}_j^t)} \alpha^T \mathbf{m} + \beta^T \mathbf{b} \quad (20)$$

- $W(\mu, \mu_t)$ is convex wrt. \mathbf{m}
- $\partial_{\mathbf{m}} W(\mu, \mu_t) = \alpha^*$
- **Entropy regularized case:** $W_\lambda(\mu, \mu_t)$ is convex and $\nabla_{\mathbf{m}} W_\lambda(\mu, \mu_t) = \lambda \log \mathbf{u}$

Case 2: fixed probability masses \mathbf{m}

Recalling the form of the primal problem

$$W_2^2(\boldsymbol{\mu}, \boldsymbol{\mu_t}) = \min_{\gamma \in \mathcal{P}} \quad \langle \gamma, \mathbf{1}_{n^s} \mathbf{1}_{n^t}^T \mathbf{X}^2 + \mathbf{X}^{t2T} \mathbf{1}_{n^t} \mathbf{1}_{n^s} - 2\mathbf{X}\mathbf{X}^t \rangle \quad (21)$$

- $W_2^2(\boldsymbol{\mu}, \boldsymbol{\mu_t})$ decreases if $\mathbf{X} \leftarrow \mathbf{X}^t \gamma^{*T} \text{diag}(\mathbf{m}^{-1})$
- explicit gradient for the regularized case.
- Barycentric interpolation !
- see Rémi next slides

General case: autodifferentiation

Automatic differentiation to the rescue !

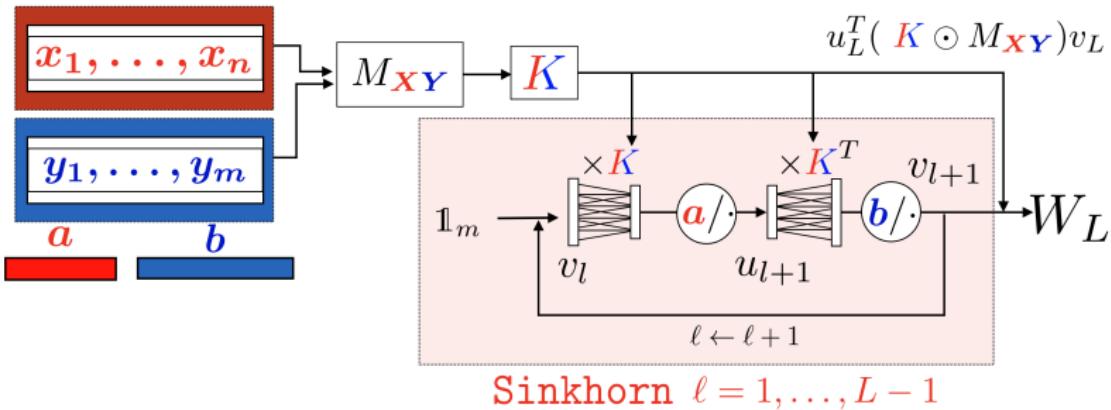


Image from Marco Cuturi

But also consider using the Enveloppe theorem, i.e. take the gradient in the optimal solution of the primal problem.

Contents

Optimal transport : introduction

- Introduction to OT

- Simple applications

Wasserstein distances

- Definition

- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Regularized optimal transport

- Dual formulation

- Minimizing the Wasserstein distance

Gradient Flows in Wasserstein Space

Outline

Optimal transport : introduction

- Introduction to OT

- Simple applications

Wasserstein distances

- Definition

- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Regularized optimal transport

- Dual formulation

- Minimizing the Wasserstein distance

Gradient Flows in Wasserstein Space

Point of View:

- Define Wasserstein gradient flows
- Analogies with gradient flows in Euclidean space
- For more abstract views, see [Santambrogio, 2017, Ambrosio et al., 2008]
- Talk from Anna Korba <https://mathtube.org/lecture/video/wasserstein-gradient-flows-machine-learning>

Let $X = \mathbb{R}^p$, d a distance (e.g. $d(x, y) = \|x - y\|_2$), $F : X \rightarrow \mathbb{R}$.

Goal:

$$\min_x F(x)$$

Let $X = \mathbb{R}^p$, d a distance (e.g. $d(x, y) = \|x - y\|_2$), $F : X \rightarrow \mathbb{R}$.

Goal:

$$\min_x F(x)$$

Definition (Gradient Flow on \mathbb{R}^p)

A gradient flow is a curve $x : [0, T] \rightarrow X$ which decreases as much as possible along the functional F .

i.e. If F is differentiable, x follows the Cauchy problem

$$\begin{cases} \frac{dx}{dt}(t) = -\nabla F(x(t)) \\ x(0) = x_0 \end{cases}$$

If F is differentiable, x follows the Cauchy problem

$$\begin{cases} \frac{dx}{dt}(t) = -\nabla F(x(t)) \\ x(0) = x_0 \end{cases}$$

Solving the ODE in practice:

- Explicit Euler scheme ($x_k = x(k\tau)$):

$$x_{k+1} = x_k - \tau \nabla F(x_k)$$

- Implicit Euler scheme:

$$\begin{aligned} x_{k+1} = x_k - \tau \nabla F(x_{k+1}) &\iff 0 = \frac{x_{k+1} - x_k}{\tau} + \nabla F(x_{k+1}) \\ &\iff x_{k+1} \in \operatorname{argmin}_{x \in X} \frac{\|x - x_k\|_2^2}{2\tau} + F(x) \\ &\iff x_{k+1} = \operatorname{prox}_{\tau F}(x_k) \end{aligned}$$

- Any ODE solver (Runge-Kutta...)

Other characterization

See [Santambrogio, 2017, Ambrosio et al., 2008]

- Energy Dissipation Equality (EDE):

$$\begin{aligned}\frac{dx}{dt}(t) = -\nabla F(x(t)) &\iff \forall 0 \leq s < t \leq 1, \\ F(x(s)) - F(x(t)) &= \int_s^t \left(\frac{1}{2}|x'(u)|^2 + \frac{1}{2}|\nabla F(x(u))|^2 \right) du,\end{aligned}\tag{22}$$

where (speed)

$$|x'(t)| = \lim_{h \rightarrow 0} \frac{d(x(t+h), x(t))}{h}$$

and (descending slope)

$$|\nabla F|(x) = \limsup_{y \rightarrow x} \frac{|F(y) - F(x)|}{d(x, y)}.$$

- Evolution Variational Inequality (EVI): For F λ -geodesically convex,

$$x'(t) \in \partial F(x(t)) \iff \forall y \in X, \frac{d}{dt} \frac{1}{2}|x(t) - y|^2 \leq F(y) - F(x(t)) - \frac{\lambda}{2}|x(t) - y|^2\tag{23}$$

Let $\mathcal{P}_2(\mathbb{R}^p) = \{\mu \in \mathcal{P}(\mathbb{R}^p), \int \|x\|^2 d\mu(x) < +\infty\}$.

Define Gradient Flows in $(\mathcal{P}_2(\mathbb{R}^p), W_2)$ via the JKO Scheme [Jordan et al., 1998]: Let $F : \mathcal{P}_2(\mathbb{R}^p) \rightarrow \mathbb{R}$, $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$,

$$\forall k \geq 0, \mu_{k+1}^\tau \in \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^p)} \frac{W_2^2(\mu, \mu_k^\tau)}{2\tau} + F(\mu) = \text{JKO}_{\tau F}(\mu_k^\tau) \quad (24)$$

Define a piecewise constant interpolation μ^τ , i.e.

$$\forall t \in [k\tau, (k+1)\tau[, \mu_t^\tau = \mu_k^\tau$$

Wasserstein gradient flows: $t \mapsto \mu_t$ such that, for $\mu^\tau \xrightarrow[\tau \rightarrow 0]{} \mu$.

(also called (generalized) minimizing movement [Bonnotte, 2013, Liutkus et al., 2019])

Gradient Flow in $(\mathcal{P}_2(\mathbb{R}^d), W_2)$:

Iterated Minimization scheme (JKO Scheme) [Jordan et al., 1998]:

$$\mu_{k+1}^\tau \in \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \frac{1}{2\tau} W_2^2(\mu, \mu_k^\tau) + F(\mu)$$

Examples

- $F(\mu) = \int \rho(x) \log \rho(x) dx + \int V(x)\rho(x)dx$ if $d\mu = \rho d\text{Leb}$
Solution in the limit $\tau \rightarrow 0$ to the PDE: (Fokker-Planck)

$$\partial_t \rho_t = \operatorname{div}(\rho_t \nabla V) + \Delta \rho_t$$

- $F(\mu) = \frac{1}{2} SW_2^2(\mu, \nu) + \lambda \mathcal{H}(\mu)$ [Bonnotte, 2013, Liutkus et al., 2019]
- $F(\mu) = \frac{1}{2} MMD^2(\mu, \nu)$ [Arbel et al., 2019]
- $F(\mu) = \frac{1}{2} KSD^2(\mu, \nu)$ [Korba et al., 2021]

- If an associated SDE is known, simulate from it
[Liu et al., 2021, Liutkus et al., 2019, Arbel et al., 2019, Korba et al., 2021]

Examples

Let $F(\mu) = \int V(x)\rho(x)dx + \int \log(\rho(x))\rho(x)dx$,

Gradient Flow solution of:

$$\partial_t \rho_t = \operatorname{div}(\rho_t \nabla V) + \Delta \rho_t$$

Associated SDE (Langevin Equation):

$$dX_t = -\nabla V(X_t)dt + \sqrt{2} dW_t$$

First Algorithms of Resolution of WGFs

If SDE is unknown, need specific numerical methods to solve for Wasserstein Gradient Flows by the JKO Scheme:

- By approximating W_2 :

- By the entropic regularized OT problem + Dykstra's algorithm + $\mu = \sum_{i=1}^n \rho_i \delta_{x_i}$, $(x_i)_i$ grid [Peyré, 2015, Carlier et al., 2017]

$$\forall k, \mu_{k+1}^\tau \in \underset{\mu \in \mathcal{P}_2(\mathbb{R}^p)}{\operatorname{argmin}} \frac{W_\epsilon^2(\mu, \mu_k^\tau)}{2\tau} + F(\mu) \quad (25)$$

- By the dual formulation of the entropic OT problem [Caluya and Halder, 2019, Frogner and Poggio, 2020]
- By using SW_2 [Bonet et al., 2021] + Neural networks g^θ and implicit modeling, i.e. $\mu = g_\#^\theta p_Z$

$$\forall k, \mu_{k+1}^\tau \in \underset{\mu \in \mathcal{P}_2(\mathbb{R}^p)}{\operatorname{argmin}} \frac{SW_2^2(\mu, \mu_k^\tau)}{2\tau} + F(\mu) \quad (26)$$

- By using the dynamic formulation of the transport + grid discretization [Laborde, 2016, Carrillo et al., 2021]
- JKO-ICNN [Alvarez-Melis et al., 2021, Mokrov et al., 2021, Bunne et al., 2021]

Theorem (Brenier's Theorem)

Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, μ absolutely continuous with respect to the Lebesgue measure.

Then, the optimal coupling γ^* is unique and of the form $\gamma^* = (\text{Id}, \nabla \varphi)_\# \mu$ with φ is a convex function.

- Reformulate the problem as:

$$u_{k+1}^\tau \in \operatorname{argmin}_{u \in \text{cvx}} \frac{1}{2\tau} \int \|\nabla u(x) - x\|_2^2 \rho_k^\tau(x) dx + F((\nabla u)_\# \rho_k^\tau)$$

- Implicitly define $\rho_{k+1}^\tau = (\nabla u_{k+1}^\tau)_\# \rho_k^\tau$
- Use Input Convex Neural Networks (ICNN) [Amos et al., 2017] to model the convex functions:

$$\theta_{k+1}^\tau \in \operatorname{argmin}_{\theta \in \{\theta, u_\theta \in \text{ICNN}\}} \frac{1}{2\tau} \int \|\nabla_x u_\theta(x) - x\|_2^2 \rho_k^\tau(x) dx + F((\nabla_x u_\theta)_\# \rho_k^\tau)$$

- Backpropagate through gradient

Resolution of Wasserstein Gradient Flows by the JKO Scheme:

- By approximating W_2
- By using the dynamic formulation of the transport + grid discretization
[Laborde, 2016, Carrillo et al., 2021]
- JKOICNN: Modeling the Monge map with ICNNs
[Alvarez-Melis et al., 2021, Mokrov et al., 2021, Bunne et al., 2021]

$$\theta_{k+1}^\tau \in \operatorname{argmin}_{\tau \in \{\theta, u_\theta \in \text{ICNN}\}} \frac{1}{2\tau} \int \|\nabla u_\theta(x) - x\|_2^2 \, d\mu_k^\tau(x) + F((\nabla_x u_\theta)_\# \mu_k^\tau) \quad (27)$$

- Modeling directly the Monge map

$$T_{k+1}^\tau \in \operatorname{argmin}_T \frac{1}{2\tau} \int \|T(x) - x\|_2^2 \, d\mu_k^\tau(x) + F(T_\# \mu_k^\tau) \quad (28)$$

Optimal transport is a well theoretically grounded ways of comparing probability distributions

- that allows to compare empirical distributions in a non-parametric ways
- that leverages on a ground metric in the embedding space
- for which exist several algorithmic solutions

It comes in several flavours:

- Monge problem: find a mapping (transport map)
- Kantorovich problem: find a coupling (transport plan)

-  Aguech, M. and Carlier, G. (2011).
Barycenters in the wasserstein space.
SIAM Journal on Mathematical Analysis, 43(2):904–924.
-  Alvarez-Melis, D., Schiff, Y., and Mroueh, Y. (2021).
Optimizing functionals on the space of probabilities with input convex neural networks.
-  Ambrosio, L., Gigli, N., and Savaré, G. (2008).
Gradient flows: in metric spaces and in the space of probability measures.
Springer Science & Business Media.
-  Amos, B., Xu, L., and Kolter, J. Z. (2017).
Input convex neural networks.
In *International Conference on Machine Learning*, pages 146–155. PMLR.

-  Arbel, M., Korba, A., Salim, A., and Gretton, A. (2019).
Maximum mean discrepancy gradient flow.
arXiv preprint arXiv:1906.04370.
-  Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015).
Iterative Bregman projections for regularized transportation problems.
SISC.
-  Bigot, J., Gouet, R., Klein, T., López, A., et al. (2017).
Geodesic pca in the wasserstein space by convex pca.
In *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, volume 53, pages 1–26. Institut Henri Poincaré.
-  Bonet, C., Berg, P., Courty, N., Septier, F., Drumetz, L., and Pham, M.-T. (2022).
Spherical sliced-wasserstein.

-  Bonet, C., Courty, N., Septier, F., and Drumetz, L. (2021).
Sliced-wasserstein gradient flows.
-  Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. (2015).
Sliced and radon Wasserstein barycenters of measures.
Journal of Mathematical Imaging and Vision, 51:22–45.
-  Bonnotte, N. (2013).
Unidimensional and evolution methods for optimal transportation.
PhD thesis, Paris 11.
-  Brenier, Y. (1991).
Polar factorization and monotone rearrangement of vector-valued functions.
Communications on pure and applied mathematics, 44(4):375–417.
-  Bunne, C., Meng-Papaxanthos, L., Krause, A., and Cuturi, M. (2021).
Jkonet: Proximal optimal transport modeling of population dynamics.

-  Caluya, K. F. and Halder, A. (2019).
Proximal recursion for solving the fokker-planck equation.
In *2019 American Control Conference (ACC)*, pages 4098–4103. IEEE.
-  Carlier, G., Duval, V., Peyré, G., and Schmitzer, B. (2017).
Convergence of entropic schemes for optimal transport and gradient flows.
SIAM Journal on Mathematical Analysis, 49(2):1385–1418.
-  Carrillo, J. A., Craig, K., Wang, L., and Wei, C. (2021).
Primal dual methods for wasserstein gradient flows.
Foundations of Computational Mathematics, pages 1–55.
-  Courty, N., Flamary, R., and Ducoffe, M. (2017).
Learning wasserstein embeddings.

-  Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2016).
Optimal transport for domain adaptation.
IEEE Transactions on Pattern Analysis and Machine Intelligence.
-  Cuturi, M. (2013).
Sinkhorn distances: Lightspeed computation of optimal transportation.
In *Neural Information Processing Systems (NIPS)*, pages 2292–2300.
-  Delon, J., Salomon, J., and Sobolevski, A. (2010).
Fast transport optimization for monge costs on the circle.
SIAM Journal on Applied Mathematics, 70(7):2239–2258.
-  Dessein, A., Papadakis, N., and Rouas, J.-L. (2016).
Regularized optimal transport and the rot mover's distance.
arXiv preprint arXiv:1610.06447.

-  Ferradans, S., Papadakis, N., Peyré, G., and Aujol, J.-F. (2014).
Regularized discrete optimal transport.
SIAM Journal on Imaging Sciences, 7(3).
-  Frogner, C. and Poggio, T. (2020).
Approximate inference with wasserstein gradient flows.
In *International Conference on Artificial Intelligence and Statistics*, pages 2581–2590. PMLR.
-  Genevay, A., Cuturi, M., Peyré, G., and Bach, F. (2016).
Stochastic optimization for large-scale optimal transport.
In *NIPS*, pages 3432–3440.
-  Hundrieser, S., Klatt, M., and Munk, A. (2021).
The statistics of circular optimal transport.
arXiv preprint arXiv:2103.15426.

-  Jordan, R., Kinderlehrer, D., and Otto, F. (1998).
The variational formulation of the fokker–planck equation.
SIAM journal on mathematical analysis, 29(1):1–17.
-  Kantorovich, L. (1942).
On the translocation of masses.
C.R. (Doklady) Acad. Sci. URSS (N.S.), 37:199–201.
-  Korba, A., Aubin-Frankowski, P.-C., Majewski, S., and Ablin, P. (2021).
Kernel stein discrepancy descent.
arXiv preprint arXiv:2105.09994.
-  Laborde, M. (2016).
Interacting particles systems, Wasserstein gradient flow approach.
PhD thesis, PSL Research University.

-  Liu, S., Sun, H., and Zha, H. (2021).
Approximating the optimal transport plan via particle-evolving method.
-  Liutkus, A., Simsekli, U., Majewski, S., Durmus, A., and Stöter, F.-R. (2019).
Sliced-wasserstein flows: Nonparametric generative modeling via optimal transport and diffusions.
In *International Conference on Machine Learning*, pages 4104–4113. PMLR.
-  McCann, R. J. (1997).
A convexity principle for interacting gases.
Advances in mathematics, 128(1):153–179.
-  Mokrov, P., Korotin, A., Li, L., Genevay, A., Solomon, J., and Burnaev, E. (2021).
Large-scale wasserstein gradient flows.

-  Monge, G. (1781).
Mémoire sur la théorie des déblais et des remblais.
De l'Imprimerie Royale.
-  Peyré, G. (2015).
Entropic approximation of wasserstein gradient flows.
SIAM Journal on Imaging Sciences, 8(4):2323–2351.
-  Rabin, J., Delon, J., and Gousseau, Y. (2011a).
Transportation distances on the circle.
Journal of Mathematical Imaging and Vision, 41(1):147–167.
-  Rabin, J., Peyré, G., Delon, J., and Bernot, M. (2011b).
Wasserstein barycenter and its application to texture mixing.
In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 435–446. Springer.

-  Rubner, Y., Tomasi, C., and Guibas, L. J. (2000).
The earth mover's distance as a metric for image retrieval.
International journal of computer vision, 40(2):99–121.
-  Santambrogio, F. (2014).
Introduction to optimal transport theory.
Notes.
-  Santambrogio, F. (2017).
{Euclidean, metric, and Wasserstein} gradient flows: an overview.
Bulletin of Mathematical Sciences, 7(1):87–154.
-  Schmitz, M. A., Heitz, M., Bonneel, N., Mboula, F. M. N., Coeurjolly, D., Cuturi, M., Peyré, G., and Starck, J.-L. (2017).
Wasserstein dictionary learning: Optimal transport-based unsupervised non-linear dictionary learning.
arXiv preprint arXiv:1708.01955.

-  Seguy, V., Bhushan Damodaran, B., Flamary, R., Courty, N., Rolet, A., and Blondel, M. (2017).
Large-scale optimal transport and mapping estimation.
-  Seguy, V. and Cuturi, M. (2015).
Principal geodesic analysis for probability measures under the optimal transport metric.
In *Advances in Neural Information Processing Systems*, pages 3312–3320.
-  Solomon, J., De Goes, F., Peyré, G., Cuturi, M., Butscher, A., Nguyen, A., Du, T., and Guibas, L. (2015).
Convolutional wasserstein distances: Efficient optimal transportation on geometric domains.
ACM Transactions on Graphics (TOG), 34(4):66.