Introduction to computational Optimal Transport

APM_52188_EP : Emerging topics in machine learning

Rémi Flamary

January 10, 2025

Distributions are everywhere



Distributions are everywhere in machine learning

- Images, vision, graphics, Time series, text, genes, proteins.
- Many datum and datasets can be seen as distributions.
- Important questions:
 - How to compare distributions?
 - How to interpret similarity between distributions?
 - How to use the geometry of distributions?
- Optimal transport provides many tools that can answer those questions.

Illustration from the slides of Gabriel Peyré.

Distributions are everywhere



Distributions are everywhere in machine learning

- Images, vision, graphics, Time series, text, genes, proteins.
- Many datum and datasets can be seen as distributions.
- Important questions:
 - How to compare distributions?
 - How to interpret similarity between distributions?
 - How to use the geometry of distributions?
- $\bullet\,$ Optimal transport provides many tools that can answer those questions.

Illustration from the slides of Gabriel Peyré.

Overview of OTML part of the course

Part 1 : Introduction to optimal transport

- Optimal transport problem
- Wasserstein distance and geometry
- Computational aspects and regularized OT
- Optimal Transport extensions

Part 2 : Learning with optimal transport

- Learning to map with OT
- Learning from histograms
- Learning from empirical distributions
- Learning from structures and across spaces

Table of content (Part 1)

Optimal transport

- Monge and Kantorovitch
- OT on discrete distributions
- Wasserstein distances
- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Special cases: OT in 1D and between Gaussian distributions
- Regularized optimal transport
- Minimizing the Wasserstein distance

Extensions of Optimal Transport

- Partial and Unbalanced Optimal Transport
- Multi-Marginal Optimal Transport (MMOT)
- Gromov-Wasserstein and transport across spaces

Optimal transport

The natural geometry of probability measures

The fathers (and grandfathers of OT):





Monge Kantorovich Koopmans





Dantzig



Brenier



Otto







Figalli Fields '18

McCann









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 $c: \Omega_s \times \Omega_t \to \mathbb{R}^+$.
- The Monge formulation [Monge, 1781] aim at finding a mapping $T: \Omega_s \to \Omega_t$

$$\inf_{T # \boldsymbol{\mu}_{\boldsymbol{s}} = \boldsymbol{\mu}_{\boldsymbol{t}}} \quad \int_{\Omega_{\boldsymbol{s}}} c(\mathbf{x}, T(\mathbf{x})) \boldsymbol{\mu}_{\boldsymbol{s}}(\mathbf{x}) d\mathbf{x}$$
(1)

• Non convex problem because of the constraint $T \# \mu_s = \mu_t$.



Pushforward operator T#

• Transfers measures from one space Ω_s to another space Ω_t

 $\mu_t(A) = \mu_s(T^{-1}(A)), \quad \forall \text{ Borel subset } A \in \Omega_s$

• 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. the change of variable formula

• For empirical measures $\mu_s = \sum_i a_i \delta_{\mathbf{x}_i}$: $T \# \mu_s = \sum_i a_i \delta_{T(\mathbf{x}_i)}$

Properties of mapping T



Non-existence / Non-uniqueness

- $T # \mu_s = \mu_t$ is a non-convex constraint.
- Existence of T is not guaranteed.
- Unicity of T is not guaranteed.
- Very difficult problem in general Prix Bordin of Académie des Sciences (3000F in 1884) never awarded.
- [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
- Focus on where the mass goes, allow splitting [Kantorovich, 1942].
- Applications mainly for resource allocation problems and economics.
- Solution can be found with the simplex algorithm from 1947 [Dantzig, 1990].

Optimal transport (Kantorovich formulation)



The Kantorovich formulation [Kantorovich, 1942] seeks for a probabilistic coupling γ ∈ P(Ω_s × Ω_t) between Ω_s and Ω_t:

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

s.t.
$$\gamma \in \mathcal{P}(\mu_{s}, \mu_{t}) = \left\{ \gamma \geq 0, \ \int_{\Omega_{t}} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mu_{s}, \int_{\Omega_{s}} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \mu_{t} \right\}$$

- γ is a joint probability measure with marginals respectively μ_s and μ_t .
- Linear Program that always has a solution $(\mu_s \otimes \mu_t \in \mathcal{P})$.

Optimal transport (Kantorovich formulation)



The Kantorovich formulation [Kantorovich, 1942] seeks for a probabilistic coupling γ ∈ P(Ω_s × Ω_t) between Ω_s and Ω_t:

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

s.t.
$$\gamma \in \mathcal{P}(\mu_{s}, \mu_{t}) = \left\{ \gamma \geq 0, \ \int_{\Omega_{t}} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mu_{s}, \int_{\Omega_{s}} \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \mu_{t} \right\}$$

- γ is a joint probability measure with marginals respectively μ_s and μ_t .
- Linear Program that always has a solution $(\mu_s \otimes \mu_t \in \mathcal{P})$.

Couplings for 1D distributions



Optimal transport (Kantorovich dual formulation)



Dual formulation of the OT linear program

$$\max_{\phi,\psi} \left\{ \int \phi d\mu_s + \int \psi d\mu_t \mid \phi(\mathbf{x}) + \psi(\mathbf{y}) \le c(\mathbf{x},\mathbf{y}) \right\}$$
(3)

- ϕ and ψ are scalar function also known as Kantorovich potentials.
- Equivalent problem by the Rockafellar-Fenchel theorem.
- Objective value separable wrt μ_s and μ_t .
- Primal-dual relation : the support of $\boldsymbol{\gamma}$ is where $\phi(\mathbf{x}) + \psi(\mathbf{y}) = c(\mathbf{x},\mathbf{y})$

Optimal transport (Kantorovich dual formulation)



Dual formulation of the OT linear program

$$\max_{\phi,\psi} \left\{ \int \phi d\mu_s + \int \psi d\mu_t \mid \phi(\mathbf{x}) + \psi(\mathbf{y}) \le c(\mathbf{x},\mathbf{y}) \right\}$$
(3)

- ϕ and ψ are scalar function also known as Kantorovich potentials.
- Equivalent problem by the Rockafellar-Fenchel theorem.
- Objective value separable wrt μ_s and μ_t .
- Primal-dual relation : the support of γ is where $\phi({\bf x})+\psi({\bf y})=c({\bf x},{\bf y})$

The linear dual constraint suggest that there exits an optimal ψ for a given ϕ .

c-transform (or c-conjugate)

$$\phi^{c}(\mathbf{y}) \stackrel{\text{def}}{=} H^{c}(\phi) = \inf_{\mathbf{x}} \quad c(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x})$$
(4)

Similar a Legendre transform (equal when $c(\mathbf{x}, \mathbf{y}) = \mathbf{x}^{\top} \mathbf{y}$).

Semi-dual formulation

$$\max_{\phi} \quad \left\{ \int \phi d\mu_{s} + \int \phi^{c} d\mu_{t} \right\}$$
(5)

- Depends only on one dual potential through the c-transform.
- Nice reformulation when H^c is easy to compute or closed form.
- Special cases when $c(\mathbf{x},\mathbf{y}) = \|\mathbf{x} \mathbf{y}\|$ and $c(\mathbf{x},\mathbf{y}) = \|\mathbf{x} \mathbf{y}\|^2$.

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



Case $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|$

- 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 dual OT problem can be reformulated as

$$\sup_{\phi \in \text{Lip}^1} \int \phi d(\boldsymbol{\mu}_s - \boldsymbol{\mu}_t) = \sup_{\phi \in \text{Lip}^1} \mathbb{E}_{\mathbf{x} \sim \boldsymbol{\mu}_s}[\phi(x)] - \mathbb{E}_{\mathbf{y} \sim \boldsymbol{\mu}_t}[\phi(y)]$$
(6)

- Also known as Kantorovich-Rubinstein duality
- Formulation used for Wasserstein GAN (more details in next part).

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



Case $c(x, y) = ||x - y||^2/2$

- When μ_s and μ_t are continuous, T(x) the OT mapping exists and is unique.
- More remarkably, it is a gradient of a convex functions $\Phi(x)$

$$T(x) = x - \nabla\phi(x) = \nabla\left(\frac{\|x\|^2}{2} - \phi(x)\right) = \nabla(\Phi(x))$$
(7)

• This is also known as Brenier's Theorem [Brenier, 1991].

Discrete distributions: Empirical vs Histogram

Discrete measure:

$$\mu = \sum_{i=1}^{n} a_i \delta_{\mathbf{x}_i}, \quad \mathbf{x}_i \in \Omega, \quad \sum_{i=1}^{n} a_i = 1$$

Lagrangian (point clouds)



- Constant weight: $a_i = \frac{1}{n}$
- Quotient space: Ω^n , Σ_n

Eulerian (histograms)



- Fixed positions \mathbf{x}_i e.g. grid
- Convex polytope Σ_n (simplex): $\{(a_i)_i \ge 0; \sum_i a_i = 1\}$

The 3 ways of optimal transport



Image from Gabriel Peyré



OT Linear Program When $\mu_s = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_{i=1}^n b_i \delta_{\mathbf{x}_i^t}$

$$\mathbf{T}_{0} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \left\{ \langle \mathbf{T}, \mathbf{C} \rangle_{F} = \sum_{i, j} T_{i, j} c_{i, j} \right\}$$

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

$$\Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \left\{ \mathbf{T} \in (\mathbb{R}^{+})^{n_{s} \times n_{t}} | \mathbf{T} \mathbf{1}_{n_{t}} = \mathbf{a}, \mathbf{T}^{T} \mathbf{1}_{n_{s}} = \mathbf{b} \right\}$$

Linear program with $n_s n_t$ variables and $n_s + n_t$ constraints. Demo



OT Linear Program When $\mu_s = \sum_{i=1}^n \frac{a_i \delta_{\mathbf{x}_i^s}}{a_i}$ and $\mu_t = \sum_{i=1}^n \frac{b_i \delta_{\mathbf{x}_i^t}}{b_i}$

$$\mathbf{T}_{0} = \operatorname*{argmin}_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})} \left\{ \langle \mathbf{T}, \mathbf{C} \rangle_{F} = \sum_{i, j} T_{i, j} c_{i, j} \right\}$$

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

$$\Pi(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \left\{ \mathbf{T} \in (\mathbb{R}^+)^{n_s \times n_t} | \mathbf{T} \mathbf{1}_{n_t} = \mathbf{a}, \mathbf{T}^T \mathbf{1}_{n_s} = \mathbf{b} \right\}$$

Linear program with $n_s n_t$ variables and $n_s + n_t$ constraints. Demo



OT Linear Program When $\mu_s = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_{i=1}^n b_i \delta_{\mathbf{x}_i^t}$

$$\mathbf{T}_{0} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \left\{ \langle \mathbf{T}, \mathbf{C} \rangle_{F} = \sum_{i,j} T_{i,j} c_{i,j} \right\}$$

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

$$\Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \left\{ \mathbf{T} \in (\mathbb{R}^{+})^{n_{s} \times n_{t}} | \mathbf{T} \mathbf{1}_{n_{t}} = \mathbf{a}, \mathbf{T}^{T} \mathbf{1}_{n_{s}} = \mathbf{b} \right\}$$

Linear program with $n_s n_t$ variables and $n_s + n_t$ constraints. Demo



- Π is the Birkhoff polytope (for uniform weights).
- No unique solution in some cases, numerical instabilities
- OT loss not differentiable !

OT Dual for discrete distributions



Discrete OT dual formulation

$$\max_{\boldsymbol{\alpha}\in\mathbb{R}^{n^s},\boldsymbol{\beta}\in\mathbb{R}^{n^t}} \quad \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$
(8)

s.t.
$$\alpha_i + \beta_j \le c_{i,j} \quad \forall i,j$$
 (9)

- With $\mu_s = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_{i=1}^n b_i \delta_{\mathbf{x}_i^t}$
- Linear program with $n_s + n_t$ variables and $n_s n_t$ constraints.
- Solved with Network Flow solver of complexity $O(n^3 \log(n))$ with $n = \max(n_s, n_t).$

Matching words embedding



Word mover's distance [Kusner et al., 2015]

- Words 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.
- Small value of the objective means similar documents.
- OT matrix provide interpretability (word correspondance).

Wasserstein distance



Wasserstein distance

$$W_p^p(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \min_{\boldsymbol{\gamma} \in \mathcal{P}} \quad \int_{\Omega_s \times \Omega_t} \|\mathbf{x} - \mathbf{y}\|^p \boldsymbol{\gamma}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \boldsymbol{\gamma}}[\|\mathbf{x} - \mathbf{y}\|^p]$$
(10)

In this case we have $c(\mathbf{x},\mathbf{y}) = \|\mathbf{x}-\mathbf{y}\|^p$

- 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).

Earth Mover's Distance (EMD)



EMD for image retrieval [Rubner et al., 2000]

- Represent images as histograms.
- Color histogram measure de color proportion
- Histogram of gradient encodes texture.
- FastEMD [Pele and Werman, 2009] is a fast approximation.

Wasserstein barycenter



Barycenters [Agueh and Carlier, 2011]

$$\bar{\mu} = \arg\min_{\mu} \quad \sum_{i}^{n} \lambda_{i} W_{p}^{p}(\mu^{i}, \mu)$$

- $\lambda_i > 0$ 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 $0 \le t \le 1$ [McCann, 1997].
- Regularized barycenters using Bregman projections [Benamou et al., 2015].
- The cost and regularization impacts the interpolation trajectory.

Wasserstein barycenter



Barycenters [Agueh and Carlier, 2011]

$$ar{\mu} = rg \min_{\mu} \quad \sum_{i}^{n} \lambda_{i} W_{p}^{p}(\mu^{i}, \mu)$$

- $\lambda_i > 0$ 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 $0 \le t \le 1$ [McCann, 1997].
- Regularized barycenters using Bregman projections [Benamou et al., 2015].
- The cost and regularization impacts the interpolation trajectory.

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

Illustration from [Kolouri et al., 2017] and maze example from [Papadakis et al., 2014]

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
- Cost between two pixels is the shortest path in the maze (Riemannian metric).

Illustration from [Kolouri et al., 2017] and maze example from [Papadakis et al., 2014]

3D Wasserstein barycenter

Shape interpolation [Solomon et al., 2015]


Wasserstein averaging of fMRI



OT averaging of neurological data [Gramfort et al., 2015]

- Average fMRI activation maps on voxels or cortical surface (natural metric).
- Classical average across subjects and gaussian blur loose information.
- OT averaging recover central activation areas with better precision.
- Can encode both geometrical (3D position) or anatomical connectivity information.
- Extension using OT-Lp seems more robust to noise [Wang et al., 2018].

Outline

Optimal transport

- Monge and Kantorovitch
- OT on discrete distributions
- Wasserstein distances
- Barycenters and geometry of optimal transport

Computational aspects of optimal transport

- Special cases: OT in 1D and between Gaussian distributions
- Regularized optimal transport
- Minimizing the Wasserstein distance
- **Extensions of Optimal Transport**
 - Partial and Unbalanced Optimal Transport
 - Multi-Marginal Optimal Transport (MMOT)
 - Gromov-Wasserstein and transport across spaces



- When c(x, y) is a strictly convex and increasing function of |x y|.
- If $x_1 < x_2$ and $y_1 < y_2$, we have $c(x_1, y_1) + c(x_2, y_2) < c(x_1, y_2) + c(x_2, y_1)$
- The OT plan respects the ordering of the elements.
- Solution is given by the monotone rearrangement of μ_1 onto μ_2 .
- Simple algorithm for discrete distribution by sorting $O(N \log N)$.



- When c(x, y) is a strictly convex and increasing function of |x y|.
- If $x_1 < x_2$ and $y_1 < y_2$, we have $c(x_1, y_1) + c(x_2, y_2) < c(x_1, y_2) + c(x_2, y_1)$
- The OT plan respects the ordering of the elements.
- Solution is given by the monotone rearrangement of μ_1 onto μ_2 .
- Simple algorithm for discrete distribution by sorting $O(N \log N)$.



Illustration with cumulative distributions

- F_{μ} cumulative distribution function of μ : $F_{\mu}(t) = \mu(-\infty, t]$.
- $F_{\mu}^{-1}(q), q \in [0,1]$ is the quantile function: $F_{\mu}^{-1}(q) = \inf\{x \in \mathbb{R} : F_{\mu}(x) \ge q\}.$
- The value of the W_1 Wasserstein distance

$$W_1(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \int_0^1 c(F_{\boldsymbol{\mu}_s}^{-1}(q), F_{\boldsymbol{\mu}_t}^{-1}(q)) dq$$

• Very fast $O(n \log(n))$ computation on discrete distributions.



Illustration with cumulative distributions

- F_{μ} cumulative distribution function of μ : $F_{\mu}(t) = \mu(-\infty, t]$.
- $F_{\mu}^{-1}(q), q \in [0,1]$ is the quantile function: $F_{\mu}^{-1}(q) = \inf\{x \in \mathbb{R} : F_{\mu}(x) \ge q\}.$
- The value of the W_1 Wasserstein distance

$$W_1(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \int_0^1 c(F_{\boldsymbol{\mu}_s}^{-1}(q), F_{\boldsymbol{\mu}_t}^{-1}(q)) dq$$

• Very fast $O(n \log(n))$ computation on discrete distributions.

Sliced Radon Wasserstein



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

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

where \mathcal{R} is the Radon transform $\mathcal{R}(\mu, \theta) = \int_{\mathbb{S}^{d-1}} \mu(\mathbf{x}) \delta(t - \theta^{\top} \mathbf{x}) d\mathbf{x} \ \forall \theta \in \mathbb{S}^{d-1}$

- Can be approximated by discrete sampling of the directions θ .
- Fast 1D wasserstein on 1D projections when d > 1, fast distance estimation and barycenter computation.
- p-sliced Wasserstein distance used for kernel learning between distributions [Kolouri et al., 2016].

Special case: OT between Gaussians (1)



Wasserstein between Gaussian distributions (Bures-Wasserstein)

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

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

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

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

Special case: OT between Gaussians (2)



OT mapping between Gaussian distributions

- $\mu_s \sim \mathcal{N}(\mathbf{m}_1, \Sigma_1)$ and $\mu_t \sim \mathcal{N}(\mathbf{m}_2, \Sigma_2)$
- The optimal map T for $c(\mathbf{x},\mathbf{y}) = \|\mathbf{x}-\mathbf{y}\|_2^2$ is given by

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

with

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

Regularized optimal transport

$$\mathbf{T}_{0}^{\lambda} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \quad \langle \mathbf{T}, \mathbf{C} \rangle_{F} + \lambda \Omega(\mathbf{T}), \qquad (1)$$

Regularization term $\Omega(\mathbf{T})$

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

Why regularize?

- Smooth the "distance" estimation: $W_{\lambda}(\mu_{s},\mu_{t}) = \left< \mathbf{T}_{0}^{\lambda},\mathbf{C} \right>_{F}$
- 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]

$$\mathbf{T}_{0}^{\lambda} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \quad \langle \mathbf{T}, \mathbf{C} \rangle_{F} + \lambda \sum_{i,j} T_{i,j} (\log T_{i,j} - 1)$$

- Regularization with the negative entropy of γ .
- Looses sparsity, gains stability.
- Strictly convex optimization problem.
- Loss and OT matrix are differentiable.

Entropic regularized optimal transport



Entropic regularization [Cuturi, 2013]

$$\mathbf{T}_{0}^{\lambda} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \quad \langle \mathbf{T}, \mathbf{C} \rangle_{F} + \lambda \sum_{i, j} T_{i, j} (\log T_{i, j} - 1)$$

- Regularization with the negative entropy of γ .
- Looses sparsity, gains stability.
- Strictly convex optimization problem.
- Loss and OT matrix are differentiable.

Lagrangian of the optimization problem

$$\mathcal{L}(\mathbf{T}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{ij} T_{ij} C_{ij} + \lambda T_{ij} (\log T_{ij} - 1) + \boldsymbol{\alpha}^{\mathbf{T}} (\mathbf{T} \mathbf{1}_{n_t} - \mathbf{a}) + \boldsymbol{\beta}^{\mathbf{T}} (\mathbf{T}^T \mathbf{1}_{n_s} - \mathbf{b})$$

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

Entropy-regularized transport

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

$$\mathbf{T}_0^{\lambda} = \mathsf{diag}(\mathbf{u}) \exp(-\mathbf{C}/\lambda)\mathsf{diag}(\mathbf{v})$$

- Through the Sinkhorn theorem $\mathsf{diag}(u)$ and $\mathsf{diag}(v)$ exist and are unique.
- Relation with dual variables: $u_i = \exp(\alpha_i/\lambda), \quad v_j = \exp(\beta_j/\lambda).$
- Can be solved by the **Sinkhorn-Knopp** algorithm.

Sinkhorn-Knopp algorithm

 $\begin{array}{l} \textbf{Algorithm 1 Sinkhorn-Knopp Algorithm (SK).} \\ \hline \textbf{Require: } \mathbf{a}, \mathbf{b}, \mathbf{C}, \lambda \\ \mathbf{u}^{(0)} = \mathbf{1}, \mathbf{K} = \exp(-\mathbf{C}/\lambda) \\ \textbf{for } i \text{ in } 1, \dots, n_{it} \textbf{ do} \\ \mathbf{v}^{(i)} = \mathbf{b} \oslash \mathbf{K}^{\top} \mathbf{u}^{(i-1)} \ // \ \text{Update right scaling} \\ \mathbf{u}^{(i)} = \mathbf{a} \oslash \mathbf{K} \mathbf{v}^{(i)} \ // \ \text{Update left scaling} \\ \textbf{end for} \\ \textbf{return } \mathbf{T} = \text{diag}(\mathbf{u}^{(n_{it})}) \mathbf{K} \text{diag}(\mathbf{v}^{(n_{it})}) \end{array}$

- The algorithm performs alternatively a scaling along the rows and columns of $\mathbf{K}=\exp(-\frac{\mathbf{C}}{\lambda})$ to match the desired marginals.
- Complexity $O(kn^2)$, where k iterations are required to reach convergence
- Fast implementation in parallel, GPU friendly
- Convolutive/Heat structure for K [Solomon et al., 2015] for solving OT and barycenters on images/tensors.

Primal formulation of entropic OT

$$\min_{\mathbf{T}\in\Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t})} \quad \langle \mathbf{T},\mathbf{C}\rangle_{F} + \lambda \sum_{i,j} \boldsymbol{\gamma}_{i,j} (\log \boldsymbol{\gamma}_{i,j} - 1)$$

Dual formulation of entropic OT

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta}} \quad \boldsymbol{\alpha}^{T} \mathbf{a} + \boldsymbol{\beta}^{T} \mathbf{b} - \frac{1}{\lambda} \exp\left(\frac{\boldsymbol{\alpha}}{\lambda}\right)^{T} \mathbf{K} \exp\left(\frac{\boldsymbol{\beta}}{\lambda}\right) \qquad \text{with } \mathbf{K} = \exp\left(-\frac{\mathbf{C}}{\lambda}\right) \quad (12)$$

- Sinkhorn algorithm is a gradient ascent on the dual variables.
- Dual problem is unconstrained: stochastic gradient descent (SGD) [Genevay et al., 2016, Seguy et al., 2017] or L-BFGS [Blondel et al., 2017].
- Semi-dual : closed form for β for a fixed α (logsumexp) leads to fast SAG algorithm [Genevay et al., 2016].

Solving entropic OT with Bregman Projections

Kullback Leibler (KL) divergence

$$\mathrm{KL}(\mathbf{T},\rho) = \sum_{ij} T_{ij} \log \frac{T_{ij}}{\rho_{ij}} = <\mathbf{T}, \log \frac{\mathbf{T}}{\rho} >_F,$$

where ${\bf T}$ anf ρ are discrete distributions with the same support.

OT as a Bregman projection [Benamou et al., 2015]

 \mathbf{T}^{\star} is the solution of the following Bregman projection

$$\mathbf{T}^{\star} = \underset{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s}, \boldsymbol{\mu}_{t})}{\operatorname{argmin}} \operatorname{KL}(\mathbf{T}, \mathbf{K}), \quad \text{where } \mathbf{K} = \exp\left(-\frac{C}{\lambda}\right)$$
(13)

- Sinkhorn is an iterative projection scheme, with alternative projections on marginal constraints.
- Generalizes to Barycenter computation [Benamou et al., 2015].
- Also generalizes to other regularization but less efficient (Dykstra's Projection algorithm [Dessein et al., 2016]).

Sinkhorn divergence

Sinkhorn loss

$$W_{\lambda}(\mu_{s},\mu_{t}) = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t})} \quad \langle \mathbf{T}, \mathbf{C} \rangle_{F} + \lambda \sum_{i,j} T_{i,j} \log T_{i,j}$$

- Entropic term has smoothing effect.
- Not a divergence $(W_{\lambda}(\mu, \mu) > 0 \text{ for } \lambda > 0).$

OT loss (aka Sharp Sinkhorn [Luise et al., 2018])

$$OT_{\lambda}(\mu_s,\mu_t) = \left\langle \mathbf{T}_0^{\lambda}, \mathbf{C} \right\rangle_F$$

- T₀^λ is the solution of entropic OT above.
- Not a divergence $(OT_{\lambda}(\mu, \mu) > 0 \text{ for } \lambda > 0).$

Sinkhorn divergence [Genevay et al., 2017]

$$SD_{\lambda}(\mu_s,\mu_t) = W_{\lambda}(\mu_s,\mu_t) - \frac{1}{2}W_{\lambda}(\mu_s,\mu_s) - \frac{1}{2}W_{\lambda}(\mu_t,\mu_t)$$

• True divergence $(SD_{\lambda}(\mu, \mu) = 0)$.

• Better statistical properties as Wasserstein distance [Genevay et al., 2018].

40 / 70

$$\gamma_0^{\lambda} = \operatorname*{argmin}_{\boldsymbol{\gamma} \in \Pi(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t)} \quad \langle \boldsymbol{\gamma}, \mathbf{C} \rangle_F + \lambda \Omega(\boldsymbol{\gamma}),$$

• Group lasso [Courty et al., 2016b]

$$\Omega(\mathbf{T}) = \sum_{g} \sqrt{\sum_{i,j \in \mathcal{G}_g} T_{i,j}^2}$$

Promotes group sparsity (also submodular reg. [Alvarez-Melis et al., 2017])

• Frobenius norm [Blondel et al., 2017]

$$\Omega(\boldsymbol{\gamma}) = \sum_{i,j} T_{i,j}^2$$

Strongly convex regularization that keeps some sparsity in the solution.

• [Dessein et al., 2016]: KL, Itakura Saito, β -divergences.

Solved with Alternative optimization techniques when projection is efficient.

Minimizing the Wasserstein distance



Minimizing the Wasserstein distance

Let $\mu_s = \sum_{i=1}^n a_i \delta_{\mathbf{x}_i^s}$. We seek the minimal Wasserstein estimator:

$$\min_{\boldsymbol{\mu}_{\boldsymbol{s}}} \quad W(\boldsymbol{\mu}_{\boldsymbol{s}}, \boldsymbol{\mu}_{\boldsymbol{t}})$$

In practice for a discrete distribution μ_s there are two ways of doing this:

- Case 1: For a fixed support $\mathbf{X}_s = {\mathbf{x}_i^s}$ find the optimal weights a (Eulerian).
- Case 2: For fixed weights a find the optimal support $\mathbf{X}_s = {\mathbf{x}_i^s}$ (Lagrangian).

Case 1: fixed support $X_s = \{x_i^s\}$



Gradient with respect to weigths a

$$W(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \max_{\boldsymbol{\alpha} \in \mathbb{R}^{n^{s}}, \boldsymbol{\beta} \in \mathbb{R}^{n^{t}}, \boldsymbol{\alpha}_{i} + \boldsymbol{\beta}_{j} \leq c(\mathbf{x}^{s}, \mathbf{x}^{t}_{j})} \boldsymbol{\alpha}^{T} \mathbf{a} + \boldsymbol{\beta}^{T} \mathbf{b}$$
(14)

- $W(\mu_s, \mu_t)$ is convex wrt. **a**
- Dual solution α^* is a sub-gradient : $\alpha^* \in \partial_{\mathbf{a}} W(\mu_s, \mu_t)$
- Entropy regularized: $W(\mu_s, \mu_t)$ is smooth, convex and $\nabla_{\mathbf{a}} W_{\lambda}(\mu_s, \mu_t) = \lambda \log \mathbf{u}$.
- OT loss: ∇_aOT_λ(μ_s, μ_t) computed using the implicit function theorem [Luise et al., 2018].

Case 2: fixed probability masses a



Gradient and update respect to weights $\mathbf{X}_s = {\mathbf{x}_i^s}$ for $c(\mathbf{x}, \mathbf{y}) = ||\mathbf{x} - \mathbf{y}||^2$

$$W_2^2(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t)} \quad \sum_{i,j} T_{i,j} \| \mathbf{x}_i^s - \mathbf{x}_j^t \|^2$$
(15)

- Gradient: $\nabla_{\mathbf{x}_i^s} W_2^2(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = 2\mathbf{x}_i^s 2\frac{1}{a_i} \sum_j T_{i,j} \mathbf{x}_j^t$
- $W_2^2(\mu_s, \mu_t)$ decreases if $\mathbf{X_s} \leftarrow \mathsf{diag}(\mathbf{a}^{-1})\mathbf{T}^*\mathbf{X_t}$
- Expression above called barycentric interpolation [Ferradans et al., 2014a].

Case 2: fixed probability masses a



Gradient and update respect to weights $\mathbf{X}_s = {\mathbf{x}_i^s}$ for $c(\mathbf{x}, \mathbf{y}) = ||\mathbf{x} - \mathbf{y}||^2$

$$W_2^2(\mu_s, \mu_t) = \min_{\mathbf{T} \in \Pi(\mu_s, \mu_t)} \quad \sum_{i,j} T_{i,j} \|\mathbf{x}_i^s - \mathbf{x}_j^t\|^2$$
(15)

- Gradient: $\nabla_{\mathbf{x}_i^s} W_2^2(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t) = 2\mathbf{x}_i^s 2\frac{1}{a_i}\sum_j T_{i,j}\mathbf{x}_j^t$
- $W_2^2(\mu_s, \mu_t)$ decreases if $\mathbf{X_s} \leftarrow \mathsf{diag}(\mathbf{a}^{-1})\mathbf{T}^*\mathbf{X_t}$
- Expression above called barycentric interpolation [Ferradans et al., 2014a].

General case for entropic OT: autodifferentiation



Image from Marco Cuturi

Sinkhorn Autodiff [Genevay et al., 2017]

- Computing gradients through implicit function theorem can be costly [Luise et al., 2018].
- Each iteration of the SInkhorn algorithm is differentiable.
- Modern neural network toolboxes can perform autodiff (Pytorch, Tensorflow).
- Fast but needs log-stabilization for numerical stability.
- At convergence, closed form solution of the gradients exist (no need to autodiff !). 45/70

Outline

Optimal transport

- Monge and Kantorovitch
- OT on discrete distributions
- Wasserstein distances
- Barycenters and geometry of optimal transport
- Computational aspects of optimal transport
 - Special cases: OT in 1D and between Gaussian distributions
 - Regularized optimal transport
 - Minimizing the Wasserstein distance

Extensions of Optimal Transport

- Partial and Unbalanced Optimal Transport
- Multi-Marginal Optimal Transport (MMOT)
- Gromov-Wasserstein and transport across spaces

Relaxation and extensions

- OT is a powerful formulation for several ML applications.
- But as illustrated by entropic regularization, one can also change the optimization problem to get a better/more representative problem.
- Several extensions and variants of OT has been studied by mathematicians and ML practitioners.

Extensions of Optimal Transport

- Partial OT, only a portion of the mass is required to be transported.
- Unbalanced OT, can transport between distributions with different total mass.
- Multi-marginal OT, searches for a transport between more than two distributions.
- Gromov-Wasserstein OT, searches for a transport across metric spaces.
- Co-Optimal Transport, searches for a transport across samples and features.

Partial Optimal Transport



Partial OT [Caffarelli and McCann, 2010, Figalli, 2010]

$$\min_{\mathbf{T}\in\Pi^m(\boldsymbol{\mu}_s,\boldsymbol{\mu}_t)} \left\{ \langle \mathbf{T},\mathbf{C} \rangle_F = \sum_{i,j} T_{i,j} c_{i,j} \right\}$$

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

$$\Pi^{m}(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \left\{ \mathbf{T} \in (\mathbb{R}^{+})^{n_{s} \times n_{t}} | \mathbf{T} \mathbf{1}_{n_{t}} \leq \mathbf{a}, \mathbf{T}^{T} \mathbf{1}_{n_{s}} \leq \mathbf{b}, \mathbf{1}_{n_{s}}^{T} \mathbf{T} \mathbf{1}_{n_{t}} = m \right\}$$

- The equality constraint is on the total transported mass that must be equal to m.
- Allows distributions with different total mass when $m \leq \min(\mathbf{1}_{n_s}^T \mathbf{a}, \mathbf{1}_{n_t}^T \mathbf{b})$ 47/70

Partial OT solver [Figalli, 2010, Chapel et al., 2020]

- Partial OT can be used solved using standard OT solvers using dummy variables.
- The problem to solve is the following

$$\min_{\widetilde{\mathbf{T}}\in\widetilde{\Pi}(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t})} \left\{ \left\langle \widetilde{\mathbf{T}},\widetilde{\mathbf{C}}\right\rangle_{F} = \sum_{i,j}\widetilde{T}_{i,j}\widetilde{c}_{i,j} \right\}$$

where
$$\widetilde{\Pi}(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \left\{ \widetilde{\mathbf{T}} \in (\mathbb{R}^{+})^{n_{s}+1 \times n_{t}+1} | \widetilde{\mathbf{T}} \mathbf{1}_{n_{t}+1} = \widetilde{\mathbf{a}}, \widetilde{\mathbf{T}}^{T} \mathbf{1}_{n_{s}+1} = \widetilde{\mathbf{b}} \right\}$$
 and
 $\widetilde{\mathbf{T}} = \begin{bmatrix} \mathbf{T} & \mathbf{q} \\ \mathbf{p}^{T} & 0 \end{bmatrix}, \ \widetilde{\mathbf{C}} = \begin{bmatrix} \mathbf{C} & \mathbf{0}_{n_{s}} \\ \mathbf{0}_{n_{t}}^{T} & c_{max} \end{bmatrix}, \ \widetilde{\mathbf{a}} = \begin{bmatrix} \mathbf{a} \\ \mathbf{a}^{T} \mathbf{1}_{n_{s}} - m \end{bmatrix}, \ \widetilde{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ \mathbf{b}^{T} \mathbf{1}_{n_{t}} - m \end{bmatrix},$

- where $c_{max} > c_{i,j}$, $\forall i, j$ and \mathbf{p}, \mathbf{q} contains the mass not transported.
- The solution T for Partial OT can be extracted from the solution of the augmented problem.

Unbalanced Optimal Transport



Unbalanced Optimal transport (UOT) [Benamou, 2003]

$$\min_{\mathbf{T} \ge 0} \quad \langle \mathbf{T}, \mathbf{C} \rangle_F + \lambda^u D_{\varphi}(\mathbf{T} \mathbf{1}_m, \mathbf{a}) + \lambda^u D_{\varphi}(\mathbf{T}^{\top} \mathbf{1}_n, \mathbf{b})$$
(16)

- D_{φ} is a Bregman divergence penalizing the violation of the marginal constraints.
- Only a portion of the total mass is transported, total mass can be unbalanced between source and target due to constraint relaxation.
- Balanced problem equivalent to Partial OT when D_{ϕ} is the total variation.
- Closed form exists between Gaussians [Janati et al., 2020, Janati, 2021].

Non regularized UOT

- Smooth convex optimization problem under positivity constraints.
- Classical approach is to use L-BFGS under box constraint [Byrd et al., 1995].
- Problem is actually equivalent to non-negative regression [Chapel et al., 2021]
 - Majorization minimization methods lead to a sinkhorn-like updates without regularization (multiplicative updates).
 - Regularization path can be done with quadratic divergence (Lasso/LARS).

Regularized UOT

$$\min_{\mathbf{T} \ge 0} \quad \langle \mathbf{T}, \mathbf{C} \rangle_F + \lambda^u D_{\varphi}(\mathbf{T} \mathbf{1}_m, \mathbf{a}) + \lambda^u D_{\varphi}(\mathbf{T}^\top \mathbf{1}_n, \mathbf{b}) + \lambda \Omega(\mathbf{T})$$
(17)

- Entropic regularization leads to convex problem [Chizat et al., 2018].
- Can be solved in the dual using block coordinate ascent.
- Algorithm similar to sinkhorn (fast, easy to implement) [Séjourné et al., 2022].
- Can be debiased to get a proper divergence [Séjourné et al., 2019].

Multi-Marginal Optimal Transport (MMOT)



- $\mu_k = \sum_{i=1}^{n_k} a_i^k \delta_{\mathbf{x}_i^k}$ with $k \in \{1, \dots, K\}$ are K discrete distributions.
- $C_{i_1,\ldots,i_K} = c(\mathbf{x}_{i_1}^1,\ldots,\mathbf{x}_{i_K}^K)$ is the MM cost and: $\Pi(\{\mu_k\}_k) = \left\{ \mathbf{T} \in (\mathbb{R}^+)^{n_1 \times \cdots \times n_K} | \sum_{q \neq k} \sum_{i_k=1}^{n_k} T_{i_1,\ldots,i_K} = a_{i_q}, \ \forall q, i_q \in \{1, n_q\} \right\}$
- Properties of MMOT (review in [Pass, 2015])
 - Search for a joint distribution (expressed as a tensor).
 - When K = 2, T is a matrix and we recover classical OT problem.
 - Link to Wasserstein barycenter for specific c [Agueh and Carlier, 2011].

Solving exact MMOT

- Linear program (LP) but with dimensionality exp.in the number of marginal.
- In the primal LP with $\prod_k n_k$ variables and $\sum_k n_k$ constraints.
- Very complex to solve for medium to large scale problems.
- [Tupitsa et al., 2020] use accelerated alternated minimization.
- For specific separable cost in 1D, fast solver [Mehta et al., 2023].

Entropic MMOT

$$\min_{\mathbf{T}\in\Pi(\{\mu_k\}_k)} \quad \sum_{k=1}^K \sum_{i_k=1}^{n_k} T_{i_1,\dots,i_K} C_{i_1,\dots,i_K} + \lambda \Omega(\mathbf{T})$$

- Problem becomes smooth and strictly convex.
- Can be solved using Bregman projections [Benamou et al., 2015].
- The solution is of the form $\mathbf{T} = \exp(-\mathbf{C}/\lambda) \odot \bigotimes_k \mathbf{u}_k$ where \mathbf{u}_k are positive scaling updated at each projections.
- Tensor extension of Sinkhorn algorithm updates \mathbf{u}_k alternatively.

Can you transport accross different spaces ?



- Ω_s : source space, Ω_t : target space.
- Both domains/spaces do not share the same variables.
- There is no $c(\mathbf{x}, \mathbf{y})$ between the two domains.
- They are related (observe similar objects) but not registered.
- Example: multi-modality with observations on different objects.

Gromov-Wasserstein divergence





Inspired from Gabriel Peyré

GW for discrete distributions [Memoli, 2011]

$$\mathcal{GW}_p(\boldsymbol{\mu_s}, \boldsymbol{\mu_t}) = \left(\min_{T \in \Pi(\boldsymbol{\mu_s}, \boldsymbol{\mu_t})} \sum_{i, j, k, l} |\boldsymbol{D}_{i, k} - \boldsymbol{D}'_{j, l}|^p T_{i, j} T_{k, l}\right)^{\frac{1}{p}}$$

with $\mu_s = \sum_i a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_j b_j \delta_{x_j^t}$ and $D_{i,k} = \|\mathbf{x}_i^s - \mathbf{x}_k^s\|, D'_{j,l} = \|\mathbf{x}_j^t - \mathbf{x}_l^t\|$

- Distance between metric measured spaces : across different spaces.
- Search for an OT plan that preserve the pairwise relationships between samples.
- Invariant to isometry in either spaces (e.g. rotations and translation).

Solving the Gromov Wasserstein optimization problem

$$\mathcal{GW}_p^p(\boldsymbol{\mu_s}, \boldsymbol{\mu_t}) = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu_s}, \boldsymbol{\mu_t})} \sum_{i,j,k,l} |D_{i,k} - D'_{j,l}|^p T_{i,j} T_{k,l}$$

with $\mu_s = \sum_i a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_j b_j \delta_{x_j^t}$ and $D_{i,k} = \|\mathbf{x}_i^s - \mathbf{x}_k^s\|, D'_{j,l} = \|\mathbf{x}_j^t - \mathbf{x}_l^t\|$

Optimization problem

- Quadratic Program (Wasserstein is a linear program).
- Nonconvex, NP-hard, related to Quadratic Assignment Problem (QAP).

Optimization algorithm

- Large problem and non convexity forbid standard QP solvers.
- Local solution can be obtained with conditional gradient (Frank-Wolfe) [Vayer et al., 2018] (each iteration is an OT problems).
- Gromov in 1D has a good approximation in close form (solved in discrete with a sort) [Vayer et al., 2019].
- Can be regularized by entropy similarly to classical OT [Peyré et al., 2016a].

Optimization Problem [Peyré et al., 2016a]

$$\mathcal{GW}_{p,\epsilon}^{p}(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t}) = \min_{\mathbf{T}\in\Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t})} \sum_{i,j,k,l} |D_{i,k} - D_{j,l}'|^{p} T_{i,j} T_{k,l} + \epsilon \sum_{i,j} T_{i,j} \log T_{i,j}$$
(18)

with $\mu_s = \sum_i a_i \delta_{\mathbf{x}_i^s}$ and $\mu_t = \sum_j b_j \delta_{x_j^t}$ and $D_{i,k} = \|\mathbf{x}_i^s - \mathbf{x}_k^s\|, D'_{j,l} = \|\mathbf{x}_j^t - \mathbf{x}_l^t\|$

• Smoothing the original GW with a convex and smooth entropic term.

Solving the entropic GW [Peyré et al., 2016a]

- Problem (18) can be solved using a KL mirror descent.
- This is equivalent to solving at each iteration \boldsymbol{t}

$$\mathbf{T}^{(t+1)} = \min_{\boldsymbol{\gamma} \in \mathcal{P}} \quad \left\langle \mathbf{T}, \mathbf{G}^{(t)} \right\rangle_F + \epsilon \sum_{i,j} T_{i,j} \log T_{i,j}$$

Where $G_{i,j}^{(t)} = 2 \sum_{k,l} |D_{i,k} - D'_{j,l}|^p T_{k,l}^{(t)}$ is the gradient of the GW loss at previous point $\mathbf{T}^{(k)}$.

- Problem above can be solved using a Sinkhorn-Knopp algorithm of entropic OT.
- Very fast approximation exist for low rank distances [Scetbon et al., 2021]. 56 / 70

Gromov-Wasserstein between graphs



Modeling the graph structure with a pairwise matrix D

- An undirected graph $\mathcal{G}:=(V,E)$ is defined by $V=\{x_i\}_{i\in[N]}$ set of the N nodes and $E=\{(x_i,x_j)|x_i\leftrightarrow x_j\}$ set of edges.
- Structure represented as a symmetric matrix D of relations between the nodes.
- Possible choices : Adjacency matrix (used in this study), Laplacian matrix, Shortest path matrix.

Graph as a distribution (D, h)



- Graph represented as $\mu_X = \sum_i h_i \delta_{x_i}$.
- The positions x_i are implicit and represented as the pairwise matrix D.
- h_i are the masses on the nodes of the graphs (uniform by default).
Applications of GW [Solomon et al., 2016]

Shape matching between 3D and 2D surfaces



Source

Targets

Multidimensional scaling (MDS) of shape collection



Labeled graphs as distributions



Graph data representation

$$\mu = \sum_{i=1}^{n} h_i \delta_{(x_i a_i)}$$

- Nodes are weighted by their mass h_i .
- But no common metric between the structure points x_i of two different graphs.
- Features values a_i can be compared through the common metric

Fused Gromov-Wasserstein distance



Fused Gromov Wasserstein distance $\mu_s = \sum_{i=1}^n h_i \delta_{x_i, a_i}$ and $\mu_t = \sum_{j=1}^m g_j \delta_{y_j, b_j}$

$$\mathcal{FGW}_{p,q,\alpha}(D,D',\boldsymbol{\mu_s},\boldsymbol{\mu_t}) = \left(\min_{\mathbf{T}\in\Pi(\boldsymbol{\mu_s},\boldsymbol{\mu_t})}\sum_{i,j,k,l} \left((1-\alpha)C_{i,j}^q + \alpha|\boldsymbol{D_{i,k}} - \boldsymbol{D'_{j,l}}|^q\right)^p T_{i,j} T_{k,l}\right)^{\frac{1}{p}}$$

with $D_{i,k} = \|x_i - x_k\|$ and $D'_{j,l} = \|y_i - y_l\|$ and $C_{i,j} = \|a_i - b_j\|$

- Parameters q > 1, $\forall p \ge 1$.
- $\alpha \in [0,1]$ is a trade off parameter between structure and features. 60 / 70

FGW Properties (1)

$$\mathcal{FGW}_{p,q,\alpha}^{p}(D,D',\boldsymbol{\mu_{s}},\boldsymbol{\mu_{t}}) = \min_{\mathbf{T}\in\Pi(\boldsymbol{\mu_{s}},\boldsymbol{\mu_{t}})} \sum_{i,j,k,l} \left((1-\alpha)C_{i,j}^{q} + \alpha |\boldsymbol{D_{i,k}} - \boldsymbol{D'_{j,l}}|^{q} \right)^{p} T_{i,j} T_{k,l}$$

Metric properties [Vayer et al., 2020]

- *FGW* defines a metric over structured data with measure and features preserving isometries as invariants.
- \mathcal{FGW} is a metric for q = 1 a semi metric for q > 1, $\forall p \ge 1$.
- The distance is nul iff :
 - There exists a Monge map $T # \mu_s = \mu_t$.
 - Structures are equivalent through this Monge map (isometry).
 - Features are equal through this Monge map.

Other properties for continuous distributions

- Interpolation between $W(\alpha = 0)$ and $\mathcal{GW}(\alpha = 1)$ distances.
- Geodesic properties (constant speed, unicity).

Computing FGW

Algorithm 2 Conditional Gradient (CG) for FGW

- 1: $\mathbf{T}^{(0)} \leftarrow \mu_X \mu_Y^\top$
- 2: for i = 1, ..., do
- 3: $\mathbf{G} \leftarrow \text{Gradient from Eq. (62) } w.r.t. \mathbf{T}^{(i-1)}$
- 4: $ilde{\mathbf{T}}^{(i)} \leftarrow \mathsf{Solve} \; \mathsf{OT} \; \mathsf{with} \; \mathsf{ground} \; \mathsf{loss} \; \mathbf{G}$
- 5: $\tau^{(i)} \leftarrow \text{Line-search for loss with } \tau \in (0,1)$

6:
$$\mathbf{T}^{(i)} \leftarrow (1 - \tau^{(i)})\mathbf{T}^{(i-1)} + \tau^{(i)}\tilde{\mathbf{T}}^{(i)}$$

7: end for

Algorithmic resolution (p = 1)

$$\mathbf{T}^* = \underset{\mathbf{T} \in \mathcal{P}(\mu_{\mathbf{s}}, \mu_{\mathbf{t}})}{\operatorname{arg\,min}} \quad \mathsf{vec}(\mathbf{T})^T \mathbf{Q} \mathsf{vec}(\mathbf{T}) + \mathsf{vec}((1-\alpha)\mathbf{C})^\top \mathsf{vec}(\mathbf{T}), \quad \text{with } \mathbf{Q} = -2\alpha \mathbf{D}' \otimes \mathbf{D}$$

- Problem is a non-convex Quadratic Program (GW with an additional linear term).
- Conditional gradient [Ferradans et al., 2014b] with network simplex solver.
- Convergence to a local minima [Lacoste-Julien, 2016].
- With entropic regularization, KL mirror descent descent [Peyré et al., 2016b]. 62 / 70

FGW barycenter





$$\min_{x} \sum_{k} \lambda_k \|x - x_k\|^2$$

 $\min_{D \in \mathbb{R}^{n \times n}, \mu} \sum_{i} \lambda_i \mathcal{FGW}(D_i, D, \mu_i, \mu)$

FGW barycenter p = 1, q = 2

- Estimate FGW barycenter using Frechet means (similar to [Peyré et al., 2016a]).
- Barycenter optimization solved via block coordinate descent (on $\mathbf{T}, \mathbf{D}, \{a_i\}_i$).
- Can chose to fix the structure (D) or the features $\{a_i\}_i$ in the barycenter.
- a_{i_i} , and D updates are weighted averages using \mathbf{T} .



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.



- We select a clean graph, change the number of nodes and add label noise and random connections.
- We compute the barycenter on n = 15 and n = 7 nodes.
- Barycenter graph is obtained through thresholding of the D matrix.

Unbalanced Gromov-Wasserstein [Séjourné et al., 2020]

$$\min_{T \in \Pi(\boldsymbol{\mu}_{s},\boldsymbol{\mu}_{t})} \sum_{i,j,k,l} |\boldsymbol{D}_{i,k} - \boldsymbol{D}'_{j,l}|^{p} T_{i,j} T_{k,l} + \lambda^{u} D_{\varphi}(\mathbf{T} \mathbf{1}_{m}, \mathbf{a}) + \lambda^{u} D_{\varphi}(\mathbf{T}^{\top} \mathbf{1}_{n}, \mathbf{b})$$

The marginal constraints are relaxed by penalizing with divergence D_{φ} .

Semi-relaxed GW [Vincent-Cuaz et al., 2022]

$$\min_{T \ge 0, \mathbf{T} \mathbf{1}_m = \mathbf{a}} \quad \sum_{i, j, k, l} \left| \mathbf{D}_{i, k} - \mathbf{D}'_{j, l} \right|^p T_{i, j} T_{k, l}$$

- $\bullet\,$ Second marginal constraint relaxed: optimal weights $\mathbf b$ w.r.t. GW.
- Very fast solver (Frank-Wolfe) because constraints are separable
- Extended to FGW, can eb used to learn a dictionary of graphs (see next course).



Heterogeneous datasets



- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ and $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{n'}]^T \in \mathbb{R}^{n' \times d'}$ contains the source and target data (heer with uniform weights on teh samples).
- Gromov-Wasserstein can be applied across different spaces (focus on pairwise distance).
- OT matrix gives a correspondances of the samples.

Heterogeneous datasets



- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ and $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{n'}]^T \in \mathbb{R}^{n' \times d'}$ contains the source and target data (heer with uniform weights on teh samples).
- Gromov-Wasserstein can be applied across different spaces (focus on pairwise distance).
- OT matrix gives a correspondances of the samples.

Heterogeneous datasets



- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ and $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{n'}]^T \in \mathbb{R}^{n' \times d'}$ contains the source and target data (heer with uniform weights on teh samples).
- Gromov-Wasserstein can be applied across different spaces (focus on pairwise distance).
- OT matrix gives a correspondances of the samples.

Joint samples/features transport





- We want to estimate simultaneously a transport matrix ${\bf T}^s$ between samples and ${\bf T}^v$ a transport matrix between variables.
- \Rightarrow Co-Optimal transport (COOT).

Joint samples/features transport



- We want to estimate simultaneously a transport matrix T^s between samples and T^v a transport matrix between variables.
- \Rightarrow Co-Optimal transport (COOT).

Joint samples/features transport



- We want to estimate simultaneously a transport matrix T^s between samples and T^v a transport matrix between variables.
- \Rightarrow Co-Optimal transport (COOT).

CO-Optimal Transport

Dataset and dimensions

- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ and $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{n'}]^T \in \mathbb{R}^{n' \times d'}$ contains the source and target data.
- $\mathbf{w} \in \Delta_n$ and $\mathbf{w}' \in \Delta_{n'}$ contain the weights of the samples in source and target.
- $\mathbf{v} \in \Delta_d$ and $\mathbf{v}' \in \Delta_{d'}$ contain the weights of the features in source and target.

COOT Optimization problem [Redko et al., 2020]

$$COOT(\mathbf{X}, \mathbf{X}', \mathbf{w}, \mathbf{w}', \mathbf{v}, \mathbf{v}') = \min_{\substack{\mathbf{T}^s \in \Pi(\mathbf{w}, \mathbf{w}') \\ \mathbf{T}^v \in \Pi(\mathbf{v}, \mathbf{v}')}} \sum_{i,j,k,l} L(\mathbf{X}_{i,k}, \mathbf{X}'_{j,l}) \mathbf{T}^s_{i,j} \mathbf{T}^v_{k,l}$$
(19)

- $L(\cdot, \cdot): \mathbb{R}^2 \to \mathbb{R}^+$ is the similarity measure.
- \mathbf{T}^s is the OT matrix between samples, \mathbf{T}^v is the OT matrix between features/variables.
- COOT entropic regularized version adds some entropic terms to the objective value.

Illustration of COOT on real data



COOT between MNIST-USPS datasets

- Sample digits from MNIST 28×28 and USPS 16×16 ordered per classes.
- Uniform weights **w**, **w**' on samples, weights **v**, **v**' on feature is average value.
- Comparison between T from Gromov Wasserstein and COOT T^s: better class correspondence.
- Visualization of \mathbf{T}^s with colors across pixels: spatial structure preserved.

Illustration of COOT on real data



COOT between MNIST-USPS datasets

- Sample digits from MNIST 28×28 and USPS 16×16 ordered per classes.
- Uniform weights \mathbf{w}, \mathbf{w}' on samples, weights \mathbf{v}, \mathbf{v}' on feature is average value.
- Comparison between ${\bf T}$ from Gromov Wasserstein and COOT ${\bf T}^s :$ better class correspondence.
- Visualization of \mathbf{T}^s with colors across pixels: spatial structure preserved.

Summary for Part 1

Optimal transport

- Theoretically grounded ways of comparing probability distributions.
- Non-parametric comparison (between empirical distributions).
- Ground metric encode the geometry of the space (barycenters, geodesic).
- Two aspects: mapping (Monge) vs coupling (Kantorovitch).
- Several variants exists depending on the application.

Optimization

- Solving OT is a linear program.
- Regularization (entropic) leads to faster algorithms.
- Minimization of Wasserstein distance can be done.
- Reference for computational OT : [Peyré et al., 2019]

Next step: how to use it in machine learning applications ?

[Agueh and Carlier, 2011] Agueh, M. and Carlier, G. (2011).

Barycenters in the wasserstein space.

SIAM Journal on Mathematical Analysis, 43(2):904–924.

[Alvarez-Melis et al., 2017] Alvarez-Melis, D., Jaakkola, T. S., and Jegelka, S. (2017). Structured optimal transport.

arXiv preprint arXiv:1712.06199.

[Benamou, 2003] Benamou, J.-D. (2003).

Numerical resolution of an "unbalanced" mass transport problem. ESAIM: Mathematical Modelling and Numerical Analysis, 37(5):851–868.

[Benamou et al., 2015] Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015).

Iterative Bregman projections for regularized transportation problems. *SISC*.

References ii

[Blondel et al., 2017] Blondel, M., Seguy, V., and Rolet, A. (2017).

Smooth and sparse optimal transport.

arXiv preprint arXiv:1710.06276.

[Bonneel et al., 2015] 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.

[Brenier, 1991] Brenier, Y. (1991).

Polar factorization and monotone rearrangement of vector-valued functions. Communications on pure and applied mathematics, 44(4):375–417.

[Byrd et al., 1995] Byrd, R. H., Lu, P., Nocedal, J., and Zhu, C. (1995).

A limited memory algorithm for bound constrained optimization. SIAM Journal on scientific computing, 16(5):1190–1208.

[Caffarelli and McCann, 2010] Caffarelli, L. A. and McCann, R. J. (2010). Free boundaries in optimal transport and monge-ampere obstacle problems. Annals of mathematics, pages 673–730.

References iii

[Chapel et al., 2020] Chapel, L., Alaya, M. Z., and Gasso, G. (2020).

Partial optimal tranport with applications on positive-unlabeled learning. Advances in Neural Information Processing Systems, 33:2903–2913.

[Chapel et al., 2021] Chapel, L., Flamary, R., Wu, H., Févotte, C., and Gasso, G. (2021).
 Unbalanced optimal transport through non-negative penalized linear regression.
 In Neural Information Processing Systems (NeurIPS).

[Chizat et al., 2018] Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. (2018). Scaling algorithms for unbalanced optimal transport problems. *Mathematics of Computation*, 87(314):2563–2609.

[Courty et al., 2016a] Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2016a). Optimal transport for domain adaptation.

IEEE Transactions on Pattern Analysis and Machine Intelligence.

[Courty et al., 2016b] Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2016b). **Optimal transport for domain adaptation.**

Pattern Analysis and Machine Intelligence, IEEE Transactions on.

References iv

[Cuturi, 2013] Cuturi, M. (2013).

Sinkhorn distances: Lightspeed computation of optimal transportation.

In Neural Information Processing Systems (NIPS), pages 2292–2300.

[Dantzig, 1990] Dantzig, G. B. (1990).

Origins of the simplex method.

In A history of scientific computing, pages 141–151.

[Dessein et al., 2016] Dessein, A., Papadakis, N., and Rouas, J.-L. (2016). Regularized optimal transport and the rot mover's distance. arXiv preprint arXiv:1610.06447.

[Ferradans et al., 2014a] Ferradans, S., Papadakis, N., Peyré, G., and Aujol, J.-F. (2014a). Regularized discrete optimal transport.

SIAM Journal on Imaging Sciences, 7(3).

[Ferradans et al., 2014b] Ferradans, S., Papadakis, N., Peyré, G., and Aujol, J.-F. (2014b). Regularized discrete optimal transport.

SIAM Journal on Imaging Sciences, 7(3):1853–1882.

References v

[Figalli, 2010] Figalli, A. (2010).

The optimal partial transport problem.

Archive for rational mechanics and analysis, 195(2):533-560.

[Genevay et al., 2018] Genevay, A., Chizat, L., Bach, F., Cuturi, M., and Peyré, G. (2018). Sample complexity of sinkhorn divergences.

arXiv preprint arXiv:1810.02733.

[Genevay et al., 2016] Genevay, A., Cuturi, M., Peyré, G., and Bach, F. (2016). Stochastic optimization for large-scale optimal transport. In NIPS, pages 3432–3440.

[Genevay et al., 2017] Genevay, A., Peyré, G., and Cuturi, M. (2017). Sinkhorn-autodiff: Tractable wasserstein learning of generative models. arXiv preprint arXiv:1706.00292.

[Gramfort et al., 2015] Gramfort, A., Peyré, G., and Cuturi, M. (2015).

Fast optimal transport averaging of neuroimaging data.

In International Conference on Information Processing in Medical Imaging, pages 261–272. Springer.

[Janati, 2021] Janati, H. (2021).

Advances in Optimal transport and applications to neuroscience. PhD thesis, Institut Polytechnique de Paris.

 [Janati et al., 2020] Janati, H., Muzellec, B., Peyré, G., and Cuturi, M. (2020).
 Entropic optimal transport between unbalanced gaussian measures has a closed form. Advances in Neural Information Processing Systems, 33.

[Kantorovich, 1942] Kantorovich, L. (1942).

On the translocation of masses.

C.R. (Doklady) Acad. Sci. URSS (N.S.), 37:199-201.

[Kolouri et al., 2017] Kolouri, S., Park, S. R., Thorpe, M., Slepcev, D., and Rohde, G. K. (2017).

Optimal mass transport: Signal processing and machine-learning applications.

IEEE signal processing magazine, 34(4):43–59.

References vii

[Kolouri et al., 2016] Kolouri, S., Zou, Y., and Rohde, G. K. (2016).

Sliced wasserstein kernels for probability distributions.

In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 5258–5267.

[Kusner et al., 2015] Kusner, M., Sun, Y., Kolkin, N., and Weinberger, K. (2015).

From word embeddings to document distances.

In International Conference on Machine Learning, pages 957-966.

[Lacoste-Julien, 2016] Lacoste-Julien, S. (2016).

Convergence rate of frank-wolfe for non-convex objectives. arXiv preprint arXiv:1607.00345.

[Luise et al., 2018] Luise, G., Rudi, A., Pontil, M., and Ciliberto, C. (2018).

Differential properties of sinkhorn approximation for learning with wasserstein distance. In Advances in Neural Information Processing Systems, pages 5864–5874.

[McCann, 1997] McCann, R. J. (1997).

A convexity principle for interacting gases. Advances in mathematics, 128(1):153–179.

References viii

[Mehta et al., 2023] Mehta, R., Kline, J., Lokhande, V. S., Fung, G., and Singh, V. (2023). Efficient discrete multi-marginal optimal transport regularization.

```
[Memoli, 2011] Memoli, F. (2011).
```

Gromov wasserstein distances and the metric approach to object matching. *Foundations of Computational Mathematics*, pages 1–71.

```
[Monge, 1781] Monge, G. (1781).
```

Mémoire sur la théorie des déblais et des remblais.

De l'Imprimerie Royale.

[Papadakis et al., 2014] Papadakis, N., Peyré, G., and Oudet, E. (2014).
 Optimal Transport with Proximal Splitting.
 SIAM Journal on Imaging Sciences, 7(1):212–238.

[Pass, 2015] Pass, B. (2015).

Multi-marginal optimal transport: theory and applications. ESAIM: Mathematical Modelling and Numerical Analysis, 49(6):1771–1790.

References ix

[Pele and Werman, 2009] Pele, O. and Werman, M. (2009).

Fast and robust earth mover's distances.

In 2009 IEEE 12th International Conference on Computer Vision, pages 460–467. IEEE.

[Peyré et al., 2019] Peyré, G., Cuturi, M., et al. (2019).

Computational optimal transport: With applications to data science.

Foundations and Trends $\widehat{\mathbb{R}}$ in Machine Learning, 11(5-6):355–607.

[Peyré et al., 2016a] Peyré, G., Cuturi, M., and Solomon, J. (2016a).

Gromov-wasserstein averaging of kernel and distance matrices.

In ICML, pages 2664-2672.

[Peyré et al., 2016b] Peyré, G., Cuturi, M., and Solomon, J. (2016b).

Gromov-Wasserstein Averaging of Kernel and Distance Matrices.

In *ICML 2016*, Proc. 33rd International Conference on Machine Learning, New-York, United States.

[Redko et al., 2020] Redko, I., Vayer, T., Flamary, R., and Courty, N. (2020).

Co-optimal transport.

In Neural Information Processing Systems (NeurIPS).

References x

 [Rubner et al., 2000] 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, 2014] Santambrogio, F. (2014).

Introduction to optimal transport theory.

Notes.

[Scetbon et al., 2021] Scetbon, M., Peyré, G., and Cuturi, M. (2021).

Linear-time gromov wasserstein distances using low rank couplings and costs. arXiv preprint arXiv:2106.01128.

[Seguy et al., 2017] Seguy, V., Bhushan Damodaran, B., Flamary, R., Courty, N., Rolet, A., and Blondel, M. (2017).

Large-scale optimal transport and mapping estimation.

[Séjourné et al., 2019] Séjourné, T., Feydy, J., Vialard, F.-X., Trouvé, A., and Peyré, G. (2019).

Sinkhorn divergences for unbalanced optimal transport.

arXiv preprint arXiv:1910.12958.

References xi

[Séjourné et al., 2020] Séjourné, T., Vialard, F.-X., and Peyré, G. (2020).

The unbalanced gromov wasserstein distance: Conic formulation and relaxation. *arXiv preprint arXiv:2009.04266.*

[Séjourné et al., 2022] Séjourné, T., Vialard, F.-X., and Peyré, G. (2022).

Faster unbalanced optimal transport: Translation invariant sinkhorn and 1-d frank-wolfe. In International Conference on Artificial Intelligence and Statistics, pages 4995–5021. PMLR.

[Solomon et al., 2015] 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.

[Solomon et al., 2016] Solomon, J., Peyré, G., Kim, V. G., and Sra, S. (2016). Entropic metric alignment for correspondence problems. ACM Transactions on Graphics (TOG), 35(4):72.

References xii

- [Tupitsa et al., 2020] Tupitsa, N., Dvurechensky, P., Gasnikov, A., and Uribe, C. A. (2020).
 Multimarginal optimal transport by accelerated alternating minimization.
 In 2020 59th IEEE Conference on Decision and Control (CDC), pages 6132–6137. IEEE.
- [Vayer et al., 2018] Vayer, T., Chapel, L., Flamary, R., Tavenard, R., and Courty, N. (2018). Fused gromov-wasserstein distance for structured objects: theoretical foundations and mathematical properties.
- [Vayer et al., 2020] Vayer, T., Chapel, L., Flamary, R., Tavenard, R., and Courty, N. (2020).
 Fused gromov-wasserstein distance for structured objects.
 Algorithms, 13 (9):212.
- [Vayer et al., 2019] Vayer, T., Flamary, R., Tavenard, R., Chapel, L., and Courty, N. (2019). Sliced gromov-wasserstein.
 - In Neural Information Processing Systems (NeurIPS).
- [Vincent-Cuaz et al., 2022] Vincent-Cuaz, C., Flamary, R., Corneli, M., Vayer, T., and Courty, N. (2022).
 - Semi-relaxed gromov wasserstein divergence with applications on graphs.
 - In International Conference on Learning Representations (ICLR).

[Wang et al., 2018] Wang, Q., Redko, I., and Takerkart, S. (2018).

Population averaging of neuroimaging data using l p distance-based optimal transport. In 2018 International Workshop on Pattern Recognition in Neuroimaging (PRNI), pages 1–4. IEEE.