

Applications of Optimal Transport

Computational Optimal Transport in Imaging Science
SIAM Conference on Imaging Science

Matthew Thorpe and Yunan Yang

28th May 2024



- ① Optimal Transport is touching a huge number of applications - impossible to talk about all of them here.
- ② There is a bias towards applications I've worked on and applications that my collaborators have worked on — but I have tried to be broader than my own interests...
- ③ ...which means some of the applications I know very little about!

Contents

- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
- 4 The Linear Wasserstein Distance
- 5 The Hellinger–Kantorovich Distance
- 6 The TL^P Distance

Contents

- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
- 4 The Linear Wasserstein Distance
- 5 The Hellinger–Kantorovich Distance
- 6 The TL^P Distance

- ① Optimal transport is a linear programme, so one can use **linear programming algorithms**, e.g. the simplex method or interior point methods.
- ② Angenent, Haker and Tannenbaum's¹ **Flow Minimization** method works directly with the Monge formulation. Starting with an initial transport map (usually the Knoth-Rosenblatt coupling) it is an iterative scheme to update the transport map with one with lower cost (following gradient descent of the Monge cost).
- ③ 1D optimal transport is very cheap. Meng et al.'s² uses iterative 1D projections in their **projection pursuit Monge map** method to find the OT map and cost.

¹ **Angenent, Haker and Tannenbaum**, *Minimizing flows for the Monge–Kantorovich problem*, SIAM journal on mathematical analysis, 35(1):61–97, 2003.

² **Meng, Ke, Zhang, Zhong and Ma**, *Large-Scale Optimal Transport Map Estimation Using Projection Pursuit*, NeurIPS 32, 2019.

- ④ One of the most popular current approaches, due to Cuturi³ is to add entropic regularisation to the cost function. The cost can then be reformulated as a Kullback-Leibler divergence which can be solved via the **Sinkhorn Algorithm**.
- ⑤ The **back-and-forth method** of Jacobs and Léger⁴ alternates between optimising two equivalent dual formulations. If ϕ and ψ are the dual optimal solutions then $\psi = \phi^c$ and $\phi = \psi^c$ leading to a dual problem that can either be written in terms of ϕ or in terms of ψ . Taking alternate H^1 gradients will find the dual optimal pair, from which the transport plans / cost can be recovered.

³Cuturi, *Sinkhorn Distances: Lightspeed Computation of Optimal Transport*, In Advances in Neural Information Processing Systems, pp. 2292–2300, 2013.

⁴Jacobs and Léger, *A Fast Approach to Optimal Transport; The Back-And-Forth Method*, Numerische Mathematik, 146(3):513–544, 2020.

Computational Methods III

- ⑤ Plenty of others: e.g. **fluid dynamics⁵, dual gradient descent⁶, Monge-Ampere solvers⁷ and semi-discrete approaches⁸.**

⁵**Benamou and Brenier**, *A Computational Fluid Mechanics Solution to the Monge-Kantorovich Mass Transfer Problem*, Numerische Mathematik, 84(3):375–393, 2000.

⁶**Chartrand, Vixie, Wohlberg and Boltt**, *A Gradient Descent Solution to the Monge-Kantorovich Problem*, Applied Mathematical Sciences, 3(22):1071–1080, 2009.

⁷**Benamou, Froese and Oberman**, *Numerical Solution of the Optimal Transportation Problem Using the Monge-Ampere Equation*, Journal of Computational Physics, 260:107–126, 2014.

⁸**Levy**, *A Numerical Algorithm for L_2 Semi-Discrete Optimal Transport in 3D*, ESAIM Math. Model. Numer. Anal., 49(6):1693–1715, 2015.

Popular Optimal Transport Packages for Python

- **GeomLoss**: Geometric loss Functions between sampled measures, images and volumes. Package authors: Jean Feydy and Pierre Roussillon.
- **OTT-Jax**: optimal transport tools in JAX. Main contributors: Marco Cuturi and Michal Klein.
- **POT**: Python Optimal Transport. Package creators and maintainers: Nicolas Courty, Remi Flamary and Cédric Vincent-Cuaz.
- **PPMP**: Projection pursuit Monge map. Package creators: Cheng Meng, Jingyi Zhang, Mengrui Zhang and Tao Li (built on POT).
- **pysdot**: Python semi-discrete optimal transportation tools. Package contributors: Hugo Leclerc and Quentin Mérigot.

Contents

1 Computational Methods

2 The Wasserstein Distance

- The Wasserstein Distance
- Seismic Imaging
- Wasserstein Generative Adversarial Networks

3 The Sliced Wasserstein Distance

4 The Linear Wasserstein Distance

5 The Hellinger–Kantorovich Distance

6 The TL^P Distance

The Wasserstein Distance

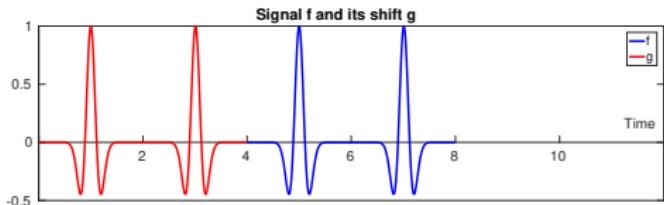
- The Wasserstein distance on the space $\mathcal{P}_p(X)$ is defined by

$$d_{W_p}(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \left(\int_{X \times X} |x - y|^p d\pi(x, y) \right)^{\frac{1}{p}}.$$

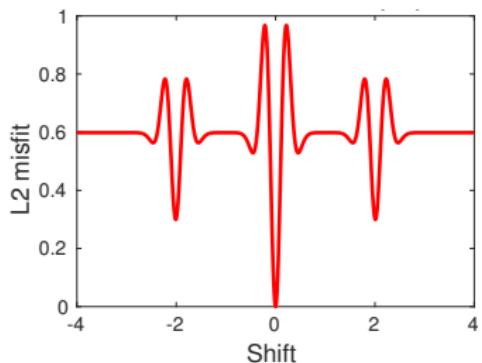
- The Wasserstein distance metrises the weak* convergence (for compact X).
- For $p = 1$ we can also represent the Wasserstein distance in the Kantorovich–Rubinstein duality form:

$$d_{W_1}(\mu, \nu) = \sup \left\{ \int_X f d(\mu - \nu) : f \text{ is 1-Lipschitz} \right\}.$$

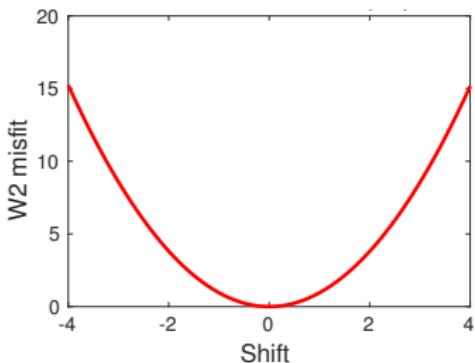
Why the Wasserstein Distance?



(a) Example Signals.



(b) L^2 distance between signals.



(c) W_2 distance between signals.

Source: **Engquist and Yang**, *Seismic Imaging and Optimal Transport*, Communications in Information and Systems, 19(2):95–145, 2019.

Example 1: Seismic Imaging

- A simple model for wave propagation is:

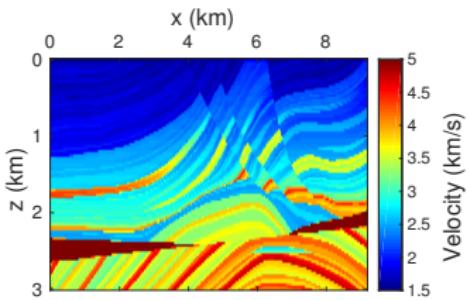
$$\begin{cases} \frac{1}{v(x)^2} \frac{\partial^2 u}{\partial t^2}(x, t) - \Delta u(x, t) = s(x, t) \\ u(x, 0) = 0 \\ \frac{\partial u}{\partial t}(x, 0) = 0 \end{cases}$$

where u is the wavefield, s is the source and $v(x)$ describes the ground material medium (i.e., the wave speed).

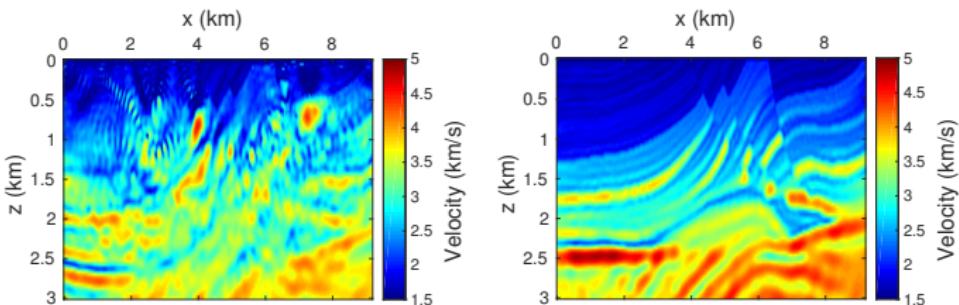
- s is known and u can be observed at locations $\{x_r\}_{r=1}^R$, the **inverse problem** is to find $v(x)$.
- The **forward model** $v \mapsto u(\cdot; v)$ is highly nonlinear.
- Let $\{g(x_r, \cdot)\}_{r=1}^R$ be observations of the wavefield. One can find v by minimising

$$\mathcal{J}(v) = \sum_{r=1}^R d(g(x_r, \cdot), u(x_r, \cdot; v))^2.$$

Seismic Imaging with the Wasserstein Distance



(a) True velocity



(b) L^2 cost:
 $d(f, g)^2 = \int |f(t) - g(t)|^2 dt$

(c) W_2 cost: $d(f, g)^2 = \inf_{\pi \in \Pi(\tilde{f}, \tilde{g})} \int |p - q|^2 d\pi(p, q)$

Source: Engquist and Yang, *Seismic Imaging and Optimal Transport*,
Communications in Information and Systems, 19(2):95–145, 2019.

Example 2: Generative Adversarial Networks

- GANs were introduced in 2014 by Goodfellow et al.⁹ for generative machine learning.
- Core idea: two neural networks compete against each other.
- Generator: creates new images $\mu_G \in \mathcal{P}(\Omega)$.
- Discriminator: tries to determine whether an image is generated or not $D : \Omega \rightarrow [0, 1]$.
- Objective functional:

$$\mathcal{L}_{\text{JS}}(\mu_G, D) = \mathbb{E}_{x \sim \mu_{\text{true}}} [\ln D(x)] + \mathbb{E}_{x \sim \mu_G} [\ln(1 - D(x))].$$

- Minimise \mathcal{L}_{JS} with respect to μ_G and maximise with respect to D .

⁹Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville and Bengio, *Generative Adversarial Networks*, arxiv:1406.2661, 2014.

- Mode collapse happens when the generator learns too quickly.
- The generator essentially creates images in the set $\operatorname{argmax}_x D(x)$.
- If the discriminator prefers a certain region of Ω then the generator will only sample from this region.
- For example, in MNIST if the discriminator prefers the digit 0 then the generator may only generate 0's.
- Even as training continues the generator can remain stuck in a local minimum and fail to generalise.

Wasserstein Generative Adversarial Networks

- Instead of \mathcal{L}_{JS} Arjovsky et al.¹⁰ proposed the following objective:

$$\mathcal{L}_W(\mu_G, D) = \mathbb{E}_{x \sim \mu_G}[D(x)] - \mathbb{E}_{x \sim \mu_{\text{true}}}[D(x)].$$

- Maximising over D that are 1-Lipschitz produces the Wasserstein distance:

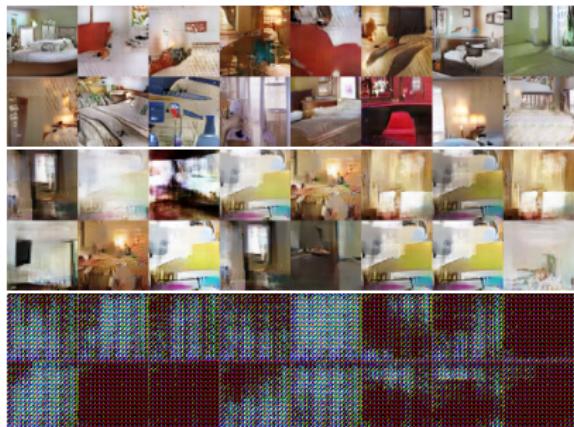
$$\sup_{\|D\|_{\text{Lip}} \leq 1} \mathcal{L}_W(\mu_G, D) = d_{W_1}(\mu_G, \mu_{\text{true}}).$$

- WGANS are empirically observed as more robust with respect to mode collapse.
- However, WGANS are a poor approximation of the Wasserstein distance, and improving the approximation often leads to worse performance.¹¹

¹⁰ **Arjovsky, Chintala and Bottou**, *Wasserstein Generative Adversarial Networks*, International Conference on Machine Learning, 2017.

¹¹ **Stanczuk, Etmann, Kreusser and Schönlieb**, *Wasserstein GANs Work Because They Fail (To Approximate the Wasserstein distance)*, arxiv:2103.01678, 2021.

Generative Adversarial Network



Wasserstein Generative Adversarial Network



Source: **Arjovsky, Chintala and Bottou**, *Wasserstein Generative Adversarial Networks*, International Conference on Machine Learning, 2017.

Contents

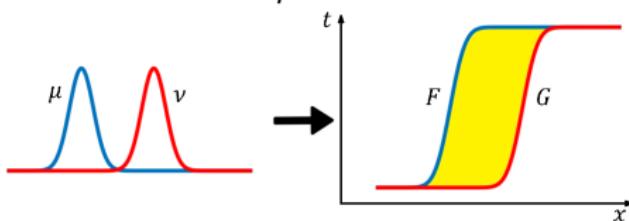
- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
 - The Sliced Wasserstein Distance
 - Sliced Wasserstein Autoencoders
- 4 The Linear Wasserstein Distance
- 5 The Hellinger–Kantorovich Distance
- 6 The TL^P Distance

1D Optimal Transport

- When we have the cumulative distribution functions:

$$d_{W_p}(\mu, \nu)^p = \int_0^1 |F^{-1}(t) - G^{-1}(t)|^p dt$$

where F, G are the cdf's of μ and ν .



- For uniform empirical distributions computing the distance is equivalent to a sorting problem.

$$d_{W_p}(\mu, \nu)^p = \frac{1}{n} \sum_{m=1}^n |x_{i[m]} - y_{j[m]}|^p$$

where $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$, $\nu = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$ and $i[m], j[m]$ are the sorted indices.

Source: **Kolouri and Rohde**, *Optimal transport a crash course*, IEEE ICIP 2016
Tutorial Slides: Part 1, 2016.

The Sliced Wasserstein Distance

- **Sliced Wasserstein Distance:**

- ① Let $\sigma \in \mathcal{P}(\mathbb{S}^{d-1})$ be a positive measure on the hypersphere $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : |x| = 1\}$.
- ② $\langle \theta, \cdot \rangle$ is the projection onto the 1D line spanned by $\theta \in \mathbb{S}^{d-1}$.
- ③ $d_{SW_p}(\mu, \nu) = \int_{\mathbb{S}^{d-1}} d_W_p(\langle \theta, \cdot \rangle_\# \mu, \langle \theta, \cdot \rangle_\# \nu) d\sigma(\theta)$.

- $d_{SW_p} : \mathcal{P}_p(\mathbb{R}^d) \times \mathcal{P}_p(\mathbb{R}^d) \rightarrow [0, \infty)$ is a metric.
- Approximate computational costs:

- ① Linear programming $O(n^3 \log(n))$,
- ② Auction algorithm $O(n^2 \log(n))$,
- ③ Entropy regularised Sinkhorn algorithm $O(n^2)$,
- ④ 1D $O(n \log(n))$.

- An autoencoder aims to generate data by learning a transformation from a samplable (latent) space.
- Let \mathcal{Z} be the latent space and \mathcal{X} be the data space.
- Let $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ be the encoder and $\psi : \mathcal{Z} \rightarrow \mathcal{X}$ be the decoder.
- Let $\mathbb{P}_{\mathcal{X}}$ be the data distribution on \mathcal{X} and $\mathbb{Q}_{\mathcal{Z}}$ be a given samplable distribution on \mathcal{Z} .
- **Properties of a good encoder:**
 - ➊ $\mathbb{Q}_{\mathcal{Z}}$ is easy to sample from;
 - ➋ $\phi_{\#}\mathbb{P}_{\mathcal{X}} \approx \mathbb{Q}_{\mathcal{Z}}$; and
 - ➌ $\psi_{\#}\phi_{\#}\mathbb{P}_{\mathcal{X}} \approx \mathbb{P}_{\mathcal{X}}$.

Example 3: Sliced Wasserstein Autoencoders

- We need a measure for $\phi_{\#}\mathbb{P}_{\mathcal{X}} \approx \mathbb{Q}_{\mathcal{Z}}$.
 - ➊ f-divergences, e.g. Nowozin et al. 2016.¹²
 - ➋ Wasserstein distances, e.g. Tolstikhin et al. 2018.¹³
 - ➌ Sliced Wasserstein distances, e.g. Kolouri et al. 2018.¹⁴
- Sliced Wasserstein autoencoder:

$$\operatorname{argmin}_{\phi, \psi} \frac{1}{n} \sum_{i=1}^n |x_i - \psi(\phi(x_i))|^2 + \lambda d_{\text{SW}_2}(\phi_{\#}\mathbb{P}_{\mathcal{X}}, \mathbb{Q}_{\mathcal{Z}})^2.$$

¹² **Nowozin, Cseke and Tomioka**, *f-GAN: Training generative neural samplers using variational divergence minimization*, NeuRIPS, 2016.

¹³ **Tolstikhin, Bousquet, Gelly and Schölkopf**, *Wasserstein Auto-Encoders*, International Conference on Learning Representations, 2018.

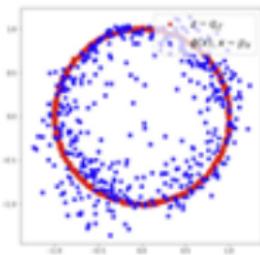
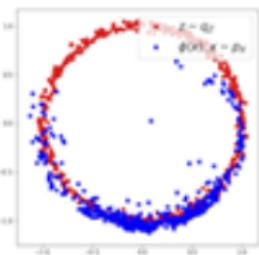
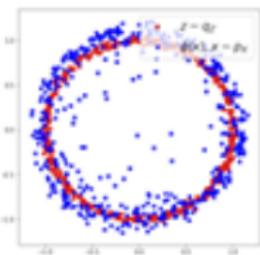
¹⁴ **Kolouri, Pope, Martin and Rohde**, *Sliced Wasserstein auto-encoders*, ICLR, 2018.

Representations of MNIST

5	9	8	9	7	9	9	7	1	2	0
0	2	9	5	1	6	8	2	6	1	
8	5	7	7	5	3	0	5	0	1	
8	2	0	6	1	0	7	1	0	9	
0	1	1	0	6	5	1	1	0	8	
0	9	5	5	0	6	9	5	1	0	
1	5	1	5	9	7	9	7	0	3	
0	1	0	9	7	1	5	2	9	2	
0	1	8	2	0	9	0	8	9	6	
1	7	1	9	8	6	9	1	2	1	

7	5	2	7	7	9	7	1	2	0	
5	3	7	3	7	4	7	3	3	3	
9	3	7	9	5	7	0	9	0	7	
9	3	0	5	1	4	7	1	0	3	
0	7	3	0	5	3	1	1	3	3	
0	9	1	3	0	6	9	3	9	3	
1	5	5	3	9	7	7	7	6	5	
0	1	0	7	9	9	3	3	3	3	
0	6	2	5	5	9	9	5	6	9	
9	7	7	9	7	7	6	9	1	6	

7	5	8	7	9	4	7	1	2	0	
5	8	7	5	9	4	9	2	5	8	
9	5	9	7	5	9	8	9	5	9	
5	2	5	8	1	4	7	8	0	8	
5	9	8	0	5	5	1	1	3	5	
5	9	1	5	0	4	8	5	9	8	
1	5	8	5	9	7	9	9	6	8	
5	1	0	9	9	4	5	2	8	8	
4	5	5	7	9	5	4	1	5		
9	7	9	4	7	7	6	9	1	6	



Left: sliced Wasserstein autoencoder; **centre:** Wasserstein generative adversarial Network; **right:** maximum mean discrepancy autoencoder.
Top: generated images; **bottom:** latent space samples.

Source: Kolouri, Pope, Martin and Rohde, *Sliced Wasserstein auto-encoders*, ICLR, 2018.

Contents

- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
- 4 The Linear Wasserstein Distance
 - The Linear Wasserstein Distance
 - Transport Based Morphometry
 - Data Generating
- 5 The Hellinger–Kantorovich Distance
- 6 The TL^P Distance

The Wasserstein distance is *great* as a distance between signals/images, because...

- ① Lagrangian modelling,
- ② simple to understand compared to other Lagrangian methods such as large deformation diffeomorphic metric mapping,
- ③ metric properties (in particular symmetry, triangle inequality,...).
- ④ geodesics and Riemannian structure,
- ⑤ theoretical and characterising properties such as existence of optimal transport maps and optimal transport plans.

But,...

- ① it places restrictive conditions on the input, in particular signals have to be probability measures,
- ② computationally expensive (despite recent advances),
- ③ there is a lack of off-the-shelf data analysis tools.

This motivates the **linear Wasserstein distance**.

The Linear Wasserstein Distance

- We recall that the Wasserstein distance can be written

$$d_{W_2}(\mu, \nu) = \sqrt{\int |v_0(x)|^2 d\mu(x)}$$

where $v_0 = v(0, \cdot)$ is the Benamou–Brenier optimal velocity.

- We also have $v_0(x) = T(x) - x$ where T is the optimal Monge map.
- We define the logarithmic map:

$$\text{Log}_{W_2}(\mu; \mu_i) := v_0.$$

- Wasserstein distance from reference:

$$d_{W_2}(\mu, \nu) = \|\text{Log}_{W_2}(\mu; \nu)\|_{L^2(\mu)}.$$

- Linear Wasserstein distance:¹⁵

$$d_{W_2,\text{lin},\mu}(\mu_1, \mu_2) = \|\text{Log}_{W_2}(\mu; \mu_1) - \text{Log}_{W_2}(\mu; \mu_2)\|_{L^2(\mu)}.$$

¹⁵ Wang, Slepčev, Basu, Ozolek and Rohde, *A Linear Optimal Transportation Framework for Quantifying and Visualizing Variations in Sets of Images*, International Journal of Computer Vision 101(2):254–269, 2013.

The Linear Wasserstein Distance Approximation

- We always have

$$d_{W_2}(\mu_1, \mu_2) \leq d_{W_2, \text{lin}, \mu}(\mu_1, \mu_2).$$

- If $\mu_i = S_i \# \mu$ where S_i is a shearing and/or translation:¹⁶

$$d_{W_2}(\mu_1, \mu_2) = d_{W_2, \text{lin}, \mu}(\mu_1, \mu_2).$$

- If μ_i live on a nice submanifold:¹⁷

$$d_{W_2}(\mu_1, \mu_2) = d_{W_2, \text{lin}, \mu}(\mu_1, \mu_2) + O(d_{W_2}(\mu_1, \mu_2)^{\frac{3}{2}}).$$

- In general:¹⁸

$$d_{W_2, \text{lin}, \mu}(\mu_1, \mu_2) \leq C d_{W_2}(\mu_1, \mu_2)^{\frac{1}{2}}.$$

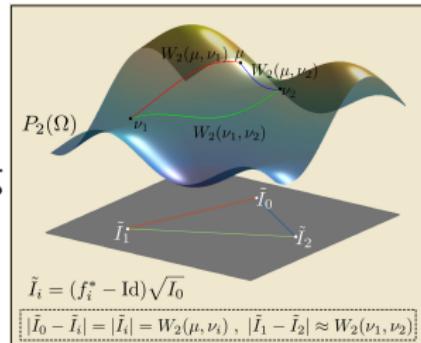


Image Source: **Kolouri, Park, T., Slepčev and Rohde**, *Optimal Mass Transport: Signal Processing and Machine Learning Applications*, IEEE Signal Processing Magazine, 34(4):43–59, 2017.

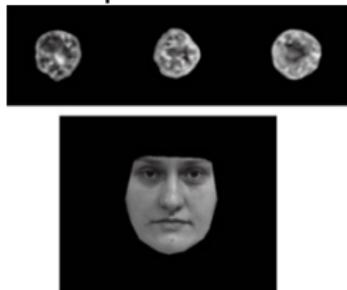
¹⁶ **Khurana, Kannan, Cloninger and Moosmüller**, *Supervised Learning of Sheared Distributions Using Linearized Optimal Transport*, Sampling Theory, Signal Processing, and Data Analysis, 21(1), 2023.

¹⁷ **Hamm, Moosmueller, Schmitzer and T.**, *Manifold Learning in Wasserstein Space*, arxiv:2311.08549, 2023.

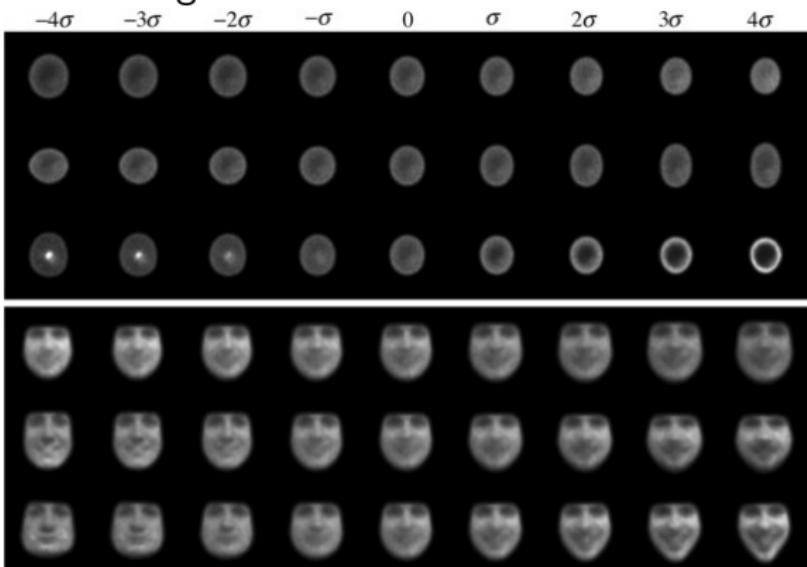
¹⁸ **Gigli**, *On Hölder Continuity-in-Time of the Optimal Transport Map Towards Measures Along a Curve*, Proceedings of the Edinburgh Mathematical Society, 54(2):401–409, 2011.

Example 4: Transport Based Morphometry

Example Data:



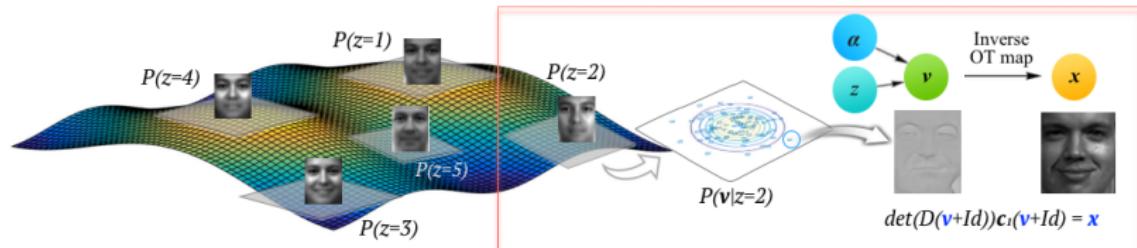
Principle Component Analysis on Linear Embedding:



Source: Wang, Slepčev, Basu, Ozolek and Rohde, *A Linear Optimal Transportation Framework for Quantifying and Visualizing Variations in Sets of Images*, International Journal of Computer Vision 101(2):254–269, 2013.

Example 5: Data Generating

Pipeline:



Results:



- ① Top row, all 19 original images.
- ② Second and third rows, generated images.

Source: Park and T., *Representing and Learning High Dimensional Data with the Optimal Transport Map from a Probabilistic Viewpoint*, CVPR, 2018.

Contents

- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
- 4 The Linear Wasserstein Distance
- 5 The Hellinger–Kantorovich Distance
 - The Hellinger–Kantorovich Distance
 - The Linear Hellinger–Kantorovich Distance
 - Point Cloud Classification
- 6 The TL^P Distance

The Hellinger–Kantorovich Distance

- ① Recall the continuity equation:

$$(\rho, \nu) \in \mathcal{CE}(\mu, \nu) \Leftrightarrow \frac{\partial \rho}{\partial t} + \nabla_x \cdot (\rho \nu) = 0, \rho_0 = \mu, \rho_1 = \nu.$$

- ② And the Wasserstein distance

$$d_{W_2}^2(\mu, \nu) = \inf_{(\rho, \nu) \in \mathcal{CE}(\mu, \nu)} \int_0^1 \int_{\Omega} \|\nu(t, \cdot)\|^2 d\rho_t dt.$$

- ③ To allow for mass creation we include a source/sink term:

$$(\rho, \nu, \zeta) \in \mathcal{CES}(\mu, \nu) \Leftrightarrow \frac{\partial \rho}{\partial t} + \nabla_x \cdot (\rho \nu) = \zeta, \rho_0 = \mu, \rho_1 = \nu.$$

- ④ The Hellinger–Kantorovich distance (a.k.a. unbalanced optimal transport and Wasserstein–Fisher–Rao):¹⁹²⁰²¹²²

$$d_{HK}^2(\mu, \nu) := \inf_{(\rho, \nu, \zeta) \in \mathcal{CES}(\mu, \nu)} \int_0^1 \int_{\Omega} \left(\|\nu(t, \cdot)\|^2 + \frac{1}{4} \left(\frac{d\zeta_t}{d\rho_t} \right)^2 \right) d\rho_t dt.$$

¹⁹ Kondratyev, Monsaingeon and Vorotnikov, *A New Optimal Transport Distance On The Space Of Finite Radon Measures*, Advances in Differential Equations, 21:1117–1164, 2016.

²⁰ Chizat, Peyré, Schmitzer and Vialard, *An Interpolating Distance Between Optimal Transport and Fisher–Rao Metrics*, Foundations of Computational Mathematics 18:1–44, 2018.

²¹ Chizat, Peyré, Schmitzer and Vialard, *Unbalanced Optimal Transport: Dynamic and Kantorovich formulations*, Journal of Functional Analysis, 274:3090–3123, 2018.

²² Liero, Mielke and Savaré, *Optimal Entropy-Transport Problems and a New Hellinger–Kantorovich Distance Between Positive Measures*, Inventiones Mathematicae 211:969–1117, 2018.

The Hellinger–Kantorovich Distance Between Diracs I

- ① Let $\mu = m_0 \delta_{x_0}$ and $\nu = m_1 \delta_{x_1}$.
- ② If $\|x_0 - x_1\| \leq \frac{\pi}{2}$ then we choose $\rho(t, x) = m(t) \delta_{x(t)}$ for some functions $m : [0, 1] \rightarrow [0, +\infty)$ and $x : [0, 1] \rightarrow \Omega$.
- ③ Further assume that trajectories are straight lines, i.e. there exists $\varphi : [0, 1] \rightarrow [0, 1]$ such that $x(t) = x_0 + (x_1 - x_0)\varphi(t)$.
- ④ The HK energy is

$$d_{HK}^2(\mu, \nu) = \inf_{m,x} \underbrace{\int_0^1 \left(|\dot{x}(t)|^2 + \frac{1}{4} \left(\frac{\dot{m}(t)}{m(t)} \right)^2 \right) m(t) dt}_{=: \mathcal{HK}(m,x)}$$

The Hellinger–Kantorovich Distance Between Diracs II

- ⑤ Taking variations in \mathcal{HK} implies the minimiser satisfies

$$0 = \frac{d}{dt} (\dot{x}m)(t)$$

$$0 = |\dot{x}(t)|^2 - \frac{1}{2} \frac{d}{dt} \left(\frac{\dot{m}}{m} \right)(t) - \frac{1}{4} \frac{\dot{m}(t)^2}{m(t)^2}$$

with $x(0) = x_0$, $x(1) = x_1$, $m(0) = m_0$ and $m(1) = m_1$.

- ⑥ The Euler-Lagrange equations are satisfied by

$$m(t) = (1-t)^2 m_0 + t^2 m_1 + 2t(1-t)\sqrt{m_0 m_1} \cos \|\|x_0 - x_1\|$$

$$\varphi(t) = \frac{1}{\|x_0 - x_1\|} \cos^{-1} \left(\frac{(1-t)\sqrt{m_0} + t\sqrt{m_1} \cos(\|x_0 - x_1\|)}{\sqrt{m(t)}} \right)$$

$$x(t) = x_0 + (x_1 - x_0)\varphi(t).$$

The Hellinger–Kantorovich Distance Between Diracs III

- 7 The Hellinger–Kantorovich distance between Diracs closer than $\pi/2$ is therefore:

$$d_{HK}^2(\mu, \nu) = m_0 + m_1 - 2\sqrt{m_0 m_1} \cos \|x_0 - x_1\|.$$

- 8 Now if $\|x_0 - x_1\| > \frac{\pi}{2}$ we consider the ‘pure creation’ and ‘pure destruction’ interpolation: $\rho(t, x) = m_0(t)\delta_{x_0} + m_1(t)\delta_{x_1}$ where $m_0(0) = m_0$, $m_0(1) = 0$, $m_1(0) = 0$ and $m_1(1) = m_1$.
- 9 The HK energy is

$$d_{HK}^2(\mu, \nu) = \inf_{m_0, m_1} \frac{1}{4} \int_0^1 \left(\frac{\dot{m}_0(t)}{m_0(t)} \right)^2 m_0(t) + \left(\frac{\dot{m}_1(t)}{m_1(t)} \right)^2 m_1(t) dt.$$

- 10 Optimising implies $m_0(t) = (1-t)^2 m_0$ and $m_1(t) = t^2 m_1$ and

$$d_{HK}^2(\mu, \nu) = m_0 + m_1.$$

The Hard Kantorovich Form of the HK Distance

- ① Define

$$\begin{aligned}\hat{c}(x_0, m_0, x_1, m_1) &:= d_{HK}(m_0 \delta_{x_0}, m_1 \delta_{x_1}) \\ &= m_0 + m_1 - 2\sqrt{m_0 m_1} \overline{\cos}(\|x_0 - x_1\|)\end{aligned}$$

where $\overline{\cos}(t) = \cos(t \wedge \frac{\pi}{2})$.

- ② Let $\pi_0 \in \mathcal{M}_+(\Omega^2)$ be the mass that leaves, i.e. the amount of mass leaving x_0 for x_1 is $\pi_0(x_0, x_1)$.
- ③ Let $\pi_1 \in \mathcal{M}_+(\Omega^2)$ be the mass that arrives, i.e. the amount of mass arriving at x_1 from x_0 is $\pi_1(x_0, x_1)$.
- ④ The amount of mass leaving x_0 is given by μ and the amount of mass arriving at x_1 is given by ν , hence $P_0 \# \pi_0 = \mu$ and $P_1 \# \pi_1 = \nu$.
- ⑤ Hard-marginal Kantorovich form:

$$d_{HK}^2(\mu, \nu) = \inf_{\substack{P_0 \# \pi_0 = \mu \\ P_1 \# \pi_1 = \nu \\ \pi_0, \pi_1 \ll \lambda}} \int_{\Omega^2} \hat{c} \left(x_0, \frac{d\pi_0}{d\lambda}(x_0, x_1), x_1, \frac{d\pi_1}{d\lambda}(x_0, x_1) \right) d\lambda(x_0, x_1).$$

The Soft Kantorovich Form of the HK Distance I

Let $C(x_0, x_1) = \overline{\cos}(\|x_0 - x_1\|)$. We can write

$$\begin{aligned} d_{HK}^2(\mu, \nu) &= \int_{\Omega^2} \frac{d\pi_0}{d\lambda} + \frac{d\pi_1}{d\lambda} - 2\sqrt{\frac{d\pi_0}{d\lambda} \frac{d\pi_1}{d\lambda}} C d\lambda \\ &= \int_{\Omega^2} \frac{d\pi_0}{d\lambda} + \frac{d\pi_1}{d\lambda} \\ &\quad - 2\sqrt{\frac{d\pi_0}{d\lambda} \frac{d\pi_1}{d\lambda}} C \underbrace{\left[1 + \log C - \frac{1}{2} \log \left(\sqrt{\frac{d\pi_1}{d\lambda}} C \right) - \frac{1}{2} \log \left(\sqrt{\frac{d\pi_0}{d\lambda}} C \right) \right]}_{=0} d\lambda \\ &= \int_{\Omega^2} -2\sqrt{\frac{d\pi_0}{d\lambda} \frac{d\pi_1}{d\lambda}} C \log C d\lambda \\ &\quad + \int_{\Omega^2} \frac{d\pi_0}{d\lambda} \left(1 - \sqrt{\frac{d\pi_1}{d\lambda}} C + \sqrt{\frac{d\pi_1}{d\lambda}} C \log \left(\sqrt{\frac{d\pi_1}{d\lambda}} C \right) \right) d\lambda \\ &\quad + \int_{\Omega^2} \frac{d\pi_1}{d\lambda} \left(1 - \sqrt{\frac{d\pi_0}{d\lambda}} C + \sqrt{\frac{d\pi_0}{d\lambda}} C \log \left(\sqrt{\frac{d\pi_0}{d\lambda}} C \right) \right) d\lambda \\ &= \int_{\Omega^2} c d\pi + \int_{\Omega^2} \varphi \left(\sqrt{\frac{d\pi_1}{d\lambda}} C \right) d\pi_0 + \int_{\Omega^2} \varphi \left(\sqrt{\frac{d\pi_0}{d\lambda}} C \right) d\pi_1 \end{aligned}$$

where $c(x_0, x_1) = -2 \log C(x_0, x_1)$, $\pi = \sqrt{\frac{d\pi_0}{d\lambda} \frac{d\pi_1}{d\lambda}} C \lambda$ and $\varphi(s) = s \log s - s + 1$.

The Soft Kantorovich Form of the HK Distance II

① Let

$$\begin{aligned}\mu &= u\bar{\mu} + \bar{\mu}^\perp & u &= \frac{d\mu}{d\bar{\mu}} \\ \nu &= w\bar{\nu} + \bar{\nu}^\perp & w &= \frac{d\nu}{d\bar{\nu}}.\end{aligned}$$

② Then

$$\begin{aligned}\pi_0 &= (u \otimes 1)\pi + (\text{Id} \times \text{Id})_{\#}\bar{\mu} \\ \pi_1 &= (1 \otimes w)\pi + (\text{Id} \times \text{Id})_{\#}\bar{\nu}.\end{aligned}$$

③ We have the densities

$$\frac{d\pi}{d\pi_0} = \frac{1}{u \otimes 1} \quad \frac{d\pi}{d\pi_0} = \frac{\frac{d\pi}{d\lambda}}{\frac{d\pi_0}{d\lambda}} = \frac{\sqrt{\frac{d\pi_0}{d\lambda} \frac{d\pi_1}{d\lambda}} C}{\frac{d\pi_0}{d\lambda}} = \sqrt{\frac{\frac{d\pi_1}{d\lambda}}{\frac{d\pi_0}{d\lambda}}} C.$$

④ Note also $\frac{dP_{0\#}\pi}{d\mu} = \frac{1}{u}$.

The Soft Kantorovich Form of the HK Distance III

⑤ So

$$\begin{aligned} \int_{\Omega^2} \varphi \left(\sqrt{\frac{\frac{d\pi_1}{d\lambda}}{\frac{d\pi_0}{d\lambda}}} c \right) d\pi_0 &= \int_{\Omega} \varphi \left(\frac{1}{u(x_0)} \right) d\pi_0(x_0, x_1) \\ &= \int_{\Omega} \varphi \left(\frac{1}{u(x_0)} \right) d\mu(x_0) \\ &= \int_{\Omega} \varphi \left(\frac{dP_{0\#}\pi}{d\mu} \right) d\mu =: \text{KL}(P_{0\#}\pi | \mu). \end{aligned}$$

⑥ Similarly,

$$\int_{\Omega^2} \varphi \left(\sqrt{\frac{\frac{d\pi_0}{d\lambda}}{\frac{d\pi_1}{d\lambda}}} c \right) d\pi_1 = \text{KL}(P_{1\#}\pi | \nu).$$

⑦ So far, for $c(x_0, x_1) = -2 \log \overline{\cos}(\|x_0 - x_1\|)$

$$d_{HK}^2(\mu, \nu) = \int_{\Omega^2} c d\pi + \text{KL}(P_{0\#}\pi | \mu) + \text{KL}(P_{1\#}\pi | \nu).$$

⑧ In fact,

$$d_{HK}^2(\mu, \nu) = \inf_{\pi \in \mathcal{M}_+(\Omega^2)} \left(\int_{\Omega^2} c d\pi + \text{KL}(P_{0\#}\pi | \mu) + \text{KL}(P_{1\#}\pi | \nu) \right).$$

Hellinger–Kantorovich Geodesics via Optimal Plans

Let $\mu, \nu \in \mathcal{M}_+(\Omega)$, π^* optimal and T^* be the Monge map $\pi^* = (\text{Id} \times T^*)_\# \bar{\mu}$. Let $\bar{\mu} = P_{1\#} \pi^*$, $\bar{\nu} = P_{2\#} \pi^*$ and write

$$\mu = u\bar{\mu} + \mu^\perp \quad \nu = w\bar{\nu} + \nu^\perp.$$

Then a geodesic is given by

$$\begin{aligned}\bar{\rho}(t, \cdot) &= X \left(t; \cdot, u(\cdot), T^*(\cdot), w \circ T^*(\cdot) \right)_\# \left[M \left(t; \cdot, u(\cdot), T^*(\cdot), w \circ T^*(\cdot) \right) \bar{\mu} \right] \\ \rho(t, \cdot) &= \bar{\rho}(t, \cdot) + (1-t)^2 \mu^\perp + t^2 \nu^\perp \\ v(t, \cdot) &= \frac{\partial X}{\partial t} \left(t; \cdot, u(\cdot), T^*(\cdot), w \circ T^*(\cdot) \right) \\ \bar{\zeta}(t, \cdot) &= X \left(t; \cdot, u(\cdot), T^*(\cdot), w \circ T^*(\cdot) \right)_\# \left[\frac{\partial M}{\partial t} \left(t; \cdot, u(\cdot), T^*(\cdot), w \circ T^*(\cdot) \right) \bar{\mu} \right] \\ \zeta(t, \cdot) &= \bar{\zeta}(t, \cdot) - 2(1-t)\mu^\perp + 2t\nu^\perp.\end{aligned}$$

where

$$\begin{aligned}M(t; x_0, m_0, x_1, m_1) &= (1-t)^2 m_0 + t^2 m_1 + 2t(1-t)\sqrt{m_0 m_1} \cos \|x_0 - x_1\| \\ \varphi(t) &= \frac{1}{\|x_0 - x_1\|} \cos^{-1} \left(\frac{(1-t)\sqrt{m_0} + t\sqrt{m_1} \cos(\|x_0 - x_1\|)}{\sqrt{M(t)}} \right)\end{aligned}$$

$$X(t; x_0, m_0, x_1, m_1) = x_0 + (x_1 - x_0)\varphi(t).$$

Time Independent Benamou–Brenier Form

- ① Let $\mu, \nu \in \mathcal{M}_+(\Omega)$ and $\pi^* = (\text{Id} \times T^*)_{\#}\bar{\mu}$ be optimal.
- ② Let (ρ, ω, ζ) be the geodesics constructed on the previous slide.
- ③ Set for $t \in [0, 1]$:

$$\alpha(t, \cdot) = \frac{d\bar{\zeta}_t}{d\rho_t} - 2(1-t)\frac{d\mu^\perp}{d\rho_t}.$$

- ④ Then²³

$$\nu(0, x) = \begin{cases} \frac{T^*(x) - x}{\|T^*(x) - x\|} \sqrt{\frac{w(T^*(x))}{u(x)}} \sin(\|T^*(x) - x\|) & \bar{\mu}\text{-a.e.}, \\ 0 & \mu^\perp\text{-a.e.}, \end{cases}$$

$$\alpha(0, x) = \begin{cases} 2 \left(\sqrt{\frac{w(T^*(x))}{u(x)}} \cos(\|T^*(x) - x\|) - 1 \right) & \bar{\mu}\text{-a.e.}, \\ -2 & \mu^\perp\text{-a.e.} \end{cases}$$

and

$$d_{HK}^2(\mu, \nu) = \int_{\Omega} \left(\|\nu(0, \cdot)\|^2 + \frac{1}{4}(\alpha(0, \cdot))^2 \right) d\mu + \|\nu^\perp\|.$$

²³Cai, Cheng, Schmitzer and T., *The Linearized Hellinger–Kantorovich Distance*, SIAM Journal on Imaging Sciences, 15(1):45–83, 2022.

The Linear Hellinger–Kantorovich Distance

- ① One can show that $\mu \perp \nu^\perp$.
- ② In particular, if $\text{spt}(\mu) = \Omega$ then $\nu^\perp = 0$, and

$$d_{HK}^2(\mu, \nu) = \int_{\Omega} \left(\|\nu(0, \cdot)\|^2 + \frac{1}{4}(\alpha(0, \cdot))^2 \right) d\mu.$$

- ③ ‘Logarithmic’ map:

$$\text{Log}_{HK}(\mu; \mu_i) := (\nu(0, \cdot), \alpha(0, \cdot)).$$

- ④ Hellinger–Kantorovich distance from reference:

$$d_{HK}(\mu, \nu) = \|\text{Log}_{HK}(\mu; \nu)\|_{L^2(\mu)}.$$

- ⑤ Linear Hellinger–Kantorovich distance:²⁴

$$d_{HK,\text{lin},\mu}(\mu_1, \mu_2) = \|\text{Log}_{HK}(\mu; \mu_2) - \text{Log}_{HK}(\mu; \mu_1)\|_{L^2(\mu)}.$$

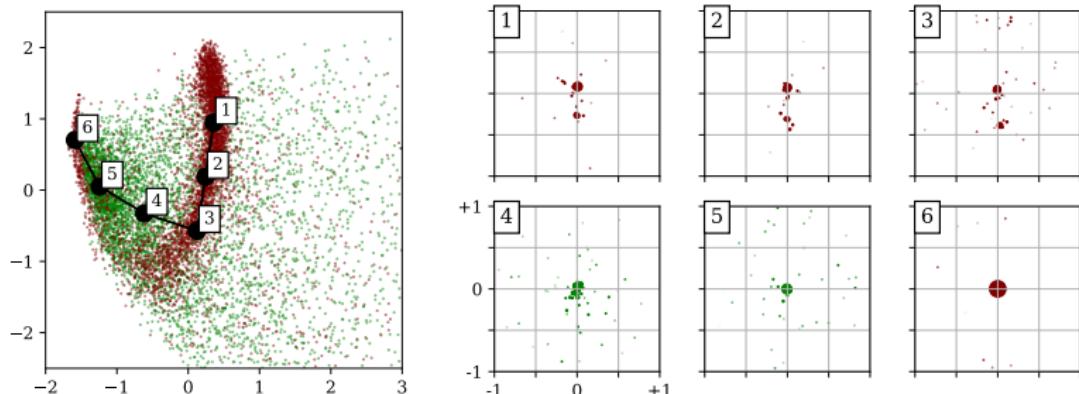
- ⑥ **Linear Hellinger–Kantorovich Assumption:**

$$d_{HK}(\mu_1, \mu_2) \approx d_{HK,\text{lin},\mu}(\mu_1, \mu_2).$$

²⁴Cai, Cheng, Schmitzer and T., *The Linearized Hellinger–Kantorovich Distance*, SIAM Journal on Imaging Sciences, 15(1):45–83, 2022.

Example 6: Particle Decay Classification

Aim: can we classify W boson jets and QCD (quark or gluon) jets?



Source: **Cai, Cheng, Schmitzer and T.**, *The Linearized Hellinger–Kantorovich Distance*, SIAM Journal on Imaging Sciences, 15(1):45–83, 2022.

- 1 Computational Methods
- 2 The Wasserstein Distance
- 3 The Sliced Wasserstein Distance
- 4 The Linear Wasserstein Distance
- 5 The Hellinger–Kantorovich Distance
- 6 The TL^P Distance
 - The TL^P Distance
 - Colour Transfer
 - The Linear TL^P Distance
 - Time Series Classification

The TL^P Distance

- ➊ The Wasserstein distance treats signals/images as probability measures. This can be restrictive. Can we generalise whilst keeping the nice properties of optimal transport?
- ➋ We want to treat the signal/image as a function rather than a measure.
- ➌ The idea is to treat signals as a pair (f, μ) where $f \in L^P(\mu)$.
- ➍ Mostly we consider when μ is the uniform measure (either continuous or discrete), but one could also trivially adapt in order to weight features of the signal, for example.
- ➎ Note that we can compare signals on different domains.
- ➏ TL^P definition (Monge formulation):²⁵ (a.k.a. graph space optimal transport²⁶)

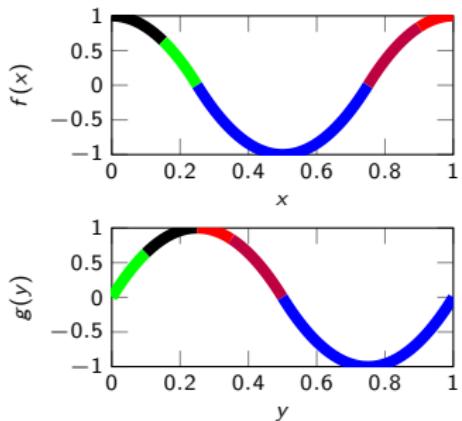
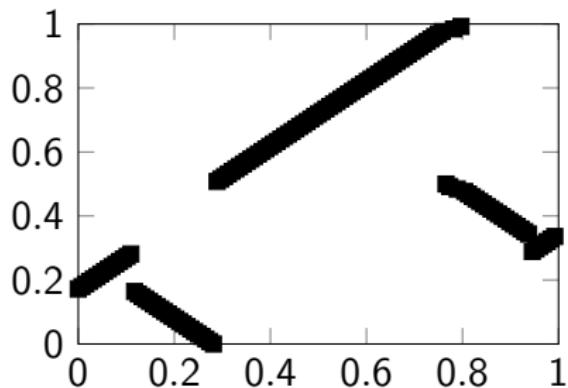
$$d_{\text{TL}^P}^P((f, \mu), (g, \nu)) = \inf_{T: T_\# \mu = \nu} \int_X |x - T(x)|^P + \lambda |f(x) - g(T(x))|^P d\mu(x).$$

²⁵ García Trillos and Slepčev, *Continuum Limit of Total Variation on Point Clouds*, ARMA, 220:193–241, 2016.

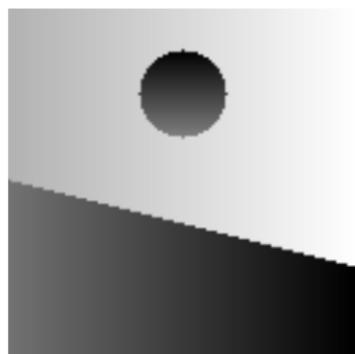
²⁶ Métivier, Brossier, Mérigot and Oudet, *A Graph Space Optimal Transport Distance as a Generalization of L^P Distances: Application to a Seismic Imaging Inverse Problem*, Inverse Problems, 35(8), 2019.

A Simple Example

For example consider the functions $f(x) = \cos(2\pi x)$ and $g(y) = \sin(2\pi y)$ defined on $[0, 1]$ with the uniform measure. The optimal plan using the TL^2 distance is given below.



Example 7: Synthetic Colour Transfer



(a) Exemplar image.



(b) Original image to be shaded.



(c) The optimal transport solution.

Source: T., Park, Kolouri, Rohde and Slepčev, *A Transportation L_p Distance for Signal Analysis*, Journal of Mathematical Imaging and Vision, 59(2):187–210, 2017.

Example 8: Real World Colour Transfer



(a) Exemplar image.



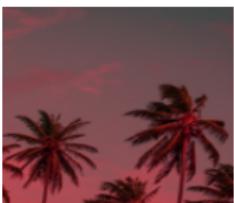
(b) Original image to be coloured.



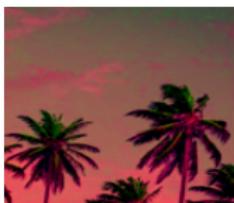
(c) The optimal transport solution.



(d) The TL^2 solution.



(e) Reinhard, Ashikhmin, Gooch and Shirley's method.



(f) Pitié and Kokaram's method.

Source: T., Park, Kolouri, Rohde and Slepčev, *A Transportation Lp Distance for Signal Analysis*, Journal of Mathematical Imaging and Vision, 59(2):187–210, 2017.

The Linear TL^2 Distance

- ① Fix a reference point $(f_0, \mu_0) \in \text{TL}^2$ and let T_i be the TL^2 -optimal transport map between (f_0, μ_0) and (f_i, μ_i) .
- ② I.e. $[T_i]_{\#}\mu_0 = \mu_i$ and

$$d_{\text{TL}^2}^2((f_i, \mu_i), (f_0, \mu_0)) = \int_X |x - T_i(x)|^2 + \lambda |f_0(x) - f_i(T_i(x))|^2 d\mu_0(x).$$

- ③ Assume μ_0 has a density ρ_0 then we define

$$P(f_i, \mu_i) = (P_1(f_i, \mu_i), P_2(f_i, \mu_i)) \in L^2(\mu_0; \mathbb{R}^d) \times L^2(\mu_0; \mathbb{R})$$

$$[P_1(f_i, \mu_i)](x) = (T_i(x) - x)\sqrt{\rho_0(x)}$$

$$[P_2(f_i, \mu_i)](x) = (f_i(T_i(x)) - f_0(x))\sqrt{\lambda\rho_0(x)}.$$

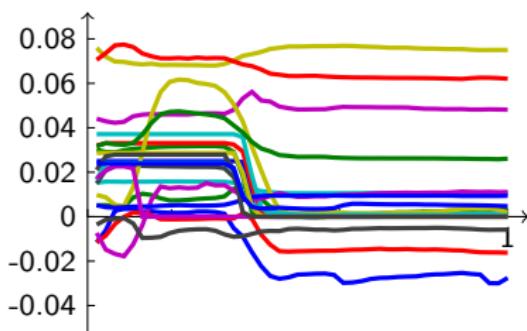
- ④ Then define

$$d_{\text{TL}^2, \text{lin}, (f_0, \mu_0)}((f_i, \mu_i), (f_j, \mu_j)) = \|P(f_i, \mu_i) - P(f_j, \mu_j)\|_{L^2}.$$

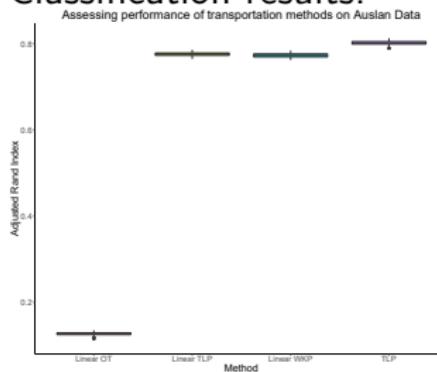
Example 9: AUSLAN Classification

Data: The AUSLAN data set is a set of 95 words ‘spoken’ by a native AUSLAN (Australian sign language) using 22 sensors on a cyberglove. There are 27 signals in each class, so a total of 2565 signals.

Example data:



Classification results:



Source: Crook, Cucuringu, Hurst, Schönlieb, T. and Zygalkis, A Linear Transportation L_p Distance for Pattern Recognition, Pattern Recognition 147, 2024.

- + Optimal transport distances are often a good model for data and using them in your method can often lead to improved performance.
- + There is a rich theory which includes geodesics and a formal Riemannian structure.
- + Restrictions, such as the data needs to be normalised, are not prohibitive as there are optimal transport distances enjoying many of the same theoretical properties as Wasserstein distances in the unbalanced setting.
- Computation is improving all the time, but is still often much more expensive than Euclidean distances.

References I

- ① **Angenent, Haker and Tannenbaum**, *Minimizing flows for the Monge–Kantorovich problem*, SIAM journal on mathematical analysis, 35(1):61–97, 2003.
- ② **Arjovsky, Chintala and Bottou**, *Wasserstein Generative Adversarial Networks*, International Conference on Machine Learning, 2017.
- ③ **Benamou and Brenier**, *A Computational Fluid Mechanics Solution to the Monge–Kantorovich Mass Transfer Problem*, Numerische Mathematik, 84(3):375–393, 2000.
- ④ **Benamou, Froese and Oberman**, *Numerical Solution of the Optimal Transportation Problem Using the Monge–Ampere Equation*, Journal of Computational Physics, 260:107–126, 2014.
- ⑤ **Cai, Cheng, Schmitzer and T.**, *The Linearized Hellinger–Kantorovich Distance*, SIAM Journal on Imaging Sciences, 15(1):45–83, 2022.
- ⑥ **Chartrand, Vixie, Wohlberg and Bollt**, *A Gradient Descent Solution to the Monge–Kantorovich Problem*, Applied Mathematical Sciences, 3(22):1071–1080, 2009.
- ⑦ **Chizat, Peyré, Schmitzer and Vialard**, *An Interpolating Distance Between Optimal Transport and Fisher–Rao Metrics*, Foundations of Computational Mathematics 18:1–44, 2018.

References II

- ⑧ **Chizat, Peyré, Schmitzer and Vialard**, *Unbalanced Optimal Transport: Dynamic and Kantorovich formulations*, Journal of Functional Analysis, 274:3090–3123, 2018.
- ⑨ **Crook, Cucuringu, Hurst, Schönlieb, T. and Zygalakis**, *A Linear Transportation L_p Distance for Pattern Recognition*, Pattern Recognition 147, 2024.
- ⑩ **Cuturi**, *Sinkhorn Distances: Lightspeed Computation of Optimal Transport*, In Advances in Neural Information Processing Systems, pp. 2292–2300, 2013.
- ⑪ **Engquist and Yang**, *Seismic Imaging and Optimal Transport*, Communications in Information and Systems, 19(2):95–145, 2019.
García Trillo and Slepčev, *Continuum Limit of Total Variation on Point Clouds*, ARMA, 220:193–241, 2016.
- ⑫ **Gigli**, *On Hölder Continuity-in-Time of the Optimal Transport Map Towards Measures Along a Curve*, Proceedings of the Edinburgh Mathematical Society, 54(2):401–409, 2011.
- ⑬ **Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville and Bengio**, *Generative Adversarial Networks*, arxiv:1406.2661, 2014.
- ⑭ **Hamm, Moosmueller, Schmitzer and T.**, *Manifold Learning in Wasserstein Space*, arxiv:2311.08549, 2023.
- ⑮ **Jacobs and Léger**, *A Fast Approach to Optimal Transport; The Back-And-Forth Method*, Numerische Mathematik, 146(3):513–544, 2020.

References III

- ⑯ **Khurana, Kannan, Cloninger and Moosmüller**, *Supervised Learning of Sheared Distributions Using Linearized Optimal Transport*, Sampling Theory, Signal Processing, and Data Analysis, 21(1), 2023.
- ⑰ **Kolouri, Park, T., Slepčev and Rohde**, *Optimal Mass Transport: Signal Processing and Machine Learning Applications*, IEEE Signal Processing Magazine, 34(4):43–59, 2017.
- ⑱ **Kolouri, Pope, Martin and Rohde**, *Sliced Wasserstein auto-encoders*, ICLR, 2018.
- ⑲ **Kolouri and Rohde**, *Optimal transport a crash course*, IEEE ICIP 2016 Tutorial Slides: Part 1, 2016.
- ⑳ **Kondratyev, Monsaingeon and Vorotnikov**, *A New Optimal Transport Distance On The Space Of Finite Radon Measures*, Advances in Differential Equations, 21:1117–1164, 2016.
- ㉑ **Levy**, *A Numerical Algorithm for L_2 Semi-Discrete Optimal Transport in 3D*, ESAIM Math. Model. Numer. Anal., 49(6):1693–1715, 2015.
- ㉒ **Liero, Mielke and Savaré**, *Optimal Entropy-Transport Problems and a New Hellinger–Kantorovich Distance Between Positive Measures*, Inventiones Mathematicae 211:969–1117, 2018.
- ㉓ **Métivier, Brossier, Mérigot and Oudet**, *A Graph Space Optimal Transport Distance as a Generalization of L^p Distances: Application to a Seismic Imaging Inverse Problem*, Inverse Problems, 35(8), 2019.

- 24 Nowozin, Cseke and Tomioka *f-GAN: Training generative neural samplers using variational divergence minimization*, NeuRIPS, 2016.
- 25 Park and T., *Representing and Learning High Dimensional Data with the Optimal Transport Map from a Probabilistic Viewpoint*, CVPR, 2018.
- 26 Stanczuk, Etmann, Kreusser and Schönlieb, *Wasserstein GANs Work Because They Fail (To Approximate the Wasserstein distance)*, arxiv:2103.01678, 2021.
- 27 T., Park, Kolouri, Rohde and Slepčev, *A Transportation L_p Distance for Signal Analysis*, Journal of Mathematical Imaging and Vision, 59(2):187–210, 2017.
- 28 Tolstikhin, Bousquet, Gelly and Schölkopf, *Wasserstein Auto-Encoders*, International Conference on Learning Representations, 2018.
- 29 Wang, Slepčev, Basu, Ozolek and Rohde, *A Linear Optimal Transportation Framework for Quantifying and Visualizing Variations in Sets of Images*, International Journal of Computer Vision 101(2):254–269, 2013.

Thank you for listening!

In theory, there is no difference between theory and practice. But in practice, there is.

— Yogi Berra