Learning with Regularized Distances: **Optimal Transport** and **Dynamic Time Warping**

Marco Cuturi



Joint work with many people, including: G. Peyré, A. Genevay (ENS), A. Doucet (Oxford), J. Solomon (MIT), J.D. Benamou, N. Bonneel, F. Bach, L. Nenna (INRIA), G. Carlier (Dauphine), M. Blondel (NTT).

What is Optimal Transport?



What is Optimal Transport?



What is Optimal Transport?













OT and data-analysis

- Key developments in (applied) maths ~'90s [McCann'95], [JKO'98], [Benamou'98], [Gangbo'98], [Ambrosio'06], [Villani'03/'09].
- Key developments in TCS / graphics since '00s [Rubner'98], [Indyk'03], [Naor'07], [Andoni'15].

Small to *no-impact* in large-scale data analysis:
 + computationally heavy;
 + Wasserstein distance is not differentiable

OT and data-analysis

Today's talk: Entropy Regularized OT

- **Very fast** compared to usual approaches, <u>GPGPU parallel</u>.
- **Differentiable**, important if we want to use OT distances as **loss functions**.
- Can be **automatically differentiated**, simple iterative process, *DL*-toolboxes compatible.
- OT can become a building block in ML.

Wasserstein distance is not differentiable

 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T: \Omega \to \Omega$ $\inf_{\mathbf{T}\neq\boldsymbol{\mu}=\boldsymbol{\nu}}\int_{\Omega}\boldsymbol{c}(x,\mathbf{T}(x))\boldsymbol{\mu}(dx)$ $\forall B \subset \Omega, \mathbf{T} \# \boldsymbol{\mu}(B) = \boldsymbol{\nu}(B)$ 8

 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T : \Omega \to \Omega$ [Brenier'87] If $\Omega = \mathbb{R}^d, c = \| \cdot - \cdot \|^2$, μ, ν a.c., then $T = \nabla u, u$ convex.



 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T : \Omega \to \Omega$ $\inf_{T \# \mu = \nu} \int_{\Omega} c(x, T(x)) \mu(dx)$



 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T: \Omega \to \Omega$



[Kantorovich'42] Relaxation

• Instead of maps $T : \Omega \to \Omega$, consider probabilistic maps, i.e. couplings $P \in \mathcal{P}(\Omega \times \Omega)$:

$$\Pi(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega, \\ \boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \\ \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B}) \}$$

[Kantorovich'42] Relaxation

 $\Pi(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega,$ $\boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B})\}$



[Kantorovich'42] Relaxation

 $\Pi(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega,$ $\boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B})\}$



Wasserstein Distance



Wasserstein between 2 Diracs



Wasserstein on Uniform Measures



Wasserstein on Uniform Measures



Optimal Assignment C Wasserstein









Consider
$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$.
 $M_{XY} \stackrel{\text{def}}{=} [D(x_i, y_j)^p]_{ij}$
 $U(a, b) \stackrel{\text{def}}{=} \{ P \in \mathbb{R}^{n \times m} | P \mathbf{1}_m = a, P^T \mathbf{1}_n = b \}$
 $\stackrel{\mathbf{y}_1 \qquad \cdots \qquad \mathbf{y}_m}{:} \stackrel{\mathbf{b}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf$

Consider
$$\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\boldsymbol{\nu} = \sum_{j=1}^{m} b_j \delta_{y_j}$.
 $M_{\boldsymbol{X}\boldsymbol{Y}} \stackrel{\text{def}}{=} [D(\boldsymbol{x}_i, \boldsymbol{y}_j)^p]_{ij}$
 $U(\boldsymbol{a}, \boldsymbol{b}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathbb{R}^{n \times m}_+ | \boldsymbol{P} \boldsymbol{1}_m = \boldsymbol{a}, \boldsymbol{P}^T \boldsymbol{1}_n = \boldsymbol{b} \}$

Def. Optimal Transport Problem $W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$









Note: flow/PDE formulations [**Beckman'61**]/[**Benamou'98**] can be used for *p*=1/*p*=2 for a sparse-graph metric/Euclidean metric.










Discrete OT Problem



Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$

$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij}(\log P_{ij})$$

Note: Unique optimal solution because of strong concavity of Entropy

Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



Note: Unique optimal solution because of strong concavity of Entropy

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle P, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} \log P_{ij} + \alpha^T (P\mathbf{1} - \mathbf{a}) + \beta^T (P^T \mathbf{1} - \mathbf{b})$$

$$\partial L/\partial P_{ij} = M_{ij} + \gamma (\log P_{ij} + 1) + \alpha_i + \beta_j$$

$$(\partial L/\partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma} + \frac{1}{2}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma} + \frac{1}{2}} = u_i K_{ij} v_j$$

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

- [Sinkhorn'64] fixed-point iterations for (*u*, *v*) *u* ← *a*/*Kv*, *v* ← *b*/*K^Tu O(nm)* complexity, GPGPU parallel [C'13].
- $O(n^{d+1})$ if $\Omega = \{1, \ldots, n\}^d$ and D^p separable. [S..C..'15]

Very Fast EMD Approx. Solver



Note. (Ω, D) is a random graph with shortest path metric, histograms sampled uniformly on simplex, Sinkhorn tolerance 10⁻².

Regularization ----> Differentiability

 $W_{\gamma}((\boldsymbol{a},\boldsymbol{X}),(\boldsymbol{b},\boldsymbol{Y})) = \min_{\boldsymbol{P}\in U(\boldsymbol{a},\boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



Regularization ----> Differentiability

 $W_{\gamma}((\boldsymbol{a} + \boldsymbol{\Delta}\boldsymbol{a}, \boldsymbol{X}), (\boldsymbol{b}, \boldsymbol{Y})) = W_{\gamma}((\boldsymbol{a}, \boldsymbol{X}), (\boldsymbol{b}, \boldsymbol{Y})) + ??$



Regularization ----> Differentiability

 $W_{\gamma}((a + \Delta a, X), (b, Y)) = W_{\gamma}((a, X), (b, Y)) + ??$



Regularization ---> Differentiability

 $W_{\gamma}((a, X + \Delta X), (b, Y)) = W_{\gamma}((a, X), (b, Y)) + ??$



Regularization ---> Differentiability

 $W_{\gamma}((a, X + \Delta X), (b, Y)) = W_{\gamma}((a, X), (b, Y)) + ??$



Crucial for "min data + W" problems

- Quantization, k-means problem [Lloyd'82] $\min_{\substack{\mu \in \mathcal{P}(\mathbb{R}^d) \\ |\operatorname{supp} \mu| = k}} W_2^2(\mu, \nu_{data})$
- [McCann'95] Interpolant

$$\min_{\boldsymbol{\mu}\in\mathcal{P}(\Omega)}(1-t)W_2^2(\boldsymbol{\mu},\boldsymbol{\nu_1})+tW_2^2(\boldsymbol{\mu},\boldsymbol{\nu_2})$$

• [JKO'98] PDE's as gradient flows in $(\mathcal{P}(\Omega), W)$.

$$\mu_{t+1} = \operatorname*{argmin}_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} J(\boldsymbol{\mu}) + \lambda_t W_p^p(\boldsymbol{\mu}, \mu_t)$$

Crucial for "min data + W" problems

• Quantization, k-means problem [Lloyd'82] min $W_2^2(\mu, \nu_{data})$ $\mu \in \mathcal{P}(\mathbb{R}^d)$

Any (ML) problem involving a KL or L2 loss between (parameterized) histograms or probabilility measures can be easily
Wasserstein-ized if we can differentiate W efficiently.

$\mu_{t+1} = \operatorname*{argmin}_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} J(\boldsymbol{\mu}) + \lambda_t W_p^p(\boldsymbol{\mu}, \boldsymbol{\mu}_t)$ $\mu \in \mathcal{P}(\Omega)$

JKO'98 PDE's as gradient flows in $(\mathcal{P}(\Sigma), W)$.

1. Differentiability of Regularized OT

Def. Dual regularized OT Problem $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\alpha, \beta} \alpha^{T} \boldsymbol{a} + \beta^{T} \boldsymbol{b} - \frac{1}{\gamma} (e^{\alpha/\gamma})^{T} \boldsymbol{K} e^{\beta/\gamma}$ **Prop.** $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu})$ is [CD'14]

1. convex w.r.t. *a* (Danskin), $\nabla_{\boldsymbol{a}} W_{\gamma} = \alpha^{\star} = \gamma \log(\boldsymbol{u}).$

2. decreased, when $p = 2, \Omega = \mathbb{R}^d$, using $\mathbf{X} \leftarrow \mathbf{Y} P_{\gamma}^T \mathbf{D}(\mathbf{a}^{-1}).$

2. Duality for Regularized OT's

Prop. Writing $H_{\boldsymbol{\nu}} : \boldsymbol{a} \mapsto W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}),$ [CP'16]

1. H_{ν} has simple Legendre transform:

$$H^*_{\boldsymbol{\nu}}: \boldsymbol{g} \in \mathbb{R}^n \mapsto \gamma \left(E(\boldsymbol{b}) + \boldsymbol{b}^T \log(\boldsymbol{K} e^{\boldsymbol{g}/\gamma}) \right)$$

2. If $A \in \mathbb{R}^{n \times d}$, f convex on \mathbb{R}^d ,

 $\min_{\boldsymbol{a}\in\Sigma_n} H_{\boldsymbol{\nu}}(\boldsymbol{a}) + f(A\boldsymbol{a}) = \max_{\boldsymbol{g}\in\mathbb{R}^d} - H_{\boldsymbol{\nu}}^*(A^T\boldsymbol{g}) - f^*(-\boldsymbol{g})$

3. Stochastic Formulation $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\alpha, \beta} \alpha^{T} \boldsymbol{a} + \beta^{T} \boldsymbol{b} - \frac{1}{\gamma} (e^{\alpha/\gamma})^{T} K e^{\beta/\gamma}$ $= \max \boldsymbol{\alpha}^T \boldsymbol{a} - \gamma (\log \boldsymbol{K} e^{\boldsymbol{\alpha}/\gamma})^T \boldsymbol{b}$ $= \max_{\boldsymbol{\alpha}} \sum_{\boldsymbol{\beta}} \boldsymbol{b}_{\boldsymbol{j}} \left(\boldsymbol{\alpha}^{T} \boldsymbol{a} - \gamma \log \left[\frac{K_{j}^{T} e^{\boldsymbol{\alpha}/\gamma}}{I_{j}} \right] \right)$ i=1m $= \max \sum f_j(\alpha)$

• **[GCPB'16]** shows that incremental gradient methods are competitive with Sinkhorn.

Def. For $L \geq 1$, define $W_L(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P}_L, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle,$ where $P_L \stackrel{\text{def}}{=} \operatorname{diag}(\boldsymbol{u}_L) K \operatorname{diag}(\boldsymbol{v}_L),$ $\boldsymbol{v_0} = \boldsymbol{1}_m; l \ge 0, \boldsymbol{u_l} \stackrel{\text{def}}{=} \boldsymbol{a}/K\boldsymbol{v_l}, \boldsymbol{v_{l+1}} \stackrel{\text{def}}{=} \boldsymbol{b}/K^T\boldsymbol{u_l}.$ **Prop.** $\frac{\partial W_L}{\partial X}, \frac{\partial W_L}{\partial a}$ can be computed recursively, in O(L) kernel $K \times$ vector products.

Algorithmic Formulation of Reg. OT



Algorithmic Formulation of Reg. OT

Example: Differentiability w.r.t.
$$a$$

 $N = K \circ M_{XY}$
 $\nabla_a W_L(\mu, \nu) = \left(\frac{\partial u_L}{\partial a}\right)^T N v_L + \left(\frac{\partial v_L}{\partial a}\right)^T N^T u_L$

Thanks to these tricks...

- [Agueh'11] Barycenters [CD'14][BCCNP'15]
 [GCP'15][S..C..'15]
- [Burger'12] TV gradient flow using duality [CP'16]
- Dictionary Learning / Latent Factors [RCP'16]
- [Bigot'15] W-PCA [SC'15]
- Inverse problems / Wasserstein regression [BPC'16]
- Density fitting / parameter estimation [MMC'16]

Wasserstein Barycenters

N $\min_{\boldsymbol{\mu}\in\mathcal{P}(\Omega)}\sum_{i=1}^{\infty}\lambda_i W_p^p(\boldsymbol{\mu},\boldsymbol{\nu_i})$ ${\cal V}_1$ Wasserstein $\mathcal{P}(\Omega)$ Barycenter [Agueh'11] ν_2 ν_3

Multimarginal Formulation

• Exact solution (W_2) using MM-OT. [Agueh'11]



Multimarginal Formulation

• Exact solution (W_2) using MM-OT. [Agueh'11]



If $|\operatorname{supp} \boldsymbol{\nu_i}| = \boldsymbol{n_i}$, LP of size $(\prod_i \boldsymbol{n_i}, \sum_i \boldsymbol{n_i})$

Finite Case, LP Formulation

• When Ω is a finite set, metric *M*, another LP.



Finite Case, LP Formulation

• When Ω is a finite set, metric *M*, another LP.



If
$$|\Omega| = n$$
, LP of size $(Nn^2, (2N - 1)n)$; unstable

Primal Descent on Regularized W



Fast Computation of Wasserstein Barycenters International Conference on Machine Learning 2014



Primal Descent on Regularized W



Fast Computation of Wasserstein Barycenters International Conference on Machine Learning 2014



Primal Descent on Regularized W



Fast Computation of Wasserstein Barycenters International Conference on Machine Learning 2014



Primal Descent on Algorithmic W



Primal Descent on Algorithmic W





Wasserstein Barycenter = KL Projections

$$\langle P, M_{XY} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P \mid \mathbf{K})$$
$$\min_{\mathbf{a}} \sum_{i=1}^{N} \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} \in [\mathbf{P}_1, \dots, \mathbf{P}_N] \\ \mathbf{P} \in \mathbf{C}_1 \cap \mathbf{C}_2}} \sum_{i=1}^{N} \lambda_i \mathbf{KL}(\mathbf{P}_i \mid \mathbf{K})$$
$$\mathbf{C_1} = \{\mathbf{P} \mid \exists \mathbf{a}, \forall i, P_i \mathbf{1}_m = \mathbf{a}\}$$
$$\mathbf{C_2} = \{\mathbf{P} \mid \forall i, P_i^T \mathbf{1}_n = \mathbf{b}_i\}$$

Wasserstein Barycenter = KL Projections

$$\begin{split} \min_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_{i} W_{\gamma}(\boldsymbol{a}, \boldsymbol{b_{i}}) &= \min_{\substack{\mathbf{P} = [\boldsymbol{P_{1}}, \dots, \boldsymbol{P_{N}}]\\ \mathbf{P} \in \boldsymbol{C_{1}} \cap \boldsymbol{C_{2}}}} \sum_{i=1}^{N} \lambda_{i} \mathbf{KL}(\boldsymbol{P_{i}} | \boldsymbol{K}) \\ \boldsymbol{C_{1}} &= \{\mathbf{P} | \exists \boldsymbol{a}, \forall i, P_{i} \mathbf{1}_{m} = \boldsymbol{a} \} \\ \boldsymbol{C_{2}} &= \{\mathbf{P} | \forall i, P_{i}^{T} \mathbf{1}_{n} = \boldsymbol{b_{i}} \} \end{split}$$



Wasserstein Barycenter = KL Projections

$$\min_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_{i} W_{\gamma}(\boldsymbol{a}, \boldsymbol{b}_{i}) = \min_{\substack{\mathbf{P} = [P_{1}, \dots, P_{N}] \\ \mathbf{P} \in \boldsymbol{C}_{1} \cap \boldsymbol{C}_{2}}} \sum_{i=1}^{N} \lambda_{i} \mathbf{KL}(\boldsymbol{P}_{i} | \boldsymbol{K})$$
$$\boldsymbol{C}_{1} = \{\mathbf{P} | \exists \boldsymbol{a}, \forall i, P_{i} \mathbf{1}_{m} = \boldsymbol{a} \}$$
$$\boldsymbol{C}_{2} = \{\mathbf{P} | \forall i, P_{i}^{T} \mathbf{1}_{n} = \boldsymbol{b}_{i} \}$$

u=ones(size(B)); % d x N matrix **BCCNP'15** while not converged v=u.*(K'*(B./(K*u))); % 2(Nd^2) cost u=bsxfun(@times,u,exp(log(v)*weights))./v; end Iterative Bregman Projections for Regularized Transportation Problems a=mean(v,2);SIAM J. on Sci. Comp. 2015

Application: Graphics



Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15 [S..C.'15]

Application: Graphics



Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15 [S..C.'15]
Application: Graphics

Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15 [S..C.'15]

Application: Graphics



Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15 [S..C.'15]

Application: Graphics



Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15 [S..C.'15]

Inverse Wasserstein Problems

• consider Barycenter operator:

$$\boldsymbol{b}(\lambda) \stackrel{\text{def}}{=} \operatorname{argmin}_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_i W_{\gamma}(\boldsymbol{a}, \boldsymbol{b}_i)$$

• address now Wasserstein inverse problems:

Given \boldsymbol{a} , find $\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}(\lambda))$

The Wasserstein Simplex



Barycenters = Fixed Points

Prop. [BCCNP'15] Consider
$$\boldsymbol{B} \in \Sigma_d^N$$

and let $\boldsymbol{U_0} = \boldsymbol{1_{d \times N}}$, and then for $l \ge 0$:
 $\boldsymbol{b}^{l \text{ def}} \exp\left(\log\left(K^T \boldsymbol{U_l}\right)\lambda\right); \begin{cases} \boldsymbol{V_{l+1}} \stackrel{\text{def}}{=} \frac{\boldsymbol{b}^{l} \boldsymbol{1}_N^T}{K^T \boldsymbol{U_l}}, \\ \boldsymbol{U_{l+1}} \stackrel{\text{def}}{=} \frac{\boldsymbol{B}}{K \boldsymbol{V_{l+1}}}. \end{cases}$

Using Truncated Barycenters

- instead of using the exact barycenter $\operatorname{argmin} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}(\lambda))$ $\lambda \in \Sigma_N$
- use instead the L-iterate barycenter

$$\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}^{(L)}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}^{(L)}(\lambda))$$

• Differente using the chain rule.

$$\nabla \mathcal{E}^{(L)}(\lambda) = [\partial \boldsymbol{b}^{(L)}]^T(\boldsymbol{g}), \ \boldsymbol{g} \stackrel{\text{def}}{=} \nabla \text{Loss}(\boldsymbol{a}, \cdot)|_{\boldsymbol{b}^{(L)}(\lambda)}.$$

Gradient / Barycenter Computation

function SINKHORN-DIFFERENTIATE(
$$(p_s)_{s=1}^S, q, \lambda$$
)
 $\forall s, b_s^{(0)} \leftarrow 1$
 $(w, r) \leftarrow (0^S, 0^{S \times N})$
for $\ell = 1, 2, ..., L$ // Sinkhorn loop
 $\forall s, \varphi_s^{(\ell)} \leftarrow K^\top \frac{p_s}{Kb_s^{(\ell-1)}}$
 $p \leftarrow \prod_s (\varphi_s^{(\ell)})^{\lambda_s}$
 $\forall s, b_s^{(\ell)} \leftarrow \frac{p}{\varphi_s^{(\ell)}}$
 $g \leftarrow \nabla \mathcal{L}(p, q) \odot p$
for $\ell = L, L - 1, ..., 1$ // Reverse loop
 $\forall s, w_s \leftarrow w_s + \langle \log \varphi_s^{(\ell)}, g \rangle$
 $\forall s, r_s \leftarrow -K^\top (K(\frac{\lambda_s g - r_s}{\varphi_s^{(\ell)}}) \odot \frac{p_s}{(Kb_s^{(\ell-1)})^2}) \odot b_s^{(\ell-1)}$
 $g \leftarrow \sum_s r_s$
return $P^{(L)}(\lambda) \leftarrow p, \nabla \mathcal{E}_L(\lambda) \leftarrow w$

Application: Volume Reconstruction









 $\lambda_0 = 0.03$

 $\lambda_1 = 0.12$



 $\lambda_2 = 0.40$



 $\lambda_{3} = 0.43$





Wasserstein Barycentric Coordinates: Histogram Regression using Optimal Transport, **SIGGRAPH'16**

[**BPC'16**]

Application: Brain Mapping



Minimum Kantorovich Estimation

$$\theta^{\star} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} W_{p}^{p}(\boldsymbol{p_{\theta}}, \boldsymbol{\nu_{data}}) \quad [\textbf{Bassetti'06}]$$

$$W_{\gamma}(p_{\theta}, \boldsymbol{\nu}_{data}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \langle \boldsymbol{\alpha}, p_{\theta} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu}_{data} \rangle - \gamma \langle e^{\boldsymbol{\alpha}/\gamma}, \ K e^{\boldsymbol{\beta}/\gamma} \rangle$$

$$\nabla_{\theta} W_{\gamma} = \left(\frac{\partial p_{\theta}}{\partial \theta}\right)^T \boldsymbol{\alpha}^{\star}$$

- Application to parameter estimation in discrete models [MMC'16].
- Stochastic methods for semi-discrete OT
 [GCPB'16] 53

To conclude on Wasserstein

- *Entropy* regularization is a very effective way to get OT to work as a generic loss.
- Many recent extensions:
 - [Schmitzer'16]: fast multiscale approaches
 - [ZFMAP'15] [CSPV'16]: Unbalanced transport
 - **[SPKS'16] [PCS'16]** extensions to *Gromov-W*.
 - [FCTR'15] Domain adaptation in ML

Dynamic Time Warping

A distance to compare time series of observations supported on a metric space.



Alignment Grid



Fill in Metric Information



Fill in Metric Information

	y_1	y_2	y_3	y_4	y_5	y_6	$\overset{\cdot}{y_7}$
x_1	D ₁₁	<i>D</i> ₁₂	D ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
x_2	D ₂₁	D ₂₂	D ₂₃	D ₂₄	D ₂₅	D ₂₆	D ₂₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D ₃₅	D ₃₆	D ₃₇
x_4	D_{41}	D ₄₂	D ₄₃	D ₄₄	D_{45}	D ₄₆	D ₄₇
x_5	D ₅₁	D ₅₂	D ₅₃	D ₅₄	D_{55}	D ₅₆	D ₅₇

Alignment Paths

start from (1,1) and ends at (5,7)

x_5	D ₅₁	D ₅₂	D ₅₃	D ₅₄	D_{55}	D ₅₆	D_{57}
x_4	D ₄₁	D_{42}	D ₄₃	D ₄₄	D_{45}	D ₄₆	D ₄₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D ₃₅	D ₃₆	D ₃₇
x_2	D ₂₁	D ₂₂	D ₂₃	D ₂₄	D_{25}	D ₂₆	D ₂₇
x_1	D_{11}	<i>D</i> ₁₂	D ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
	y_1	y_2	y_3	y_4	y_5	y_6	y_7

Three Possible Directions



Example Path

	y_1	y_2	y_3	y_4	y_5	y_6	y_7
x_1	D ₁₁	<i>D</i> ₁₂	D ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
x_2	D ₂₁	D ₂₂	D ₂₃	D ₂₄	D_{25}	D ₂₆	D ₂₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D ₃₅	D ₃₆	D ₃₇
x_4	D_{41}	D_{42}	D_{43}	D ₄₄	D_{45}	D ₄₆	D ₄₇
x_5	D_{51}	D_{52}	D ₅₃	D ₅₄	D_{55}	D ₅₆	D ₅₇

 $C = D_{11} + D_{21}.$

	y_1	y_2	y_3	y_4	y_5	y_6	y_7
x_1	D ₁₁	<i>D</i> ₁₂	D ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
x_2	<i>D</i> ₂₁	D ₂₂	D ₂₃	D ₂₄	D ₂₅	D ₂₆	D ₂₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D ₃₅	D ₃₆	D ₃₇
x_4	D ₄₁	D_{42}	D ₄₃	D ₄₄	D_{45}	D ₄₆	D ₄₇
x_5	D ₅₁	D ₅₂	D ₅₃	D ₅₄	D ₅₅	D ₅₆	D ₅₇





 $C = D_{11} + D_{21} + D_{32} + D_{33} + D_{34} + D_{35} + D_{45} + D_{46} + D_{57}.$

x_5	D ₅₁	D ₅₂	D ₅₃	D ₅₄	D ₅₅	D ₅₆	D 57
x_4	D ₄₁	D_{42}	D ₄₃	D ₄₄	D ₄₅	D 46	D ₄₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D 35	D ₃₆	D ₃₇
x_2	<i>D</i> ₂₁	D ₂₂	D ₂₃	D ₂₄	D ₂₅	D ₂₆	D ₂₇
x_1	D ₁₁	<i>D</i> ₁₂	<i>D</i> ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
	y_1	y_2	y_3	y_4	y_5	y_6	y_7

#All Paths = Delannoy(5,7)

Delannoy(5,7) = 2,241 ; *Delannoy*(20,20)= 4.53e13



Dynamic Time Warping (Distance)

$$d_{\text{DTW}}(\boldsymbol{X}, \boldsymbol{Y}) = \min_{\pi \in \mathcal{A}(\boldsymbol{X}, \boldsymbol{Y})} \sum_{i=1}^{N} d\left(\boldsymbol{x}_{\pi_{1}(i)}, \boldsymbol{y}_{\pi_{2}(i)} \right)$$

$$x_{4} \quad D_{41} \quad D_{42} \quad D_{43} \quad D_{44} \quad D_{45} \quad D_{40} \quad D_{47}$$

$$x_{3} \quad D_{31} \quad D_{32} \quad D_{33} \quad D_{31} \quad D_{35} \quad D_{36} \quad D_{37}$$

$$x_{2} \quad D_{21} \quad D_{22} \quad D_{23} \quad D_{24} \quad D_{25} \quad D_{26} \quad D_{27}$$

$$x_{1} \quad D_{12} \quad D_{13} \quad D_{14} \quad D_{15} \quad D_{16} \quad D_{17}$$

$$y_{1} \quad y_{2} \quad y_{3} \quad y_{4} \quad y_{5} \quad y_{6} \quad y_{7}$$

DP Computation

$$C_{ij}^{\star} = \min_{\pi \in \mathcal{A}(i,j)} C_{\mathbf{x}_1^i, \mathbf{y}_1^j}(\pi).$$

x_5	D ₅₁	D ₅₂	D ₅₃	D ₅₄	D ₅₅	D ₅₆	D ₅₇
x_4	D ₄₁	$egin{array}{ccc} D_{42} \ C_{42}^{\star} \end{array}$	D ₄₃	D ₄₄	D ₄₅	D ₄₆	D ₄₇
x_3	D ₃₁	D ₃₂	D ₃₃	D ₃₄	D ₃₅	D ₃₆	D ₃₇
x_2	<i>D</i> ₂₁	D ₂₂	D ₂₃	D ₂₄	D ₂₅	D ₂₆	D ₂₇
x_1	<i>D</i> ₁₁	<i>D</i> ₁₂	D ₁₃	D ₁₄	D ₁₅	D ₁₆	D ₁₇
	y_1	y_2	y_3	y_4	y_5	y_6	y_7

Bellman Recursion



Bellman recursion: for all $i \leq n-1, j \leq m-1$,

 $C_{i+1,j+1}^{\star} = \min(C_{i+1,j}^{\star}, C_{ij}^{\star}, C_{i,j+1}^{\star}) + D_{i+1,j+1}$










DTW Strengths



Dynamic Time Warping Matching

DTW as LP

Let
$$X = (x_1, \dots, x_n)$$
 and $Y = (y_1, \dots, y_m)$.
 $U(n, m) \stackrel{\text{def}}{=} \operatorname{co}\{\pi, (n, m) \text{ alig. mat.}\} \subset [0, 1]^{n \times m}$
 $M_{XY} \stackrel{\text{def}}{=} [D(x_i, y_j)^p]_{ij}$

$$\boldsymbol{\pi} = \begin{bmatrix} \mathbf{1} & \mathbf{0} & \cdots & \cdots & \cdots \\ \mathbf{1} & \mathbf{0} & \cdots & \cdots & \cdots \\ \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \cdots \\ \vdots & \vdots & \mathbf{1} & \ddots & \vdots \\ \cdots & \cdots & \mathbf{1} & \mathbf{1} \end{bmatrix} \xrightarrow{\boldsymbol{x_1}} \begin{bmatrix} \boldsymbol{x_1} \\ \vdots \\ M_{\boldsymbol{X}\boldsymbol{Y}} \stackrel{\text{def}}{=} \boldsymbol{x_i} \\ \vdots \\ \boldsymbol{x_n} \end{bmatrix} \begin{bmatrix} D(\boldsymbol{x_i}, \boldsymbol{y_j})^p \\ \vdots \\ \boldsymbol{x_n} \end{bmatrix}$$

 $y_1 \quad \dots \quad y_j \quad \dots \quad y_m$

DTW Problem







$$d\mathbf{tw}(\boldsymbol{X}, \boldsymbol{Y}) = \min_{\boldsymbol{\pi} \in \boldsymbol{U}(\boldsymbol{n}, \boldsymbol{m})} \langle \boldsymbol{\pi}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$
$$\min^{\gamma} \{a_{1}, \dots, a_{n}\} := \begin{cases} \min_{i \leq n} a_{i}, & \gamma = 0, \\ -\gamma \log \sum_{i=1}^{n} e^{-a_{i}/\gamma}, & \gamma > 0. \end{cases}$$
$$d\mathbf{tw}^{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = \min_{\boldsymbol{\pi} \in \boldsymbol{U}(\boldsymbol{n}, \boldsymbol{m})}^{\gamma} \langle \boldsymbol{\pi}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$d\mathbf{tw}(\boldsymbol{X}, \boldsymbol{Y}) = \min_{\boldsymbol{\pi} \in \boldsymbol{U}(\boldsymbol{n}, \boldsymbol{m})} \langle \boldsymbol{\pi}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$
$$\min^{\gamma} \{a_{1}, \dots, a_{n}\} := \begin{cases} \min_{i \leq n} a_{i}, & \gamma = 0, \\ -\gamma \log \sum_{i=1}^{n} e^{-a_{i}/\gamma}, & \gamma > 0. \end{cases}$$
$$d\mathbf{tw}^{\gamma}(\boldsymbol{X}, \boldsymbol{Y}) = \min_{\boldsymbol{\pi} \in \boldsymbol{U}(\boldsymbol{n}, \boldsymbol{m})}^{\gamma} \langle \boldsymbol{\pi}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

Differentiation of sDTW

Algorithm 2 Computes $dtw_{\gamma}(\mathbf{x}, \mathbf{y})$ and $\nabla_{\mathbf{x}} dtw_{\gamma}(\mathbf{x}, \mathbf{y})$

```
1: Inputs: x, y, smoothing \gamma \ge 0, distance function \delta.
  2: \Delta = [\delta(x_i, y_j)]_{i,j}.
  3: r_{0,0} = 0; r_{i,0} = r_{0,j} = \infty; i \in [[n]], j \in [[m]].
                                               Forward recursion
  4: for j = 1, ..., m do
           for i = 1, ..., n do
  5:
         r_{i,j} = \delta_{i,j} + \min^{\gamma} \{ r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1} \}
  6:
            end for
  7:
  8: end for
  9: \delta_{i,m+1} = \delta_{n+1,j} = 0, i \in [n], j \in [m]
10: e_{i,m+1} = e_{n+1,j} = 0, i \in [[n]], j \in [[m]]
11: r_{i,m+1} = r_{n+1,j} = -\infty, i \in [[n]], j \in [[m]]
12: \delta_{n+1,m+1} = 0, e_{n+1,m+1} = 1, r_{n+1,m+1} = r_{n,m}
13: for j = m, ..., 1 do
                                           ▷ Backward recursion
           for i = n, ..., 1 do
14:
15: a = \exp \frac{1}{\gamma} (r_{i+1,j} - r_{i,j} - \delta_{i+1,j})
16: b = \exp \frac{1}{\gamma} (r_{i,j+1} - r_{i,j} - \delta_{i,j+1})

r_{i-1,j} = \frac{1}{r_{i-1,j}} \exp \frac{1}{\gamma} (r_{i+1,j+1} - r_{i,j} - \delta_{i+1,j+1}) r_{i-1,j+1}
           e_{i,j} \stackrel{}{=} \vartheta_{i,j+1,j} \cdot a + e_{i,j+1} \cdot \vartheta_{i,j} \cdot e_{i+1,j+1} \cdot c
18:
19: r_{i,j} end for
                             r_{i,j}
                                                                         r_{i,j+1}
20: end for
21: Output: \delta_{i \neq 1, j}(\mathbf{x}, \mathbf{y}) = r_{n,m} \delta_{i+1, j+1}
22: r_{i+1,j-1} = r_{i+1,j} \left( \frac{\partial \Delta(\mathbf{x},\mathbf{y})}{\partial \mathbf{x}} \right)^T E^{r_{i+1,j+1}}
```

Automatic Differentiation





Applications: sDTW as Loss



Applications: sDTW as Loss

