# Fast Computation of Wasserstein Distances and Applications to Parameter Estimation

#### Marco Cuturi Kyoto University





Joint work with G. Peyré, G. Carlier (Dauphine), J. Solomon (Princeton/MIT), F. de Goes (Pixar), A. Gramfort (ParisTech), J.D. Benamou (INRIA), V. Seguy, A. Rolet (Kyoto), G. Montavon and K.R. Müller (TU Berlin) and more...

#### A geometric toolbox to compare probability measures supported on a metric space.



Monge



Kantorovich



Dantzig



Wasserstein









McCann

Villani















# Why is it relevant to ML?

- New geometry for statistical modeling
  - Information geometry is crucial in stats.
  - That geometry is often KL (e.g. MLE).
- New algorithms to study histogram data
  - Bags-of-features are everywhere.
  - Knowledge on these features is often known but not used.

# Why now?

- Key results in maths since '95~
  - [McCann'95], [JKO'98], [Benamou'98],
     [Ambrosio'06], [Villani'03/'09]
- More work in CV/TCS/Graphics since '98~
  - Earth Mover's Distance [Rubner'98], Embeddings [Indyk'03] Google "earth mover"
    - Scholar

About 8,640 results

- Longstanding roadblock: computation
  - Regularization [C.'13] can provide the key

# Outline

- Definitions: The Wasserstein Distances
- Fast computations with regularization
- Wasserstein variational problems
  - barycenters
  - dictionary learning
  - PCA
  - minimum Kantorovich estimation

Definitions: Couplings & Wasserstein

Assume  $(\Omega, \mathbf{D})$  is a probability space endowed with a metric.

For  $\mu, \nu$  probability measures in  $\mathcal{P}(\Omega)$ ,

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

## Couplings



## Couplings



#### Wasserstein Distance

**Def.** For 
$$p \ge 1$$
, the *p*-Wasserstein distance  
between  $\boldsymbol{\mu}, \boldsymbol{\nu}$  in  $\mathcal{P}(\Omega)$  is  
 $W_p(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \left( \inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \mathbb{E}_{\boldsymbol{P}}[D(X, Y)^p] \right)^{1/p}.$ 

#### Wasserstein on 2 Diracs



### Wasserstein on Uniform Measures



### Wasserstein on Uniform Measures



# Optimal Assignment C Wasserstein



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

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



Assume 
$$\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 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  
 $W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$ 











### Regularized Optimal Transport



# Entropic Regularization [Wilson'62]

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

**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

# Entropic Regularization [Wilson'62]

![](_page_30_Figure_1.jpeg)

**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  
 $P_{\gamma} = \mathbf{D}(\boldsymbol{u}) \ K \mathbf{D}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$ 

• [Sinkhorn'64] fixed-point iterations for (u, v)

$$\boldsymbol{u} \leftarrow \boldsymbol{a} / \boldsymbol{K} \boldsymbol{v}, \quad \boldsymbol{v} \leftarrow \boldsymbol{b} / \boldsymbol{K}^T \boldsymbol{u}$$

• Fast, O(nm) or less, GPU parallel [C'13].

# Regularized Transport: Fast

![](_page_32_Figure_1.jpeg)

# Regularized Transport: differentiable

Prop. Gradients w.r.t 
$$\boldsymbol{a}, \boldsymbol{X}$$
: [CD'14]  
1.  $W_{\gamma} = \max_{\alpha,\beta} \alpha^{T} \boldsymbol{a} + \beta^{T} \boldsymbol{b} - \frac{1}{\gamma} (e^{\alpha/\gamma})^{T} K e^{\beta/\gamma}$   
2.  $W_{\gamma}$  is convex w.r.t.  $\boldsymbol{a}$ ;  $\nabla_{\boldsymbol{a}} W_{\gamma} = \gamma \log(\boldsymbol{u})$ .  
3. If  $p = 2, \Omega = \mathbb{R}^{d}$ ,  
 $\nabla_{\boldsymbol{X}} W_{\gamma} = \boldsymbol{X} \mathbf{D}(\boldsymbol{a}^{\frac{1}{2}}) - \boldsymbol{Y} P_{\gamma}^{T} \mathbf{D}(\boldsymbol{a}^{-\frac{1}{2}})$ 

# Regularized Transport: duality

**Prop.** Writing  $H_{\nu} : a \mapsto W_{\gamma}(\mu, \nu), [CP'15]$ 

**1.** The Legendre transform of  $H_{\nu}$  has a **closed form**:

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

**2.** By Fenchel duality, if f concave on  $\Sigma_n$ ,

$$\min_{\boldsymbol{a}\in\Sigma_n} W_{\gamma}(\boldsymbol{\mu},\boldsymbol{\nu}) - f(\boldsymbol{a}) = \max_{\boldsymbol{g}\in\mathbb{R}^n} f_*(\boldsymbol{g}) - H_{\boldsymbol{\nu}}^*(\boldsymbol{g})$$

# Regularized Transport: duality

**Prop.** Writing  $H_{\nu} : a \mapsto W_{\gamma}(\mu, \nu), [CP'15]$ 

1. The Legendre transform of  $H_{\nu}$  has a closed form:

Optimizing over measures with the Wasserstein metric is crucial to use OT in statistics / machine learning.

### Variational Wasserstein Problems

$$\min_{\boldsymbol{\mu}\in Q\subset\mathcal{P}(\Omega)}F\left(\boldsymbol{\mu},W_{p}^{p}(\boldsymbol{\mu},\boldsymbol{\nu_{1}}),W_{p}^{p}(\boldsymbol{\mu},\boldsymbol{\nu_{2}}),\cdots,W_{p}^{p}(\boldsymbol{\mu},\boldsymbol{\nu_{N}})\right)$$

• k-means Algorithm [Lloyd'82]

$$\min_{\substack{\boldsymbol{\mu} \in \mathcal{P}(\mathbb{R}^d) \\ |\operatorname{supp} \boldsymbol{\mu}| = k}} W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}_{data})$$

• [McCann'95] Interpolant

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

# Variational Wasserstein Problems

• [JKO'98] gradient flow

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

[Agueh'11] Wasserstein barycenters
 Wasserstein Dictionary Learning [RCP'15]
 [Bigot'15] Wasserstein PCA [SC'15]
 [Bassetti'06] Min. Kantorovich Estimation
 [MMC'15]

### 1. Wasserstein Barycenters

N $\min_{\boldsymbol{\mu}\in\mathcal{P}(\Omega)}\sum_{i=1}^{\infty}\lambda_i W_p^p(\boldsymbol{\mu},\boldsymbol{\nu_i})$  $p_{\theta'}$ Wasserstein  $P(\Omega)$ Barycenter [Agueh'11]  $p_{ heta^{\prime\prime}}$ 28

# LP Formulations

Can solve exactly this problem with empirical measures in 2-Wasserstein case, MM-OT:
 If | supp *ν<sub>i</sub>* | = *n<sub>i</sub>*, LP of size (∏<sub>i</sub> *n<sub>i</sub>*, ∑<sub>i</sub> *n<sub>i</sub>*)

![](_page_39_Figure_2.jpeg)

## LP Formulations

• If solving on a grid (all locations fixed), LP:

$$\min_{P_1, \cdots, P_N, \boldsymbol{a}} \sum_{i=1}^N \lambda_i \langle P_i, M \rangle$$
  
s.t.  $P_i^T \mathbf{1}_d = \boldsymbol{b_i}, \forall i \leq N,$   
 $P_1 \mathbf{1}_d = \cdots = P_N \mathbf{1}_d = \boldsymbol{a}.$ 

# Primal Descent on Regularized W

![](_page_41_Picture_1.jpeg)

# Primal Descent on Regularized W

[CD'14]

![](_page_42_Picture_2.jpeg)

![](_page_42_Figure_3.jpeg)

# Primal Descent on Regularized W

![](_page_43_Figure_1.jpeg)

![](_page_43_Figure_2.jpeg)

### Regularized OT as KL Projection

$$\mathbf{KL}(P \mid \mathbf{K}) = \sum_{ij} P_{ij} \log (P_{ij} / \mathbf{K}_{ij})$$
$$\langle P, M_{\mathbf{XY}} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P \mid \mathbf{K})$$

Prop. 
$$P_{\gamma} = \operatorname{Proj}_{C_{a} \cap C_{b}'}(K)$$
  
 $C_{a} = \{P | P\mathbf{1}_{m} = a\}, C_{b}' = \{P | P^{T}\mathbf{1}_{n} = b\}$ 

$$\begin{vmatrix} \mathbf{Prop.} \ P_{\gamma} = \operatorname{Proj}_{C_{a} \cap C_{b}'}(\mathbf{K}) \\ C_{a} = \{P | P\mathbf{1}_{m} = \mathbf{a}\}, \ C_{b}' = \{P | P^{T}\mathbf{1}_{n} = \mathbf{b}\} \end{vmatrix}$$

$$\operatorname{Proj}_{C_{a}}(P) = \mathbf{D}\left(\frac{a}{P\mathbf{1}_{m}}\right)P,$$
$$\operatorname{Proj}_{C_{b}'}(P) = P \mathbf{D}\left(\frac{b}{P^{T}\mathbf{1}_{n}}\right).$$

Sinkhorn = Dykstra's alternate projection K •
 Only need to store & update diagonal multipliers

### Wasserstein Barycenter = KL Projections

$$\langle P, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P \mid \boldsymbol{K})$$
$$\min_{\boldsymbol{a}} \sum_{i=1}^{N} W_{\gamma}(\boldsymbol{a}, \boldsymbol{b}_{i}) = \min_{\substack{\mathbf{P} = [\boldsymbol{P}_{1}, \dots, \boldsymbol{P}_{N}]\\ \mathbf{P} \in \boldsymbol{C}_{1} \cap \boldsymbol{C}_{2}}} \sum_{i=1}^{N} \lambda_{i} \mathbf{KL}(\boldsymbol{P}_{i} \mid \boldsymbol{K})$$
$$\boldsymbol{C_{1}} = \{\mathbf{P} \mid \exists \boldsymbol{a}, \forall i, P_{i} \mathbf{1}_{m} = \boldsymbol{a}\}$$
$$\boldsymbol{C_{2}} = \{\mathbf{P} \mid \forall i, P_{i}^{T} \mathbf{1}_{n} = \boldsymbol{b}_{i}\}$$

### Wasserstein Barycenter = KL Projections

![](_page_47_Figure_1.jpeg)

![](_page_47_Figure_2.jpeg)

# Wasserstein Barycenter = KL Projections

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

u=ones(size(B)); % d x N matrix
while not converged
v=u.\*(K'\*(B./(K\*u))); % 2(Nd^2) cost
u=bsxfun(@times,u,exp(log(v)\*weights))./v;
end
a=mean(v,2);

![](_page_49_Picture_1.jpeg)

Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15

![](_page_50_Picture_1.jpeg)

Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15

![](_page_51_Picture_1.jpeg)

35

![](_page_52_Picture_1.jpeg)

Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15

![](_page_53_Picture_1.jpeg)

Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains, SIGGRAPH'15

# Applications: Brain Imaging

![](_page_54_Picture_1.jpeg)

Fast Optimal Transport Averaging of Neuroimaging Data Information Processing in Medical Imaging 2015

# Applications: Brain Imaging

![](_page_55_Figure_1.jpeg)

A Smoothed Dual Approach for Variational Wasserstein Problems to appear in SIAM Imaging Sciences 37

#### 2. Dictionary Learning

$$\min_{\boldsymbol{A}\in(\Sigma_n)^K,\boldsymbol{\Lambda}\in(\Sigma_K)^N}\sum_{i=1}^N W_{\gamma}\left(\boldsymbol{b_i},\sum_{k=1}^K\boldsymbol{\Lambda_k^i\boldsymbol{a_k}}\right)$$

![](_page_56_Figure_2.jpeg)

#### 2. Dictionary Learning

![](_page_57_Figure_1.jpeg)

### 2. Dictionary Learning

![](_page_58_Picture_1.jpeg)

#### 3. Wasserstein PCA

![](_page_59_Figure_1.jpeg)

# Wasserstein PCA vs. Euclidean PCA

![](_page_60_Picture_1.jpeg)

### Generalized Principal Geodesics

![](_page_61_Figure_1.jpeg)

# [Ambrosio'06] Generalized Geodesics

$$\min_{\boldsymbol{v_1}, \boldsymbol{v_2} \in L^2(\bar{\boldsymbol{\nu}}, \Omega) } \sum_{i=1}^N \min_{t \in [0,1]} W_2^2 \left( g_t(\boldsymbol{v_1}, \boldsymbol{v_2}), \boldsymbol{\nu_i} \right) + \lambda R(\boldsymbol{v_1}, \boldsymbol{v_2}),$$
subject to 
$$\begin{cases} g_t(\boldsymbol{v_1}, \boldsymbol{v_2}) = \left( \operatorname{Id} - \boldsymbol{v_1} + t(\boldsymbol{v_1} + \boldsymbol{v_2}) \right) \# \bar{\boldsymbol{\nu}} \\ \operatorname{Id} - \boldsymbol{v_1} \text{ and } \operatorname{Id} + \boldsymbol{v_2} \end{cases}$$
are Monge maps from  $\bar{\boldsymbol{\nu}}$ 

![](_page_62_Figure_2.jpeg)

# Generalized Principal Geodesics

![](_page_63_Picture_1.jpeg)

MNIST data

### **4.** Minimum Kantorovich Estimation

$$\theta^{\star} = \operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} W_p^p(\boldsymbol{p_{\theta}}, \boldsymbol{\nu_{\text{data}}})$$

![](_page_64_Figure_2.jpeg)

### 4. Minimum Kantorovich Estimation

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

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

- Application to Boltzmann machines [MMC'15] using contrastive divergence, better performance in denoising, data completion.
- Important statistical regularization / stochastic optimization problems.

# Minimum Kantorovich Estimation

![](_page_66_Figure_1.jpeg)

# To conclude

- OT has deep/rich mathematical foundations.
- Regularized OT [C.'13] provides a convenient idea to import these ideas into stats/ML.
- Adopted in graphics/imaging, now in stats, ML, data fusion, Bayesian computation.