Nasserstein Dictionary Learning

Applications

Outlook 0000

Wasserstein Dictionary Learning

Collaborators: Matthieu Heitz, Nicolas Bonneel, Fred Ngolè, David Coeurjolly, Marco Cuturi, Gabriel Peyré, Jean-Luc Starck morgan.schmitz@cea.fr

arXiv:1708.01955

26/01/2018



Morgan A. Schmitz et al.

CosmoStat Day on Machine Learning in Astrophysics

1/25

Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook 0000
Dictionary Lear	rning		

- Data $X \in \mathbb{R}^{N \times M}$, dictionary $D \in \mathbb{R}^{N \times S}$, codes $\Lambda \in \mathbb{R}^{S \times M}$
- Learn *D* and Λ such that $X \approx D\Lambda$
- Throw in whichever requirements suits you:
 - Compact representation? Take *S* < *N* (dimensionality reduction)
 - Sparsity? Add *I*₀ or *I*₁ constraint
 - Positivity? Add a constraint! (NMF)
- Ultimately, data is reconstructed through *linear* combinations of dictionary atoms



Optimal Transport

Wasserstein Dictionary Learning

Applications

Outlook

Dictionary Learning Toy example



9

Optimal	

Outline



2 Wasserstein Dictionary Learning

3 Applications

- Toy examples
- Point Spread Function





Optimal Transport ●000	Wasserstein Dictionary Learning	Applications	Outlook 0000
Optimal Transp	ort		

Overview



Figure: Graphical representation of the mass transportation problem: find the optimal way of moving a heap of sand μ into a hole ν knowing the cost of moving grains of sand to and from any position



Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook 0000
Wasserstein Di Discrete setting	stances		

- Knowledge of the cost of moving: cost matrix $C \in \mathbb{R}^{N \times N}$
- Optimal Transport distance defined as:

$$W(\mu,
u) = \min_{T\in\Pi(\mu,
u)} \langle T, C
angle$$

- Where $\Pi(\mu, \nu) := \left\{ T \in \mathbb{R}^{N \times N}_+, T \mathbb{1}_N = \mu, T^\top \mathbb{1}_N = \nu \right\}$ is the set of acceptable *transport plans*
- Recent numerical schemes allow for fast computation of (approximate) Wasserstein distances¹



¹Cuturi, 2013

Optimal Transport oo●o	Wasserstein Dictionary Learning	Applications	Outlook 0000
Barycenters			

Inputs D = (d₁,..., d_S) and weights λ = (λ₁,..., λ_S)
Euclidean barycenter:

$$\mathcal{B}(D,\lambda) := \operatorname{argmin}_{u \in \mathbb{R}^N} \sum_{s=1}^S \lambda_s \|u - d_s\|_2^2 = \sum_{s=1}^S \lambda_s d_s$$



Figure: Euclidean simplex



1

Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook 0000
Barycenters			

• Similarly, define Wasserstein barycenters as:

$$P(D, \lambda) = \operatorname*{argmin}_{u} \sum_{s=1}^{S} \lambda_{s} W_{\gamma}(u, d_{s}) = \ldots?$$

 Iterative scheme to compute approximate Wasserstein barycenter: P^(L)(D, λ) ≈ P(D, λ) after L iterations



Figure: Wasserstein simplex

Outlook

Outline

Optimal Transport

2 Wasserstein Dictionary Learning

3 Applications

- Toy examples
- Point Spread Function

4 Outlook



Optimal	Transport

Wasserstein Dictionary Learning

Applications

Outlook

Barycenters Simplices comparison





10/25

Wasserstein Dictionary Learning

- Idea: instead of $X \approx D\Lambda$, how about $X \approx P(D, \Lambda)$?
- Optimization problem:

$$\underset{D,\Lambda}{\operatorname{argmin}} \mathcal{E}(D,\Lambda) := \sum_{i=1}^{M} L(x_i, P(D,\lambda_i))$$

• Gradient descent in both dictionary and weights: need $\partial_D P, \partial_\lambda P$



Optimal Transport

Wasserstein Dictionary Learning

Application

Outlook 0000

WDL & deep learning Automatic differentiation

- Automatic differentiation: differentiate algorithm instead of analytical operator
- WDL: gradients are obtained by automatic differentiation of iterative scheme to compute Wasserstein barycenters
- Deep learning: gradients are obtained by backprop!



Optimal Transport

Wasserstein Dictionary Learning

Applications

Outlook 0000

WDL & deep learning WDL & autoencoders

- Autoencoders: learn both encoder and decoder nets
- Learning a Wasserstein dictionary similar to learning one very specific such pair:
 - Encode new datapoint through its barycentric weights in the atoms' simplex;
 - The Wasserstein barycenter operator $P^{(L)}(D,.)$ is the decoder.



13/25

Optimal	

Outline

Optimal Transport

2 Wasserstein Dictionary Learning

3 Applications

- Toy examples
- Point Spread Function

Outlook



Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook 0000
Tov examples			





Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook 0000
Toy examples Handwritten digits			



• WDL recovers Principal Wasserstein Geodesics²!

²Seguy & Cuturi 2015

9

Optimal Transport

Wasserstein Dictionary Learning

Application:

Outlook

Chromatic variation of the PSF



Figure: Simulated Euclid-like PSF variation for different wavelengths (400 to 900nm)



Figure: Euclidean barycenters





Figure: Wasserstein displacement interpolation

17/25

Morgan A. Schmitz et al.

CosmoStat Day on Machine Learning in Astrophysics

Chromatic variation of the PSF

- Training data: set of simulated Euclid-like PSFs at different wavelengths
- Learn 2 atoms on a set of monochromatic PSFs, initialize as constant images
- Initialize weights as projected wavelength



Optimal Transport

Wasserstein Dictionary Learning

Application

Outlook

Chromatic variation of the PSF





19/2

Morgan A. Schmitz et al.

CosmoStat Day on Machine Learning in Astrophysics

Optimal Transport

Wasserstein Dictionary Learning

Application

Outlook

Chromatic variation of the PSF Reconstruction



Figure: Original PSFs (top row) and their reconstruction with our method (bottom row)



Optimal	

Outline

Optimal Transport

2 Wasserstein Dictionary Learning

3 Applications

- Toy examples
- Point Spread Function





Toward chromatic decomposition of star images Real-life setting

- Stars give measurement of the PSF integrated with their own SED S_{*}
- Our objective is to obtain the integrated PSF for a given galaxy's SED S_● ≠ S_∗
- Need to decompose PSF into 'monochromatic' components:

$$Y = \sum_{i=1}^n S_*(\lambda_i) H_i$$

- Y: observed star image;
- $S_*(\lambda_1), \ldots, S_*(\lambda_n)$: discretized star SED;
- H_1, \ldots, H_n : monochromatic PSFs.

22/25

Toward chromatic decomposition of star images

- Use Optimal Transport to break degeneracy
- Reconstruct star image as weighted sum of intermediary steps (Wasserstein barycenters) in the transportation of extreme-wavelengths PSFs
- Use whole field of stars to learn basis set of transport plans



Optimal Transport

Wasserstein Dictionary Learning

Applications

Outlook

Toward chromatic decomposition of star images λRCA (preliminary results)



Figure: From top to bottom: observed polychromatic PSF, true monochromatic PSFs, reconstructed monochromatic PSFs



Optimal Transport	Wasserstein Dictionary Learning	Applications	Outlook ooo●
Summary			

- Wasserstein Dictionary Learning: a new, Optimal Transport-based representation learning method
- Acknowledgements and funding:



morgan.schmitz@cea.fr
http://www.cosmostat.org/people/mschmitz/



Appendix



- 6 Additional results
- Automatic differentiation
- 8 Going further: unbalanced WDL



Optimal Transport - WDL Graphical representation



Figure: Graphical representation of the WDL approach. Source: Bonneel et al, 2016

9

Wasserstein Distances Continuous setting

 The optimal cost achieved in the mass transportation problem defines an Optimal Transport distance between any two measures μ, ν on some space Ω:

$$W(\mu,
u) \mathrel{\mathop:}= \min\left\{\int_{\Omega imes\Omega} oldsymbol{c}(oldsymbol{x},oldsymbol{y}) oldsymbol{d} \pi(oldsymbol{x},oldsymbol{y}), \ \pi\in \Pi(\mu,
u)
ight\}$$

 Where Π(μ, ν) is the set of bivariate measures with marginals μ, ν:

•
$$\int_{\Omega} \pi(x, y) dx = \nu(y)$$

• $\int_{\Omega} \pi(x, y) dy = \mu(x)$



Wasserstein Distances Discrete setting

• Discrete case: if $|\Omega| = N$, *c* reduces to $C \in \mathbb{R}^{N \times N}$ and:

$$W(\mu,
u) = \min_{T\in \Pi(\mu,
u)} \langle T, C
angle$$

- Where $\Pi(\mu, \nu) := \left\{ T \in \mathbb{R}^{N \times N}_+, T \mathbb{1}_N = \mu, T^\top \mathbb{1}_N = \nu \right\}$ is the set of acceptable *transport plans*
- Recent numerical schemes allow for fast computation of (approximate) Wasserstein distances³



³Cuturi, 2013

Morgan A. Schmitz et al.

Appendix



- 6 Additional results
- Automatic differentiation
- 8 Going further: unbalanced WDL



Numerical Optimal Transport Entropic penalty

- Adding an entropy penalty term:
 W_γ(μ, ν) := min_{T∈Π(μ,ν)}⟨T, C⟩ + γH(T)
- Leads to the Sinkhorn algorithm being useable for linear convergence to W_γ:

•
$$b^{(0)} := \mathbb{1}_N$$

• $a^{(l)} := \frac{\nu}{K^{\top} b^{(l-1)}}$
• $b^{(l)} := \frac{\mu}{Ka^{(l)}}$
• $T^{(L)} := \Delta (b^{(L)}) K\Delta (a^{(L)})$

• Where
$$K = \exp(-C/\gamma)$$



Numerical Optimal Transport Entropic penalty & barycenter computation

 Generalized Sinkhorn algorithm for barycenter computation:

•
$$P^{(l)}(D,\lambda) = \prod_{s=1}^{S} \left(K^{\top} a_{s}^{(l)}\right)^{\lambda_{s}}$$

• $a_{s}^{(l)} = \frac{d_{s}}{K b_{s}^{(l-1)}}$
• $b_{s}^{(l)} = \frac{P^{(l)}(D,\lambda)}{K^{\top} a_{s}^{(l)}}$
Where $K = \exp(-C/\gamma)$



Entropy parameter

Additional results

Automatic differentiation

Going further: unbalanced WDL

Wasserstein simplices Effect of γ



Figure: Wasserstein simplices for different γ values



Appendix





- Automatic differentiation
- 8 Going further: unbalanced WDL



Chromatic variation of the PSF PCA components





Morgan A. Schmitz et a

35/25

Automatic differentiation

Going further: unbalanced WDL

Chromatic variation of the PSF Weights



Figure: Left-hand side: weights reached by our method at convergence. Right-hand side: PCA-learned codes corresponding to the first two principal components



Chromatic variation of the PSF

Back to PSF Application



9

377

Cardiac sequence



Figure: Left: True frames (bottom row) and reconstructions after application of our method. Right: Projection of the weights in the 4 atom basis.



Face editing MUG Dataset



Figure: Top row: isobarycenter of 5 learnt atoms. Bottom row: extrapolation giving 30% weight importance to each of the 5.



Appendix



- 6 Additional results
- Automatic differentiation
 - B Going further: unbalanced WDL



5

Wasserstein Dictionary Learning Automatic differentiation

- Differentiating true P operator impossible and/or comes at prohibitive cost
- *P*^(*L*) is what we are going to use in practice, and is just the result of iteratively applying simple mathematical operations:

 $P^{(L)} = f(f(\ldots(f(\mathbb{1}_N, D, \lambda)), D, \lambda), D, \lambda))$

- Automatic differentiation⁴: differentiate *algorithm* by applying chain rule
- Can be done by specialized libraries⁵ or https://www.specialized-libraries-5

⁵http://deeplearning.net/software/theano/

⁴Griewank & Walther, 2008

Handmade automatic differentiation

$$P^{(l)}(D,\lambda) = \Psi(b^{(l-1)}(D,\lambda), D,\lambda)$$
(1)

$$b^{(l)}(D,\lambda) = \Phi(b^{(l-1)}(D,\lambda), D,\lambda)$$
(2)

where:

$$\begin{split} \Psi(b, D, \lambda) &:= \prod_{s} \left(K^{\top} \frac{d_{s}}{K b_{s}} \right)^{\lambda_{s}} \\ \Phi(b, D, \lambda) &:= \left[\left(\frac{\Psi(b, D, \lambda)}{K^{\top} \frac{d_{1}}{K b_{1}}} \right)^{\top}, \dots, \left(\frac{\Psi(b, D, \lambda)}{K^{\top} \frac{d_{s}}{K b_{s}}} \right)^{\top} \right]^{\top} \end{split}$$

42/25

Handmade automatic differentiation

Introduce:

$$\xi_{y}^{(l)} := \left[\partial_{y}\xi(\boldsymbol{b}^{(l)},\boldsymbol{D},\lambda)\right]^{\top} \qquad \boldsymbol{B}_{y}^{(l)} := \left[\partial_{y}\boldsymbol{b}^{(l)}(\boldsymbol{D},\lambda)\right]^{\top}$$

where ξ can be Ψ or Φ , *y* can be *D* or λ

• And denote:

$$\boldsymbol{v}^{(L-1)} := \Psi_b^{(L-1)} \left(\nabla L(\boldsymbol{P}^{(L)}(\boldsymbol{D}, \boldsymbol{\lambda}), \boldsymbol{x}) \right) \quad (3)$$

$$\forall l < L-1, v^{(l-1)} := \Phi_b^{(l-1)} \left(v^{(l)} \right)$$
 (4)



Morgan A. Schmitz et al.

Handmade automatic differentiation

Then by total differentiation and the chain rule, differentiating (1) yields:

$$\left[\partial_D P^{(l)}(D,\lambda)\right]^{\top} = \Psi_D^{(l-1)} + B_D^{(l-1)} \Psi_b^{(l-1)}$$
(5)

And, differentiating (2):

$$B_D^{(l)} = \Phi_D^{(l-1)} + B_D^{(l-1)} \Phi_b^{(l-1)}$$
(6)



5

Handmade automatic differentiation

We then have, by definitions (3)-(4) and by (5) and (6):

$$\begin{aligned} \nabla_D \mathcal{E}_L(D,\lambda) &= \Psi_D^{(L-1)} \left(\nabla \mathcal{L}(\mathcal{P}^{(L)}(D,\lambda),x) \right) + \mathcal{B}_D^{(L-1)} v^{(L-1)} \\ &= \Psi_D^{(L-1)} \left(\nabla \mathcal{L}(\mathcal{P}^{(L)}(D,\lambda),x) \right) + \Phi_D^{(L-2)} \left(v^{(L-1)} \right) + \\ & \mathcal{B}_D^{(L-2)} \left(v^{(L-2)} \right) \end{aligned}$$

$$\nabla_{D}\mathcal{E}_{L}(D,\lambda) = \Psi_{D}^{(L-1)}\left(\nabla\mathcal{L}(\mathcal{P}^{(L)}(D,\lambda),x)\right) + \sum_{l=0}^{L-2} \Phi_{D}^{(l)}\left(\mathbf{v}^{(l+1)}\right)$$

where the sum starts at 0 because $B_D^{(0)} = 0$ since we initialized $b^{(0)}$ as a constant vector

45/25

Appendix



- 6 Additional results
- Automatic differentiation
- Boing further: unbalanced WDL



Going further: unbalanced WDL

• Recall discrete Wasserstein distance definition:

$$\begin{split} \mathcal{W}(\boldsymbol{\rho},\boldsymbol{q}) &:= \min_{\boldsymbol{T} \in \Pi(\boldsymbol{\mu},\boldsymbol{\nu})} \langle \boldsymbol{T},\boldsymbol{C} \rangle \\ \Pi(\boldsymbol{\mu},\boldsymbol{\nu}) &:= \Big\{ \boldsymbol{T} \in \mathbb{R}_{+}^{N \times N}, \boldsymbol{T} \mathbb{1}_{N} = \boldsymbol{\mu}, \boldsymbol{T}^{\top} \mathbb{1}_{N} = \boldsymbol{\nu} \Big\} \end{split}$$

Identical to:

$$W(p,q) := \min_{T \in \mathbb{R}^{N imes N}_+} \langle T, \mathcal{C}
angle + \iota_{\{\mu\}}(T \mathbb{1}_N) + \iota_{\{
u\}}(T^ op \mathbb{1}_N)$$

Unbalanced Optimal Transport: relaxation by replacing *i* with other similarity criterion, *e.g.*:

$$W_{\mathrm{KL}}(\boldsymbol{\rho},\boldsymbol{q}) \coloneqq \min_{\boldsymbol{\mathcal{T}} \in \mathbb{R}^{N \times N}_{+}} \langle \boldsymbol{\mathcal{T}}, \boldsymbol{\mathcal{C}} \rangle + \mathrm{KL}\left(\mu | \boldsymbol{\mathcal{T}} \mathbb{1}_{N}\right) + \mathrm{KL}\left(\nu | \boldsymbol{\mathcal{T}}^{\top} \mathbb{1}_{N}\right)$$

Going further: unbalanced WDL



Figure: Example datapoints from set of bimodal Gaussians (in blue) and reconstructions from our method (in yellow) in both the balanced and the unbalanced cases

