

Stochastic Optimization for Regularized Wasserstein Estimators

ICML 2020

Marin Ballu



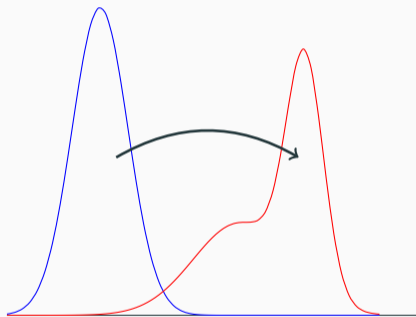
Quentin Berthet



Francis Bach

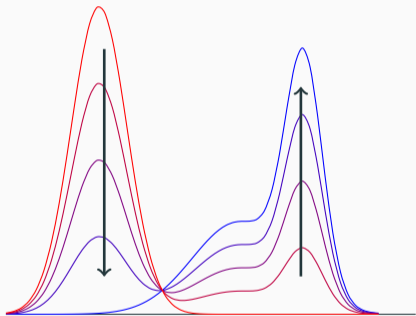


Wasserstein Distance: a natural geometry for distributions



How does one compute the distance between two data distributions?

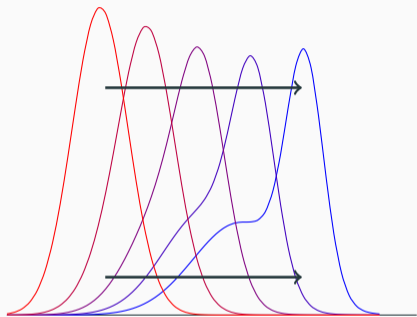
Wasserstein Distance: a natural geometry for distributions



How does one compute the distance between two data distributions?

- Relative entropy and other f-divergences allow classical statistical approaches.

Wasserstein Distance: a natural geometry for distributions

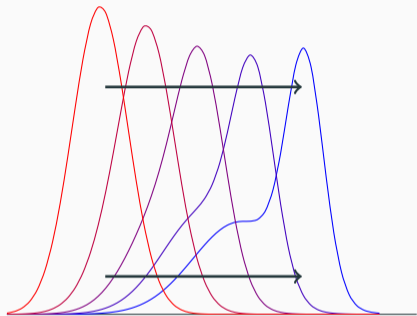


How does one compute the distance between two data distributions?

- Relative entropy and other f-divergences allow classical statistical approaches.
- Optimal transport theory allows us to capture the geometry of the data distributions, with the *Wasserstein distance*.

$$W_c(\mu, \nu) = \text{OT}(\mu, \nu) = \min_{T \# \mu = \nu} \mathbb{E}_{X \sim \mu} [c(X, T(X))]$$

Wasserstein Distance: a natural geometry for distributions

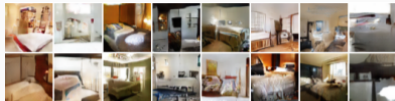


How does one compute the distance between two data distributions?

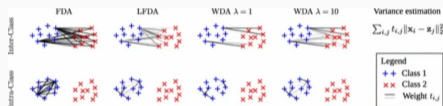
- Relative entropy and other f-divergences allow classical statistical approaches.
- Optimal transport theory allows us to capture the geometry of the data distributions, with the *Wasserstein distance*.

$$W_c(\mu, \nu) = \text{OT}(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)]$$

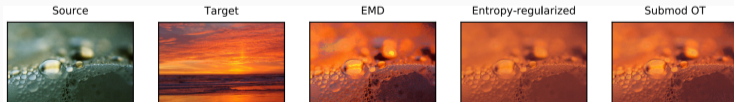
Wasserstein distance in machine learning



Wasserstein GAN (Arjovsky et al., 2017)

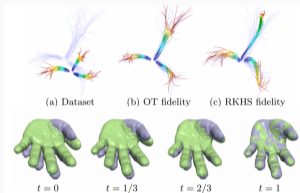


Wasserstein Discriminant Analysis (Flamary et al., 2018)



Clustered point-matching (Alvarez-Melis et al., 2018)

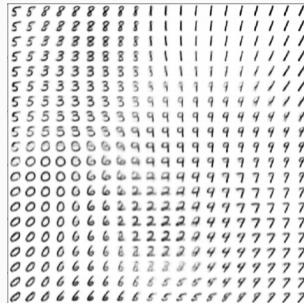
Wasserstein distance in machine learning



Diffeomorphic registration (Feydy et al., 2017)



Alignment of embeddings (Grave et al., 2019)



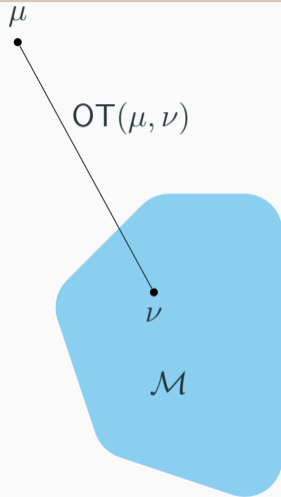
Sinkhorn divergence for generative models (Genevay et al., 2019)

Our contribution

We consider the minimum Kantorovich estimator (Bassetti et al., 2006), or *Wasserstein estimator* of the measure μ :

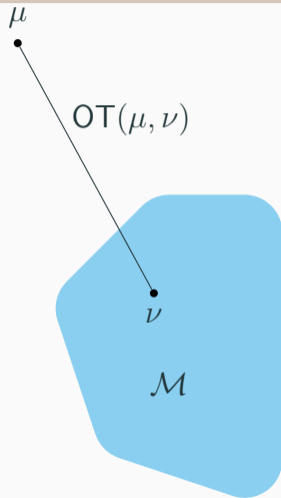
$$\min_{\nu \in \mathcal{M}} \text{OT}(\mu, \nu),$$

which is often used for $\mu = \sum_i \delta_{x_i}$ to fit a parametric model \mathcal{M} (as with MLE, where KL divergence replaces OT).



Our contribution

- We add two layers of entropic regularization.
- We propose a new stochastic optimization scheme to minimize the regularized problem.
- Time per step is sublinear in the natural dimension of the problem.
- We provide theoretical guarantees, and simulations.



Regularized Wasserstein Distance

Wasserstein distance

$$W_c(\mu, \nu) = \text{OT}(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)]$$

Regularized Wasserstein Distance

Wasserstein distance

$$W_c(\mu, \nu) = \text{OT}(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)]$$

Regularized Wasserstein distance

$$\text{OT}_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)] + \varepsilon \text{KL}(\pi, \mu \otimes \nu)$$

Computed at lightspeed by Sinkhorn algorithm (Cuturi 2013)

SGD on dual problem (Genevay et al. 2016)

Wasserstein estimator

$$\min_{\nu \in \mathcal{M}} \text{OT}(\mu, \nu)$$

Regularized Wasserstein Estimator

Wasserstein estimator

$$\min_{\nu \in \mathcal{M}} \text{OT}(\mu, \nu)$$

First layer of regularization

$$\min_{\nu \in \mathcal{M}} \text{OT}_{\varepsilon}(\mu, \nu)$$

Regularized Wasserstein Estimator

Wasserstein estimator

$$\min_{\nu \in \mathcal{M}} \text{OT}(\mu, \nu)$$

First layer of regularization

$$\min_{\nu \in \mathcal{M}} \text{OT}_{\varepsilon}(\mu, \nu)$$

Second layer of regularization

$$\min_{\nu \in \mathcal{M}} \text{OT}_{\varepsilon}(\mu, \nu) + \eta \text{KL}(\nu, \beta)$$

First layer: Gaussian deconvolution

This is a recent interpretation (Rigollet, Weed 2018).

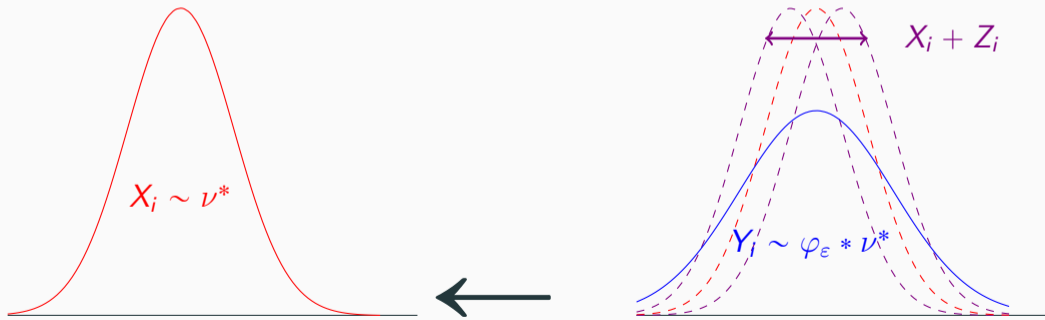
Let X_i be iid random variables following ν^* , $Z_i \sim \varphi_\varepsilon = \mathcal{N}(0, \varepsilon \text{Id})$ an iid gaussian noise and $Y_i = X_i + Z_i$ the perturbed observation with distribution μ .



First layer: Gaussian deconvolution

For $c(x, y) = \|x - y\|^2$, the MLE for ν^* is

$$\hat{\nu} := \arg \max_{\nu \in \mathcal{M}} \sum_i \log(\varphi_\epsilon * \nu)(X_i) \Leftrightarrow \hat{\nu} = \arg \min_{\nu \in \mathcal{M}} \text{OT}_\epsilon(\mu, \nu).$$



First layer: adds entropy to the transport matrix

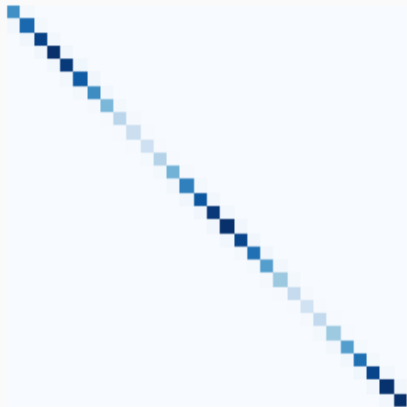


Figure 1: Small regularization $\varepsilon = 0.01$



Figure 2: Big regularization $\varepsilon = 0.1$

Second Layer: Interpolation with likelihood estimators

Wasserstein Estimator

$$\min_{\nu \in \mathcal{M}} \text{OT}(\mu, \nu)$$

Maximum Likelihood Estimator

$$\min_{\nu \in \mathcal{M}} \text{KL}(\nu, \beta)$$

Regularized Wasserstein Estimator

$$\min_{\nu \in \mathcal{M}} \text{OT}_{\epsilon}(\mu, \nu) + \eta \text{KL}(\nu, \beta)$$

Second Layer: adds entropy to the target measure

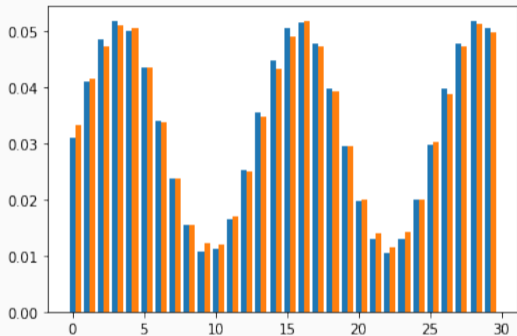


Figure 3: Small regularization $\eta = 0.02$

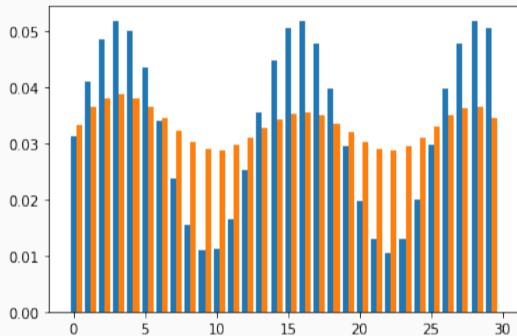


Figure 4: Big regularization $\eta = 0.2$

Dual Formulation of the problem

$$\min_{\nu \in \mathcal{M}} \text{OT}_\varepsilon(\mu, \nu) + \eta \text{KL}(\nu, \beta)$$

with

$$\text{OT}_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)] + \varepsilon \text{KL}(\pi, \mu \otimes \nu)$$

is

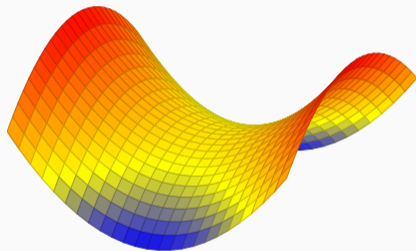
$$\min_{\nu \in \mathcal{M}} \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \pi} [c(X, Y)] + \varepsilon \text{KL}(\pi, \mu \otimes \nu) + \eta \text{KL}(\nu, \beta).$$

We consider the dual of the second min.

Dual Formulation

The dual problem can be written as a saddle point problem, where the min and the max can be swapped. The final formulation is of the form

$$\max_{(a,b) \in \mathbb{R}^I \times \mathbb{R}^J} F(a, b).$$



Properties of the function F in the discrete case

1. F is λ -strongly convex on the hyperplane $E = \{\sum_i \mu_i a_i = \sum_j \beta_j b_j\}$.
2. There exists a solution of $\max_{(a,b) \in \mathbb{R}^I \times \mathbb{R}^J} F(a, b)$, which is in E , and it is unique.
3. The gradients of F can be written as expectations

$$\nabla_a F = \mathbb{E} [(1 - D_{i,j}) e_i],$$

$$\nabla_b F = \mathbb{E} [(f_j - D_{i,j}) e_j].$$

$$\text{with } D_{i,j}(a, b) = \exp\left(\frac{a_i + b_j - C_{i,j}}{\varepsilon}\right) \text{ and } f_j = \frac{v_j(b)}{\beta_j}.$$

Stochastic Gradient Descent

We have stochastic gradients for F

$$G_a = (1 - D_{i,j})e_i$$

$$G_b = (f_j - D_{i,j})e_j.$$

SGD algorithm:

- Sample $i \in \{1, \dots, I\}$ with probability μ_i ,
- Sample $j \in \{1, \dots, J\}$ with probability β_j ,
- Compute G_a and G_b
- $a \leftarrow a + \gamma_t G_a$,
- $b \leftarrow b + \gamma_t G_b$.

Stochastic Gradient Descent

We only have to compute a and b one coefficient at a time

- Sample $i \in \{1, \dots, I\}$ with probability μ_i ,
- Sample $j \in \{1, \dots, J\}$ with probability β_j ,
- Compute f_j and $D_{i,j}$
- $a_i \leftarrow a_i + \gamma_t(1 - D_{i,j})$,
- $b_j \leftarrow b_j + \gamma_t(f_j - D_{i,j})$.

The sum memorization trick

The computation of $D_{i,j}(a, b) = \exp\left(\frac{a_i + b_j - C_{i,j}}{\varepsilon}\right)$ and $f_j = \frac{\nu_j(b)}{\beta_j}$ is $O(1)$.

However

$$\nu_j^* = \frac{\beta_j e^{-b_j/(\eta-\varepsilon)}}{\sum_k \beta_k e^{-b_k/(\eta-\varepsilon)}},$$

but we can do it in $O(1)$ if we update

$$S^{(t)} = \sum_k \beta_k e^{-b_k^{(t)}/(\eta-\varepsilon)},$$

with

$$S^{(t+1)} = S^{(t)} + \beta_j e^{-b_j^{(t+1)}/(\eta-\varepsilon)} - \beta_j e^{-b_j^{(t)}/(\eta-\varepsilon)}.$$

Convergence Bounds

With stepsize $\gamma_t = \frac{1}{\lambda t}$, the estimator verifies

$$\mathbb{E} [\text{KL}(\nu^*, \nu^t)] \leq \frac{C_1}{(\eta - \varepsilon)\lambda^2} \frac{1 + \log t}{t}.$$

With stepsize $\gamma_t = \frac{C_2}{\sqrt{t}}$, the estimator verifies the following bound:

$$\mathbb{E} [\text{KL}(\nu^*, \nu^t)] \leq \frac{C_3}{(\eta - \varepsilon)\lambda} \frac{2 + \log t}{\sqrt{t}}.$$

Simulations

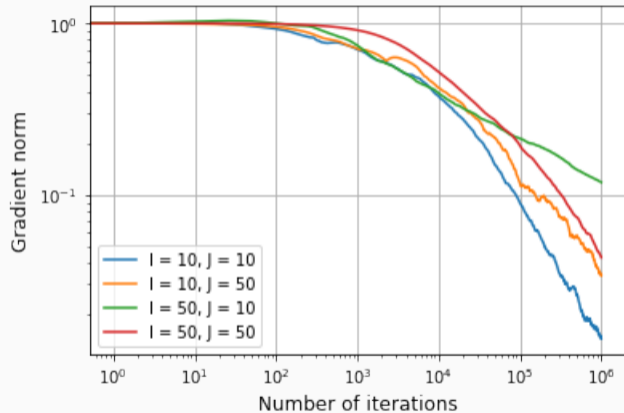


Figure 5: Convergence of the gradient norm for different dimensions.

Wasserstein barycenter

$$\min_{\nu} \sum_{k=1}^K \theta_k \text{OT}(\mu^k, \nu).$$

Doubly regularized Wasserstein barycenter

$$\min_{\nu} \sum_{k=1}^K \theta_k \text{OT}_{\epsilon}(\mu^k, \nu) + \eta \text{KL}(\nu, \beta).$$

Takeaways:

- Wasserstein estimators are "projections" according to Wasserstein distances,
- Two layers of entropic regularization are used here,
- It is then possible to compute stochastic gradients in $O(1)$ for this problem,
- The results are also valid for Wasserstein barycenters.

Thank you for your attention!