

Stochastic Methods for Optimal Transport and Applications in Machine Learning

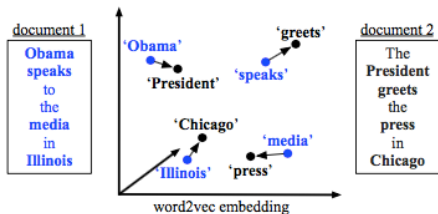
Aude Genevay

CEREMADE - Université Paris Dauphine
INRIA - Mokaplan project-team
DMA - Ecole Normale Supérieure

ISMP - July 2018

Joint work with F. Bach, M. Cuturi, G. Peyré

Comparing High Dimensional Cloud Points



- Document d_i = histogram of words
- Word w_k = point in \mathbb{R}^d for a certain embedding (usually learnt with neural networks, e.g. Word2Vec)
- Document \sim weighted cloud of points in $\mathbb{R}^d \Rightarrow$
 $d_i \sim \mu_i = \sum \alpha_{k,i} \delta_{w_k}$
- Distance between 2 documents d_1, d_2 is the optimal transport distance between the associated point clouds μ_1, μ_2 .

Fitting data to a probabilistic model

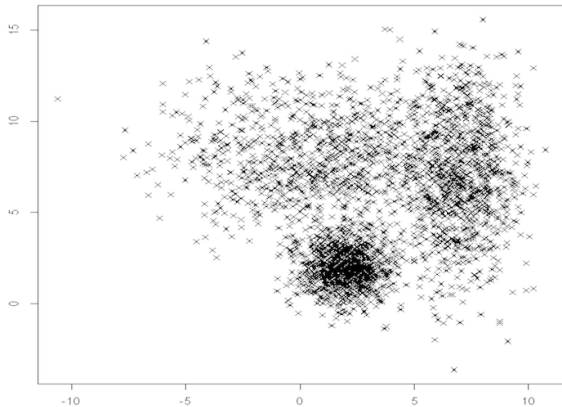


Figure 1: Data points in 2D

Recurrent issue in ML : Fitting data to a probabilistic model

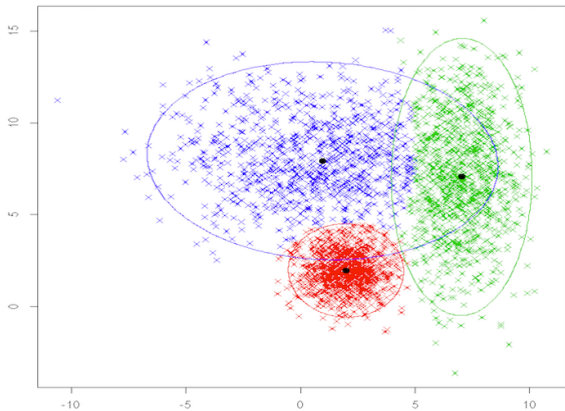


Figure 2: Gaussian Mixture Model

Density Fitting with MLE

- Observed dataset $(y_1, \dots, y_n) \in \mathcal{X}$ (IID assumption)
- Empirical measure $\hat{\nu} = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$
- Parametric model $(\mu_\theta)_{\theta \in \Theta}$ measure with density $(f_\theta)_{\theta \in \Theta}$
- Goal : find $\hat{\theta} = \arg \min_{\theta \in \Theta} \mathcal{L}(\mu_\theta, \hat{\nu})$ where \mathcal{L} is a loss on measures.
- **Maximum Likelihood Estimator**

$$\hat{\theta} \stackrel{\text{def.}}{=} \arg \min_{\theta \in \Theta} - \sum_{i=1}^n \log f(y_i | \theta)$$

Generative Models

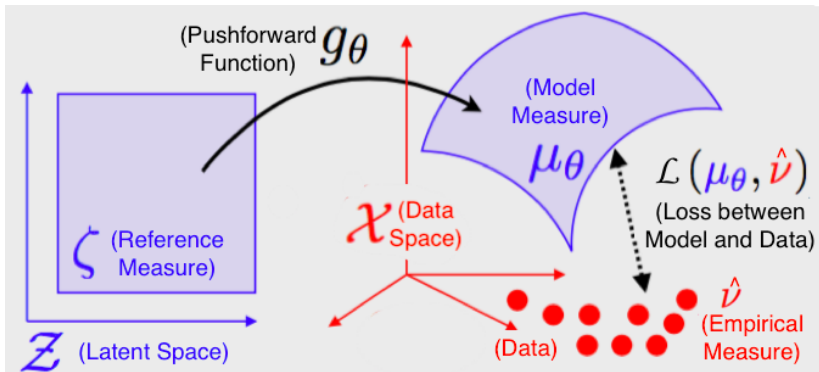
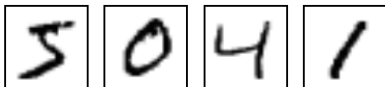


Figure 3: Illustration of Density Fitting on a Generative Model

Density Fitting for Generative Models I

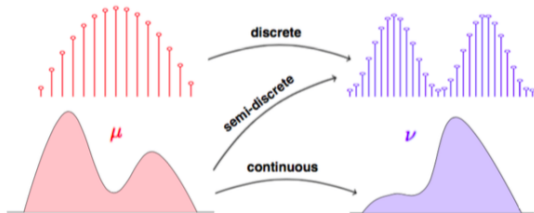
Very popular topic in ML : image generation



- Parametric model : $\mu_\theta = g_{\theta\#}\zeta$
- ζ reference measure on (low dimensional) latent space \mathcal{Z}
- $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ from latent space to data space
- Sampling procedure : $x \sim \mu_\theta$ obtained by $x = g_\theta(z)$ where $z \sim \zeta$
- $\dim \mathcal{Z} \ll \dim \mathcal{X} \Rightarrow \mu_\theta$ doesn't have density wrt Lebesgue measure

\Rightarrow MLE can't be applied in this context!

Optimal Transport I



- Optimal Transport : find coupling that minimizes total cost of moving μ to ν with unit cost function c
- Constrained problem : coupling has fixed marginals
- Minimal cost of moving μ to ν (e.g. solution of the OT problem) is called the **Wasserstein distance** (it's an actual distance!)

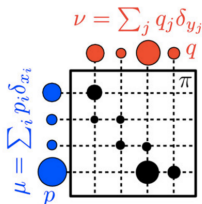
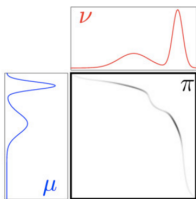
Optimal Transport II

Cost $c(x, y)$ to move a unit of mass from x to y

Constrained set of couplings $\Pi(\mu, \nu)$ with marginals μ and ν

$$W(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y)$$

What's the coupling that minimizes the total cost?



Kantorovitch Formulation of OT

The optimal overall cost for transporting μ to ν is given by

$$W(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\mathcal{P}_\varepsilon)$$

Kantorovitch Formulation of OT

The optimal overall cost for transporting μ to ν is given by

$$W_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \mu \otimes \nu) \quad (\mathcal{P}_\varepsilon)$$

where

$$\text{KL}(\pi | \mu \otimes \nu) \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{Y}} \left(\log \left(\frac{d\pi}{d\mu d\nu}(x, y) \right) - 1 \right) d\pi(x, y)$$

Kantorovitch Formulation of OT

The optimal overall cost for transporting μ to ν is given by

$$W_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \mu \otimes \nu) \quad (\mathcal{P}_\varepsilon)$$

where

$$\text{KL}(\pi | \mu \otimes \nu) \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{Y}} \left(\log \left(\frac{d\pi}{d\mu d\nu}(x, y) \right) - 1 \right) d\pi(x, y)$$

Adding an entropic regularization smoothes the constraint. In particular it yields an unconstrained dual problem.

Dual formulation of OT

$$W(\mu, \nu) = \max_{u \in \mathcal{C}(\mathcal{X}), v \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} u(x) d\mu(x) + \int_{\mathcal{Y}} v(y) d\nu(y) \quad (u, v) \in (\mathcal{D}_\varepsilon)$$

under the constraint that

$$\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, u(x) + v(y) \leq c(x, y)$$

Dual formulation of OT (with entropy)

$$W_\varepsilon(\mu, \nu) = \max_{u \in \mathcal{C}(\mathcal{X}), v \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} u(x) d\mu(x) + \int_{\mathcal{Y}} v(y) d\nu(y) - \iota_{U_c}^\varepsilon(u, v)$$

and the smoothed indicator is

$$\iota_{U_c}^\varepsilon(u, v) \stackrel{\text{def.}}{=} \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right) d\mu(x) d\nu(y)$$

Semi-Dual formulation of OT (with entropy)

The dual problem is convex in u and v . We fix v and minimize over u .

Semi-Dual formulation of OT (with entropy)

The dual problem is convex in u and v . We fix v and minimize over u . This yields :

$$u(x) \stackrel{\text{def.}}{=} -\varepsilon \log \left(\int_{\mathcal{Y}} \exp\left(\frac{v(y) - c(x, y)}{\varepsilon}\right) d\nu(y) \right)$$

Semi-Dual formulation of OT (with entropy)

The dual problem is convex in u and v . We fix v and minimize over u . This yields :

$$u(x) \stackrel{\text{def.}}{=} -\varepsilon \log \left(\int_{\mathcal{Y}} \exp\left(\frac{v(y) - c(x, y)}{\varepsilon}\right) d\nu(y) \right)$$

Plugging back in the dual :

$$\begin{aligned} W_{\varepsilon}(\mu, \nu) &= \max_{\nu \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} -\varepsilon \log \left(\int_{\mathcal{Y}} \exp\left(\frac{v(y) - c(x, y)}{\varepsilon}\right) d\nu(y) \right) d\mu(x) \\ &\quad + \int_{\mathcal{Y}} v(y) d\nu(y) - \varepsilon \\ &= \max_{\nu \in \mathcal{C}(\mathcal{Y})} \mathbb{E}_{\mu} \left[-\varepsilon \log \left(\int_{\mathcal{Y}} \exp\left(\frac{v(y) - c(x, y)}{\varepsilon}\right) d\nu(y) \right) \right] \\ &\quad + \int_{\mathcal{Y}} v(y) d\nu(y) - \varepsilon \end{aligned}$$

We consider 2 frameworks :

- Semi-Discrete : μ is continuous and $\nu = \sum_{j=1}^M \nu_j \delta_{y_j}$ The optimization problem is

$$\max_{\nu \in \mathbb{R}^M} \mathbb{E}_{\mu} \left[-\varepsilon \log \left(\sum_{j=1}^M \exp\left(\frac{\nu(y_j) - c(x, y_j)}{\varepsilon}\right) \right) + \sum_{j=1}^M \nu(y_j) \nu_j^{-\varepsilon} \right]$$

We consider 2 frameworks :

- Semi-Discrete : μ is continuous and $\nu = \sum_{j=1}^M \nu_j \delta y_j$ The optimization problem is

$$\max_{\nu \in \mathbb{R}^M} \mathbb{E}_{\mu} \left[-\varepsilon \log \left(\sum_{j=1}^M \exp\left(\frac{\nu(y_j) - c(x, y_j)}{\varepsilon}\right) \right) + \sum_{j=1}^M \nu(y_j) \nu_j - \varepsilon \right]$$

- Discrete : $\mu = \sum_{i=1}^N \mu_i \delta x_i$ and $\nu = \sum_{j=1}^M \nu_j \delta y_j$ The optimization problem is

$$\max_{\nu \in \mathbb{R}^M} \sum_{i=1}^N \left[-\varepsilon \log \left(\sum_{j=1}^M \exp\left(\frac{\nu(y_j) - c(x_i, y_j)}{\varepsilon}\right) \right) + \sum_{j=1}^M \nu(y_j) \nu_j - \varepsilon \right] \mu_i$$

Stochastic Optimization

Computing the full gradient is

- Hard in the semi-discrete setting (even impossible if we don't know μ explicitly)

Stochastic Optimization

Computing the full gradient is

- Hard in the semi-discrete setting (even impossible if we don't know μ explicitly)
- Very costly in the discrete case since we need to compute N gradients and sum them.

The idea of stochastic optimization is to use approximate gradients so that each iteration is inexpensive.

Stochastic Optimization I

- **Goal** : maximize $H_\varepsilon(\mathbf{v}) = \mathbb{E}_\mu [h_\varepsilon(\mathbf{X}, \mathbf{v})]$ over \mathbf{v} in \mathbb{R}^M .
- Standard gradient ascent :

$$\mathbf{v}^{(k)} = \mathbf{v}^{(k-1)} + \nabla_{\mathbf{v}} H_\varepsilon(\mathbf{v}^{(k-1)})$$

- The whole gradient $\nabla_{\mathbf{v}} H_\varepsilon(\mathbf{v})$ is too costly/complicated to compute
- **Idea** : Sample x from μ and use $\nabla_{\mathbf{v}} h_\varepsilon(x, \mathbf{v})$ as a proxy for the full gradient in the gradient ascent.

Stochastic Optimization II

Algorithm 1 Averaged SGD

Input: C

Output: \bar{v}

$$v \leftarrow \mathbb{0}_M, \bar{v} \leftarrow v$$

for $k = 1, 2, \dots$ do

 Sample x_k from μ

$$v \leftarrow v + \frac{C}{\sqrt{k}} \nabla_v h_\varepsilon(x_k, v) \quad (\text{gradient ascent step})$$

$$\bar{v} \leftarrow \frac{1}{k} v + \frac{k-1}{k} \bar{v} \quad (\text{averaging})$$

end for

- cost of each iteration M
- convergence rate $O(1/\sqrt{k})$

Stochastic Optimization : Case of a Finite Sum I

In the specific case where μ is also a discrete measure, we are minimizing a finite sum of N functionals :

$$\max_{\mathbf{v} \in \mathbb{R}^M} \sum_{i=1}^N \left[-\varepsilon \log \left(\sum_{j=1}^M \exp\left(\frac{\mathbf{v}(y_j) - c(x_i, y_j)}{\varepsilon}\right) \right) + \sum_{j=1}^M \mathbf{v}(y_j) \nu_j - \varepsilon \right] \mu_i$$

Variance reduction algorithms (e.g. SAGA) can be used to improve speed of convergence:

- cost of each iteration M
- convergence rate $O(1/k)$

Numerical Results for Word Mover's Distance (Discrete OT)

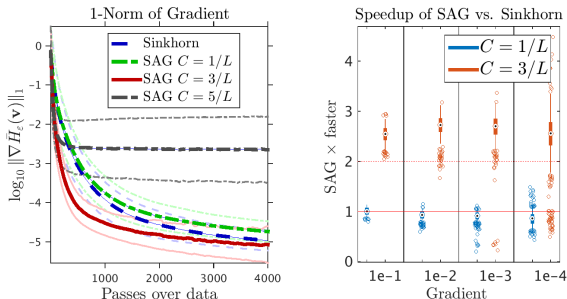
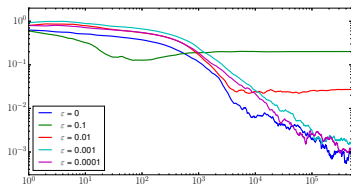
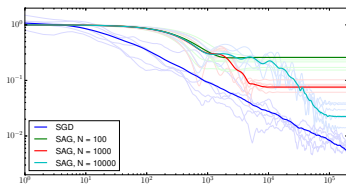


Figure 4: Results for the computation of 595 pairwise word mover's distances between 35 very large corpora of text, each represented as a cloud of $l = 20,000$ word embeddings.

Numerical Results for Density Fitting (Semi-discrete OT)



(a) SGD



(b) SGD vs. SAG

Figure 5: (a) Effect of regularization parameter ϵ (b) Effect of sampling (discrete algo) vs. using semi-discrete algo (blue)

Dual Formulation as an Expectation

Recall the dual objective function to be maximized, for $\varepsilon > 0$

$$F_\varepsilon(\mathbf{u}, \mathbf{v}) = \int_{\mathcal{X}} u(x) d\mu(x) + \int_{\mathcal{Y}} v(y) d\nu(y) - \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right) d\mu(x) d\nu(y)$$

Dual Formulation as an Expectation

Recall the dual objective function to be maximized, for $\varepsilon > 0$

$$F_\varepsilon(\mathbf{u}, \mathbf{v}) = \int_{\mathcal{X}} \mathbf{u}(x) d\mu(x) + \int_{\mathcal{Y}} \mathbf{v}(y) d\nu(y) - \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} \exp\left(\frac{\mathbf{u}(x) + \mathbf{v}(y) - c(x, y)}{\varepsilon}\right) d\mu(x) d\nu(y)$$

Let $X \sim \mu$ and $Y \sim \nu$ be two independent random variables, we get

$$F_\varepsilon(\mathbf{u}, \mathbf{v}) = \mathbb{E}_{\mu \otimes \nu} [f_\varepsilon(X, Y, \mathbf{u}, \mathbf{v})]$$

where $\forall \varepsilon > 0$,

$$f_\varepsilon(x, y, \mathbf{u}, \mathbf{v}) \stackrel{\text{def.}}{=} \mathbf{u}(x) + \mathbf{v}(y) - \varepsilon \exp\left(\frac{\mathbf{u}(x) + \mathbf{v}(y) - c(x, y)}{\varepsilon}\right).$$

Reminder on RKHS I

We consider two reproducing kernel Hilbert spaces (RKHS) \mathcal{H} and \mathcal{G} on \mathcal{X} and on \mathcal{Y} , with kernels κ and ℓ .

Properties of RKHS

(a) if $u \in \mathcal{H}$, then $u(x) = \langle u, \kappa(\cdot, x) \rangle_{\mathcal{H}}$

(b) $\kappa(x, x') = \langle \kappa(\cdot, x), \kappa(\cdot, x') \rangle_{\mathcal{H}}$.

The Gaussian Kernel

For the Gaussian Kernel i.e. $\kappa(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$ the associated RKHS is dense in the space of continuous functions. This means that any continuous function can be approximated by a linear combination of Gaussian Kernels.

Reminder on RKHS II

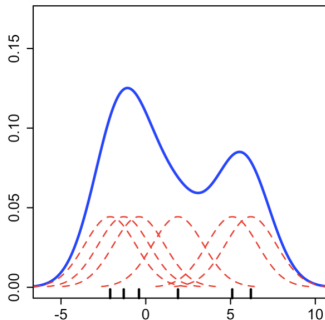


Figure 6: Approximation of a function by a sum of gaussian kernels. The choice of the bandwidth is crucial.

Continuous OT I

$$f_\varepsilon(x, y, \mathbf{u}, \mathbf{v}) \stackrel{\text{def.}}{=} \mathbf{u}(x) + \mathbf{v}(y) - \varepsilon \exp\left(\frac{\mathbf{u}(x) + \mathbf{v}(y) - c(x, y)}{\varepsilon}\right).$$

Rewriting $\mathbf{u}(x)$ and $\mathbf{v}(y)$ as scalar products in \mathcal{H} and \mathcal{G} we get

$$f_\varepsilon(x, y, \mathbf{u}, \mathbf{v}) \stackrel{\text{def.}}{=} \langle \mathbf{u}, \kappa(\cdot, x) \rangle_{\mathcal{H}} + \langle \mathbf{v}, \ell(\cdot, y) \rangle_{\mathcal{G}} - \varepsilon \exp\left(\frac{\langle \mathbf{u}, \kappa(\cdot, x) \rangle_{\mathcal{H}} + \langle \mathbf{v}, \ell(\cdot, y) \rangle_{\mathcal{G}} - c(x, y)}{\varepsilon}\right).$$

we can apply the SGD algorithm in the RKHS :

$$(\mathbf{u}_k, \mathbf{v}_k) \stackrel{\text{def.}}{=} (\mathbf{u}_{k-1}, \mathbf{v}_{k-1}) + \frac{C}{\sqrt{k}} \nabla f_\varepsilon(x_k, y_k, \mathbf{u}_{k-1}, \mathbf{v}_{k-1}) \in \mathcal{H} \times \mathcal{G}, \quad (1)$$

where (x_k, y_k) are i.i.d. samples from $\mu \otimes \nu$.

Continuous OT : Numerical Results

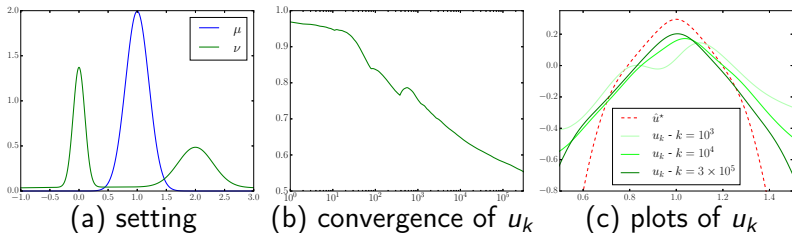


Figure 7: (a) Plot of $\frac{d\mu}{dx}$ (blue) and $\frac{d\nu}{dx}$ (green). (b) Plot of $\|\mathbf{u}_k - \hat{\mathbf{u}}^*\|_2 / \|\hat{\mathbf{u}}^*\|_2$ as a function of k with SGD in the RKHS, for regularized OT using $\varepsilon = 10^{-1}$. (c) Plot of the iterates u_k for $k = 10^3, 10^4, 10^5$ and the proxy for the true potential $\hat{\mathbf{u}}^*$, evaluated on a grid where μ has non negligible mass.

Continuous OT : Theory in progress

We recently proved that the dual potentials are in a Sobolev ball (and thus bounded in a certain RKHS)

- We get convergence of kernel SGD in the continuous setting
- We can use standard results on RKHS to prove regularized OT has sample complexity in $O(n^{-1/2})$, similar to MMD / much better than standard OT

Conclusion

- Dual formulations of OT can be rewritten as expectation maximization problems
- This allows the use of stochastic optimization methods
- Surpass Sinkhorn in the discrete setting (online method more efficient than batch)
- Tackle semi-discrete and continuous problems without requiring discretization