

# Learning Generative Models with Optimal Transport

Aude Genevay

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

Journée YSP - 26 Janvier 2018

*Joint work with M.Cuturi and G. Peyré*

## Recurrent issue in ML : Fitting data to a probabilistic model

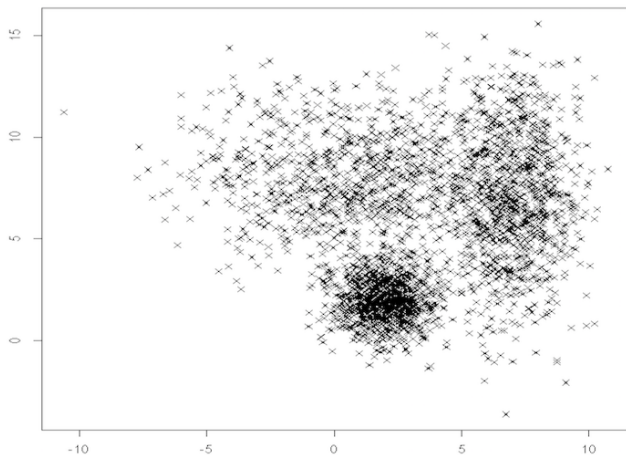


Figure 1: Density Fitting with a Gaussian Mixture

## Recurrent issue in ML : Fitting data to a probabilistic model

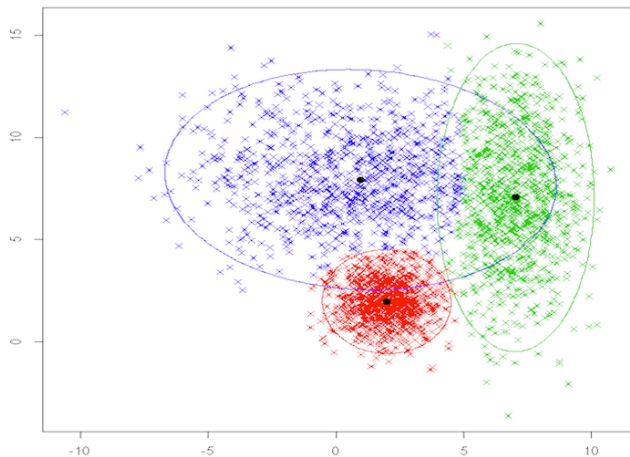


Figure 2: Density Fitting with a Gaussian Mixture

# 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

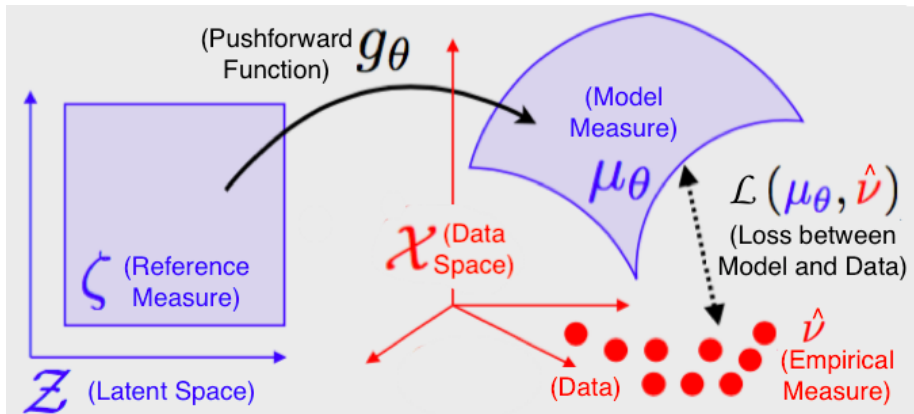
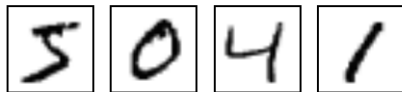


Figure 3: Illustration of Density Fitting on a Generative Model

# Density Fitting for Generative Models I

- Parametric model :  $\mu_\theta = g_{\theta\#}\zeta$
- $\zeta$  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  $z \sim \zeta$
- Very popular topic in ML : image generation



## Density Fitting for Generative Models II

- Generative Models usually supported on low dimensional manifolds ( $\dim \mathcal{Z} < \dim \mathcal{X}$ )
- $\mu_\theta$  doesn't have density wrt Lebesgue measure

⇒ **MLE can't be applied in this context!**

- 2 natural candidates emerge for  $\mathcal{L}$ 
  - ▶ Maximum Mean Discrepancy (based on Reproducing Kernel Hilbert Spaces) → Hilbertian norm
  - ▶ The Wasserstein Distance (based on Optimal Transport) → Non-Hilbertian distance

## Reminders on Maximum Mean Discrepancy I

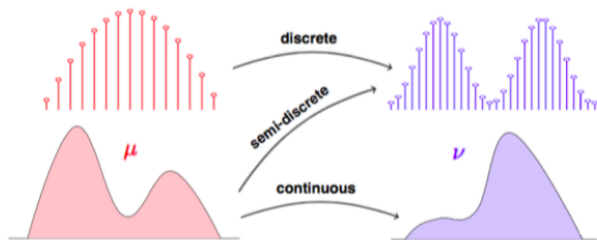
- Consider Reproducing Kernel Hilbert Space  $\mathcal{H}$  with kernel  $k$
- $f \in \mathcal{H} \Rightarrow f(x) = \langle f, k(\cdot, x) \rangle$
- MMD [Gretton et al. '12]:

$$\|\mu - \nu\|_k = \mathbb{E}_{\mu \otimes \mu}[k(x, x')] + \mathbb{E}_{\nu \otimes \nu}[k(y, y')] - 2\mathbb{E}_{\mu \otimes \nu}[k(x, y)]$$

- Usual kernels
  - ▶ Gaussian kernel :  $k(x, y) = \exp\left(\frac{\|x-y\|^2}{\sigma}\right)$
  - ▶ Energy distance kernel :  $k(x, y) = d(x, 0) + d(y, 0) - d(x, y)$



# Reminders on Optimal Transport I



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

## Reminders on Optimal Transport II

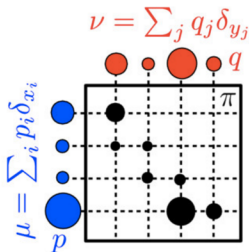
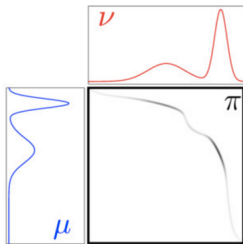
Two positive Radon measures  $\mu$  on  $\mathcal{X}$  and  $\nu$  on  $\mathcal{Y}$  of mass 1

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

Set of couplings with marginals  $\mu$  and  $\nu$

$\Pi(\mu, \nu) \stackrel{\text{def.}}{=} \{ \pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y}) \mid \pi(A \times \mathcal{Y}) = \mu(A), \pi(\mathcal{X} \times B) = \nu(B) \}$

*What's the coupling that minimizes the total cost?*



## Reminders on Optimal Transport III

Main issues of Wasserstein distance :

- Computationally Expensive : need to solve LP (in discrete case)
- Poor Sample Complexity :  $W(\mu, \hat{\mu}_n) \sim n^{-\frac{1}{d}}$ 
  - scales exponentially with dimension (more on that in Francis' talk)
  - need a lot of samples to get a good approximation of  $\mathcal{W}$

# Entropy!

- Basically : Adding an entropic regularization smoothes the constraint
- Makes the problem easier :
  - ▶ yields an unconstrained dual problem
  - ▶ discrete case can be solved efficiently with alternate maximizations on the dual variables : Sinkhorn's algorithm (more on that later)
- For ML applications, regularized Wasserstein is better than standard one
- In high dimension, helps avoiding overfitting

## Entropic Relaxation of OT [Cuturi '13]

Add entropic Penalty to Kantorovitch formulation of OT

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

where

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

Regularized loss :

$$W_{c, \varepsilon}(\mu, \nu) \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi_{\varepsilon}(x, y)$$

where  $\pi_{\varepsilon}$  solution of  $(\mathcal{P}_{\varepsilon})$

# Sinkhorn Divergences : interpolation between OT and MMD

## Theorem

*The Sinkhorn loss between two measures  $\mu, \nu$  is defined as:*

$$\bar{W}_{c,\varepsilon}(\mu, \nu) = 2W_{c,\varepsilon}(\mu, \nu) - W_{c,\varepsilon}(\mu, \mu) - W_{c,\varepsilon}(\nu, \nu)$$

*with the following limiting behavior in  $\varepsilon$ :*

- 1 as  $\varepsilon \rightarrow 0$ ,  $\bar{W}_{c,\varepsilon}(\mu, \nu) \rightarrow 2W_c(\mu, \nu)$
- 2 as  $\varepsilon \rightarrow +\infty$ ,  $\bar{W}_{c,\varepsilon}(\mu, \nu) \rightarrow \|\mu - \nu\|_{-c}$

*where  $\|\cdot\|_{-c}$  is the MMD distance whose kernel is minus the cost from OT.*

**Remark :** Some conditions are required on  $c$  to get MMD distance when  $\varepsilon \rightarrow \infty$ . In particular,  $c = \|\cdot\|_p^p, 0 < p < 2$  is valid.

## Density Fitting with Sinkhorn loss "Formally"

Solve  $\min_{\theta} E(\theta)$

where  $E(\theta) \stackrel{\text{def.}}{=} \bar{W}_{c,\varepsilon}(\mu_{\theta}, \nu)$

$\Rightarrow$  Issue : untractable gradient

## Approximating Sinkhorn loss

- Rather than approximating the gradient approximate the loss itself
- Minibatches :  $\hat{E}(\theta)$ 
  - ▶ sample  $x_1, \dots, x_m$  from  $\mu_\theta$
  - ▶ use empirical Wasserstein distance  $W_{c,\varepsilon}(\hat{\mu}_\theta, \hat{\nu})$  where  $\hat{\mu}_\theta = \frac{1}{N} \sum_{i=1}^m \delta_{x_i}$
- Use  $L$  iterations of Sinkhorn's algorithm :  $\hat{E}^{(L)}(\theta)$ 
  - ▶ compute  $L$  steps of the algorithm
  - ▶ use this as a proxy for  $W(\hat{\mu}_\theta, \nu)$



# Sinkhorn's Algorithm

- State of the art solver for discrete regularized OT
- Two equivalent views
  - ▶ Alternate projections on the constraints of the primal
  - ▶ Alternate minimizations on the dual

- Iterates  $(a, b)$  : 
$$\begin{cases} a \leftarrow \frac{1}{K(b \odot \nu)} \\ b \leftarrow \frac{1}{K^T(a \odot \mu)} \end{cases}$$

where  $K \stackrel{\text{def.}}{=} \exp \frac{-c}{\varepsilon}$  and  $\odot$  is coordinatewise vector multiplication.

- Primal solution  $\pi_\varepsilon = \text{diag}(a)K\text{diag}(b)$
- Linear convergence of the iterates to the optimizers
- Number of iterations needed for convergence increases when  $\varepsilon$  decreases

# Computing the Gradient in Practice

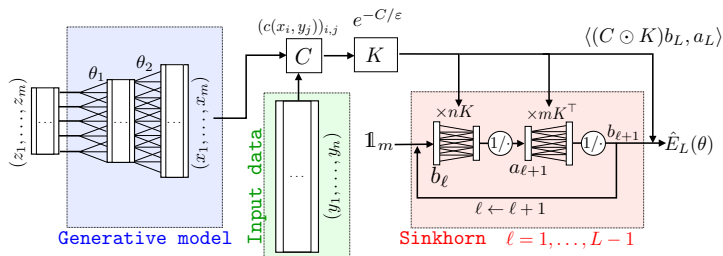
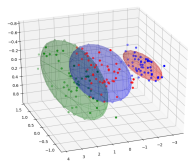


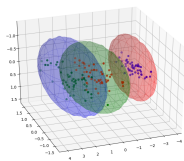
Figure 4: Scheme of the loss approximation

- Compute *exact* gradient of  $\hat{E}^{(L)}(\theta)$  with autodiff
- Backpropagation through above graph
- Same computational cost as evaluation of  $\hat{E}^{(L)}(\theta)$

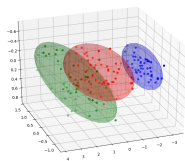
## Numerical Results : a toy example



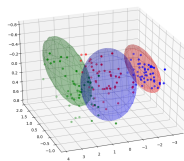
(a) MMD



(b)  $\varepsilon = 1$



(c)  $\varepsilon = 0.1$



(d)  $\varepsilon = 0.01$

Figure 5: Ellipses after convergence of the stochastic gradient descent with  $L = 20$ ,  $m = 200$

## Numerical Results on MNIST (L2 cost)



Figure 6: Samples from MNIST dataset

## Numerical Results on MNIST (L2 cost)

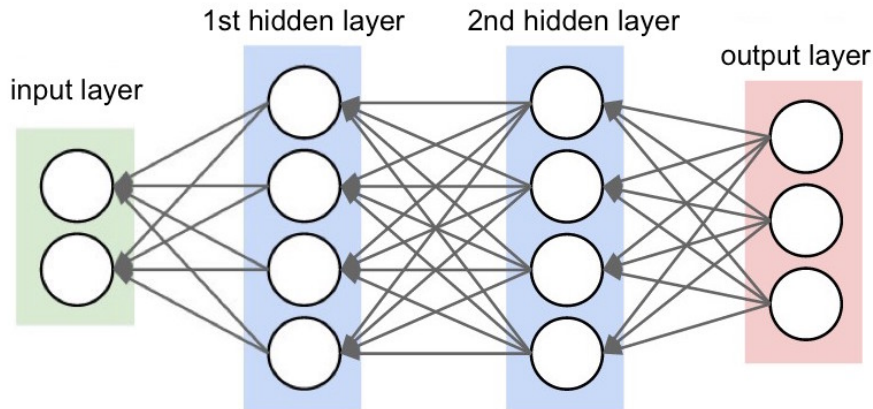
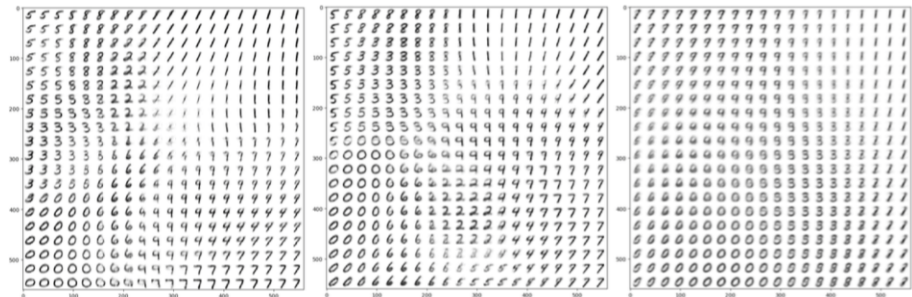


Figure 7: Fully connected NN with 2 hidden layers

# Numerical Results on MNIST (L2 cost)



(a)  $\epsilon = 1, m = 200, L = 10$       (b)  $\epsilon = 10^{-1}, m = 200, L = 100$       (c)  $\epsilon = 10^{-1}, m = 10, L = 300$

Figure 8: Manifolds in the latent space for various parameters

## Learning the cost [Li et al. '17, Bellemare et al. '17]

- On complex data sets, choice of a good ground metric  $c$  is not trivial
- Use parametric cost function  $c_\phi(x, y) = \|f_\phi(x) - f_\phi(y)\|_2^2$   
(where  $f_\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ )
- Optimization problem becomes minmax (like GANs)

$$\min_{\theta} \max_{\phi} \bar{W}_{c_{\phi}, \varepsilon}(\mu_{\theta}, \nu)$$

- Same approximations but alternate between updating the cost parameters  $\phi$  and the measure parameters  $\theta$

## Numerical Results on CIFAR (learning the cost)



Figure 9: Samples from CIFAR dataset



# Numerical Results on CIFAR (learning the cost)

## Deep convolutional GANs (DCGAN) [1511.06434]

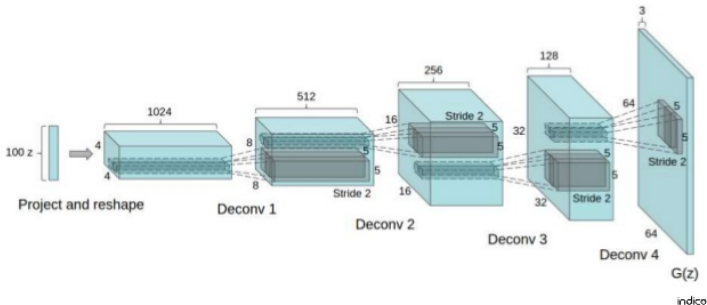
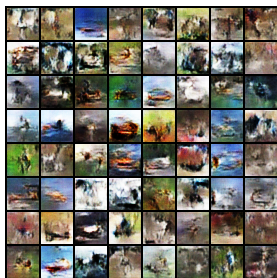


Figure 10: Fully connected NN with 2 hidden layers

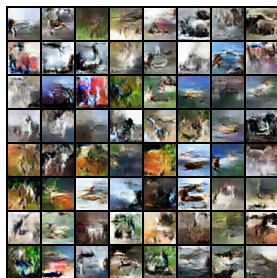
## Numerical Results on CIFAR (learning the cost)



(a) MMD



(b)  $\varepsilon = 1000$



(c)  $\varepsilon = 10$

Figure 11: Samples from the generator trained on CIFAR 10 for MMD and Sinkhorn loss (coming from the same samples in the latent space)

Which is better? Not just about generating nice images, but more about capturing a high dimensional distribution... Hard to evaluate.

# Conclusion

- Take Home message : Sinkhorn Divergences allow to interpolate between OT and MMD
- Future Work :
  - ▶ Theory of Sinkhorn Divergences (positivity / sample complexity)
  - ▶ Evaluation of generative models to study optimal choice of epsilon