

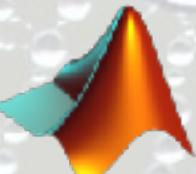
Numerical Optimal Transport

<http://optimaltransport.github.io>

Density Fitting

Gabriel Peyré

www.numerical-tours.com



ENS

ÉCOLE NORMALE
SUPÉRIEURE

Weak vs Strong Topology

Random vectors

$$\mathbb{P}(\textcolor{red}{X} \in A)$$

Convergence in law:

\forall set A

$$\mathbb{P}(\textcolor{red}{X}_n \in A) \xrightarrow{n \rightarrow +\infty} \mathbb{P}(\textcolor{red}{X} \in A)$$

Radon measures

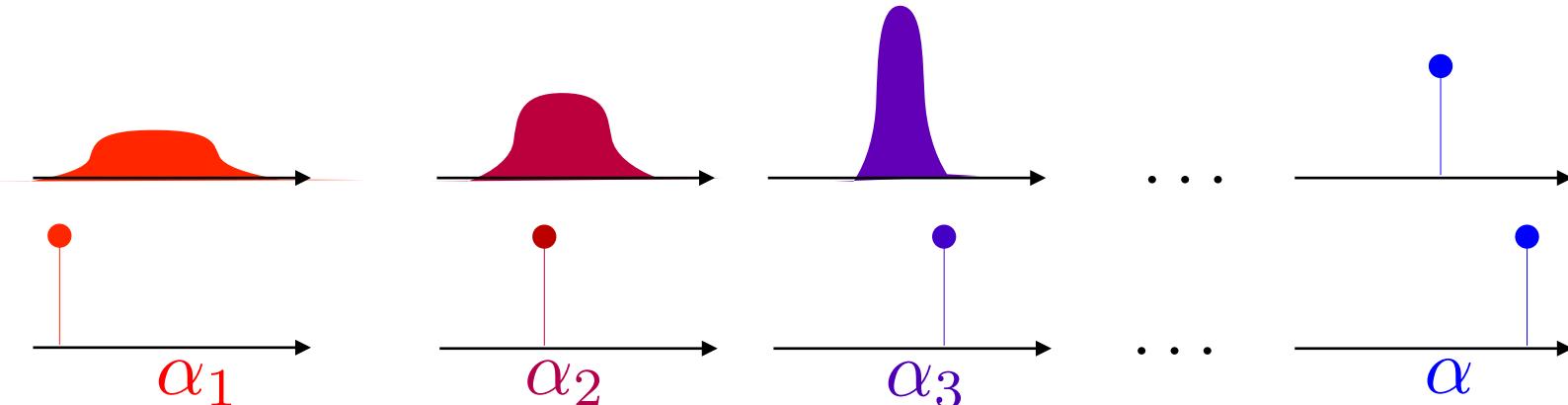
$$\int_A d\alpha(x)$$

Weak* convergence:

\forall continuous function f

$$\int f d\alpha_n \xrightarrow{n \rightarrow +\infty} \int f d\alpha$$

Weak convergence:



Overview

- **Csiszar Divergences**
- Dual Norms and MMD
- Minimum Kantorovitch Estimators
- Deep Generative Models Fitting

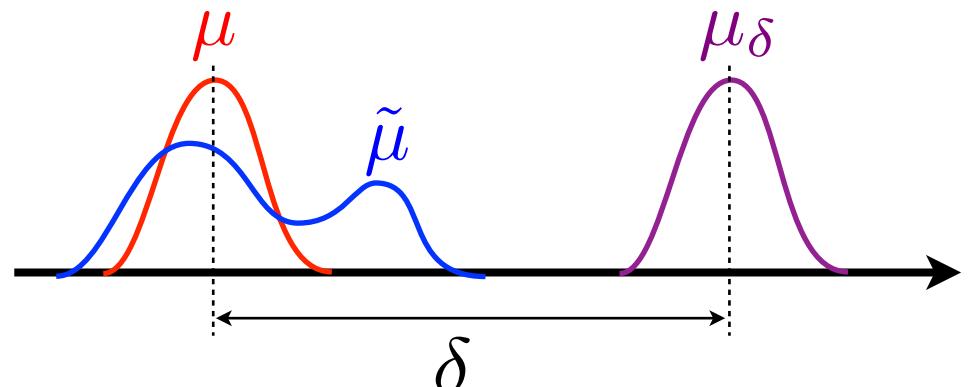
Metrics on the Space of Measures

$$d\mu(x) = \rho(x)dx$$

$$d\tilde{\mu}(x) = \tilde{\rho}(x)dx$$

$$d\mu_\delta(x) = \rho(x - \delta)dx$$

Bins-to-bins metrics:



Kullback-Leibler divergence:

$$D_{\text{KL}}(\mu, \tilde{\mu}) = \int \rho(x) \log \frac{\rho(x)}{\tilde{\rho}(x)} dx$$

Hellinger distance:

$$D_{\text{H}}(\mu, \tilde{\mu})^2 = \int \left(\sqrt{\rho(x)} - \sqrt{\tilde{\rho}(x)} \right)^2 dx$$

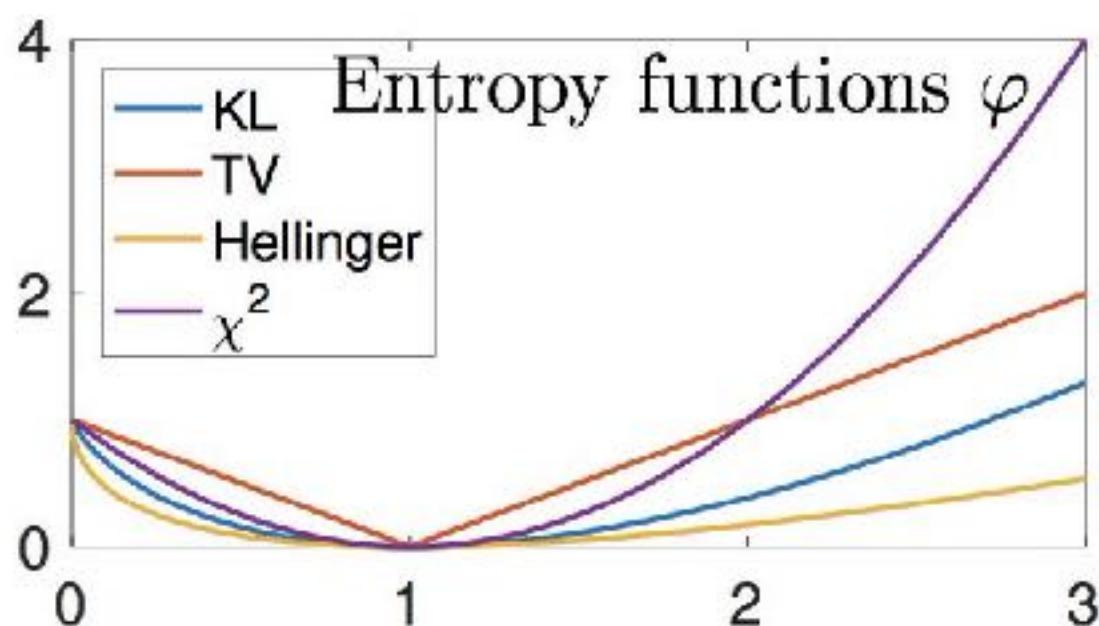
Effect of translation:

$$D(\mu, \mu_\delta) \approx \text{cst}$$

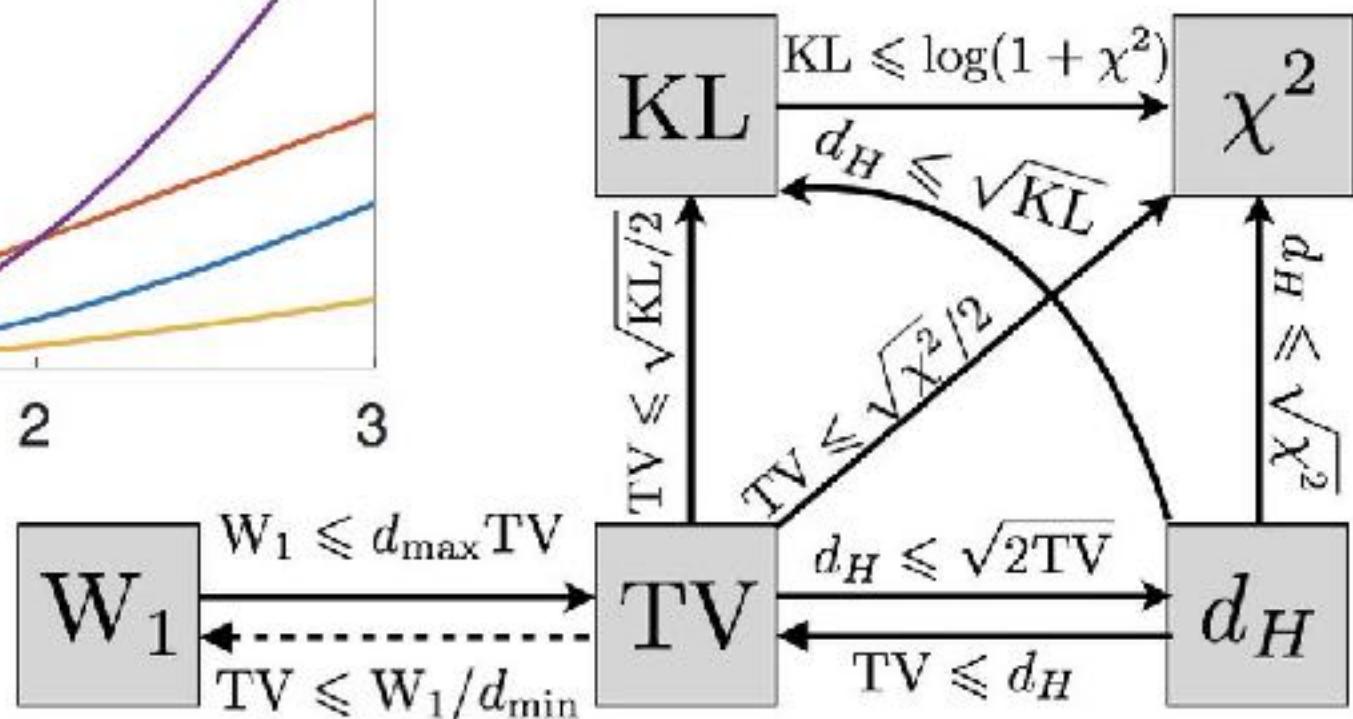
$$W_2(\mu, \mu_\delta) = \delta$$

Csiszar Divergence

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right) d\beta + \varphi'_\infty \alpha^\perp(\mathcal{X})$$



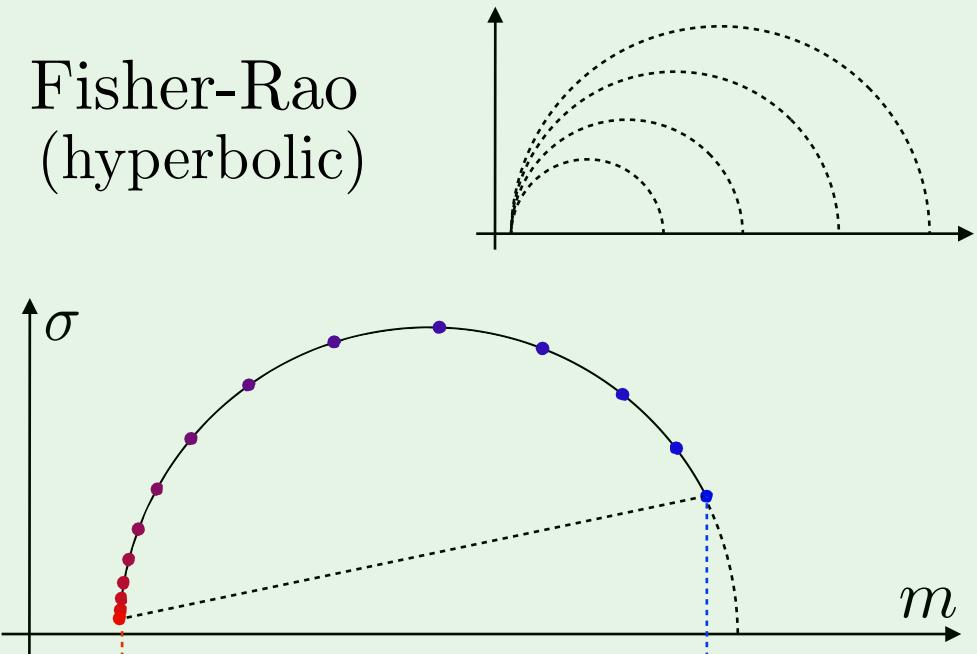
$$\varphi'_\infty = \lim_{x \uparrow +\infty} \varphi(x)/x \in \mathbb{R} \cup \{\infty\}$$



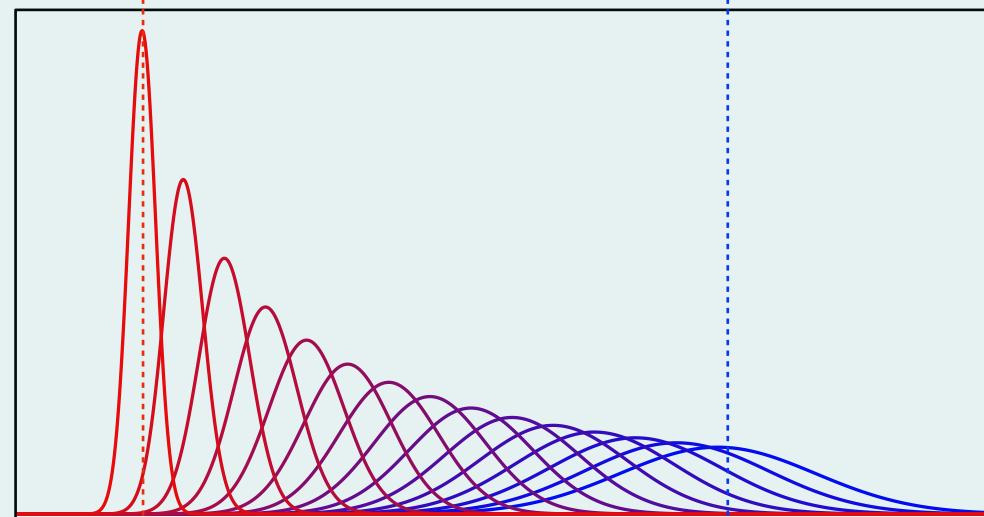
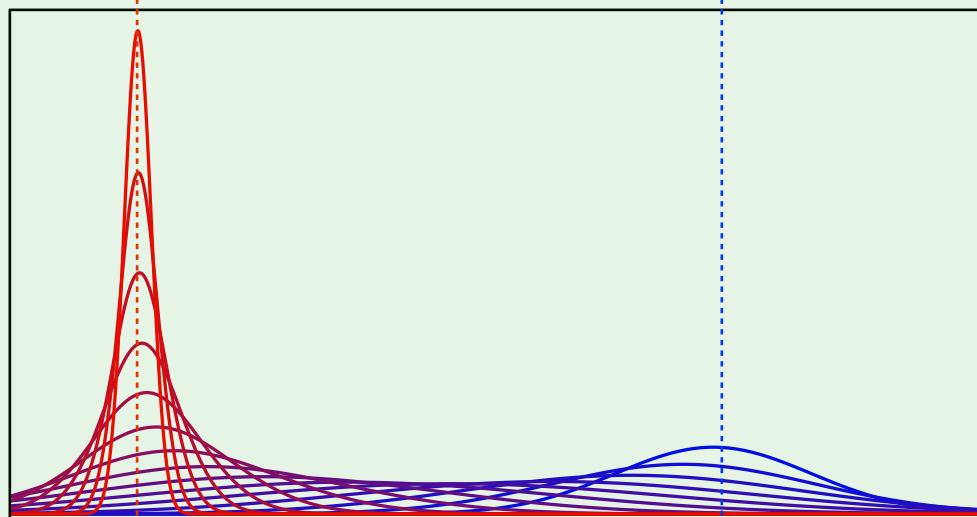
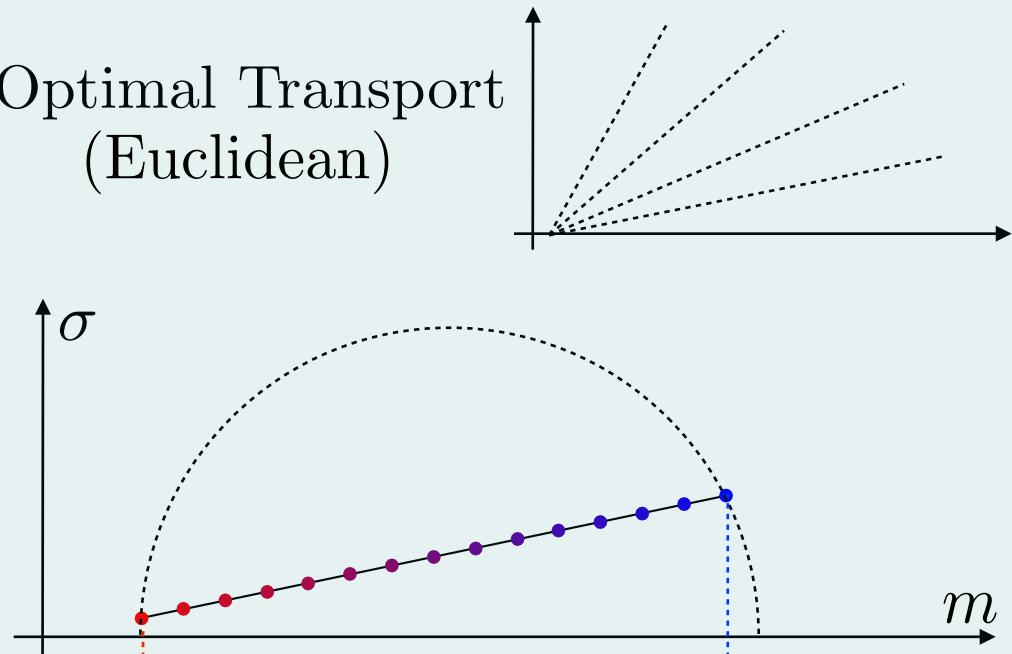
Csiszár divergences, a unifying way to define losses between arbitrary positive measures (discrete & densities). [https://en.wikipedia.org/wiki/F-divergence ...](https://en.wikipedia.org/wiki/F-divergence)

OT vs. KL (Fisher-Rao)

Fisher-Rao
(hyperbolic)



Optimal Transport
(Euclidean)



Overview

- Csiszar Divergences
- Dual Norms and MMD
- Minimum Kantorovitch Estimators
- Deep Generative Models Fitting

Dual Norms

Dual norms: (aka Integral Probability Metrics)

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max \left\{ \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x)) ; f \in B \right\}$$

Wasserstein 1: $B = \{f ; \|\nabla f\|_\infty \leq 1\}$.

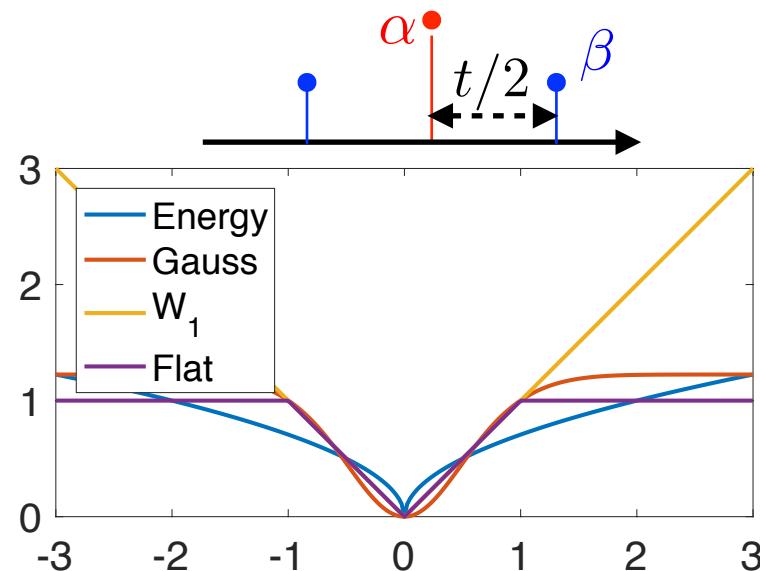
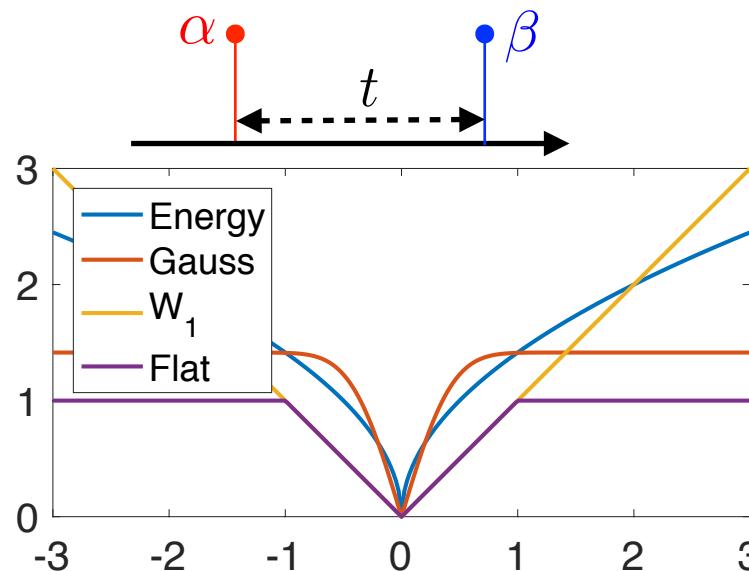
Flat norm: $B = \{f ; \|f\|_\infty \leq 1, \|\nabla f\|_\infty \leq 1\}$.

RKHS: $B = \{f ; \|f\|_k^2 \leq 1\}$.

$$\|\alpha - \beta\|_B^2 = \int k(x, x') d\alpha(x) d\alpha(x') + \int k(x, x') d\beta(y) d\beta(y') - 2 \int k(x, y) d\alpha(x) d\beta(y)$$

Energy distance: $k(x, y) = -\frac{\|x - y\|^2}{2}$

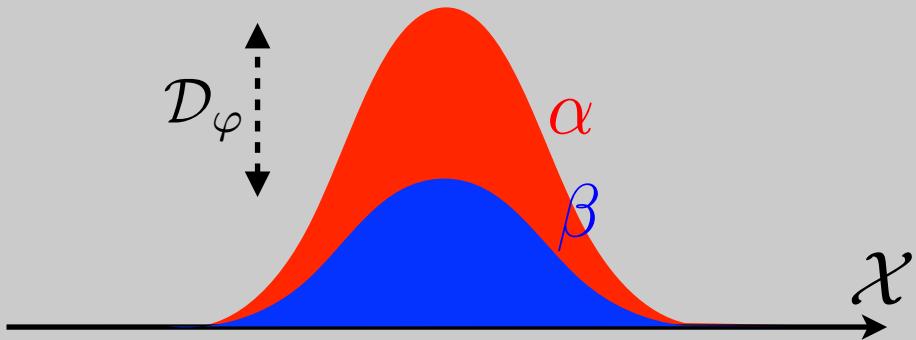
Gaussian: $k(x, y) = e^{-\frac{\|x - y\|^2}{2\sigma^2}}$



Csiszar Divergence vs Dual Norms

Csiszár divergences:

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi \left(\frac{d\alpha}{d\beta} \right) d\beta$$

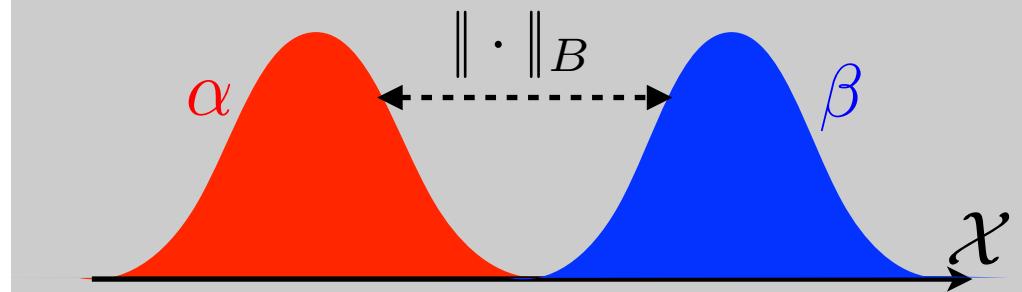


Strong topology

→ KL, TV, χ^2 , Hellinger ...

Dual norms:

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max_{f \in B} \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x))$$



Weak topology

→ W_1 , flat, RKHS*, energy dist, ...

RKHS Norms aka Maximum Mean Discrepancy

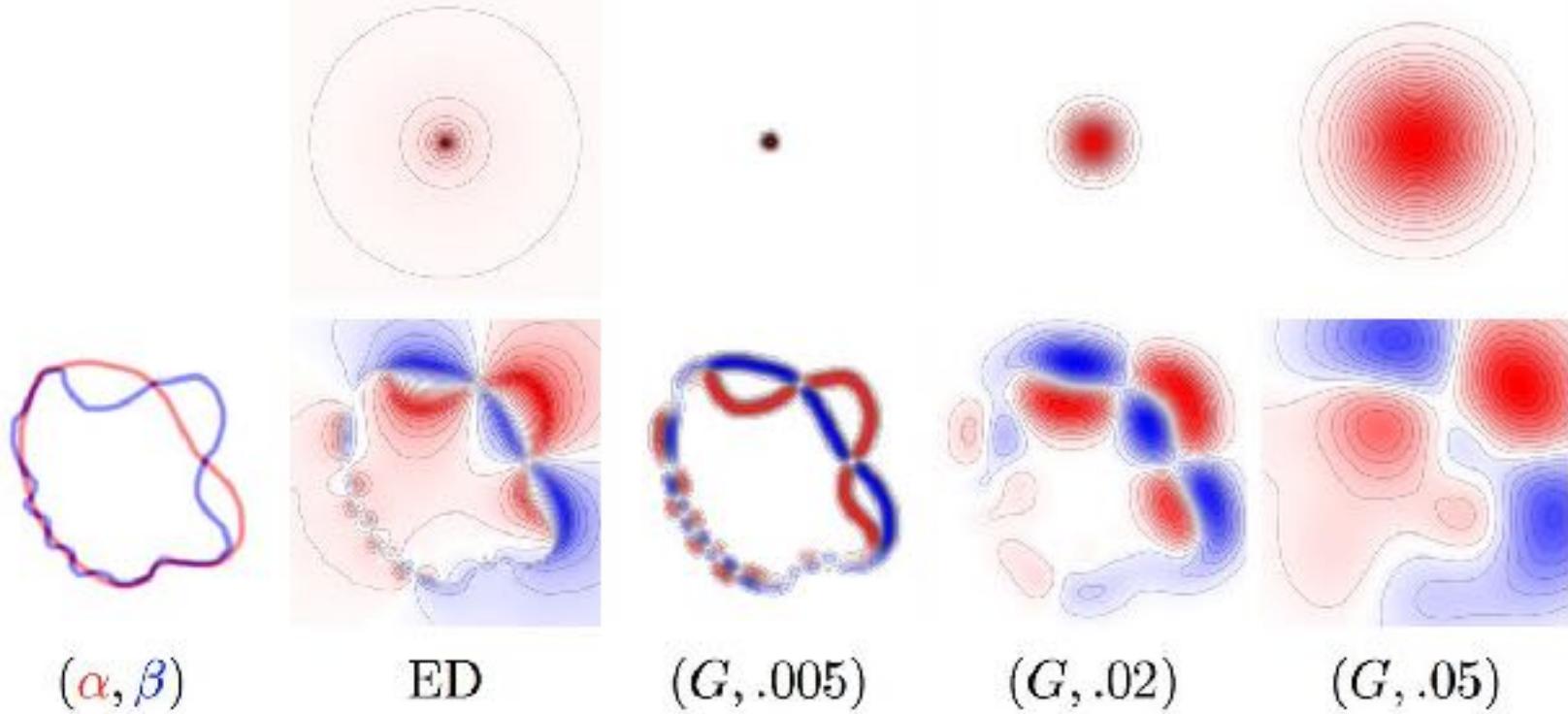


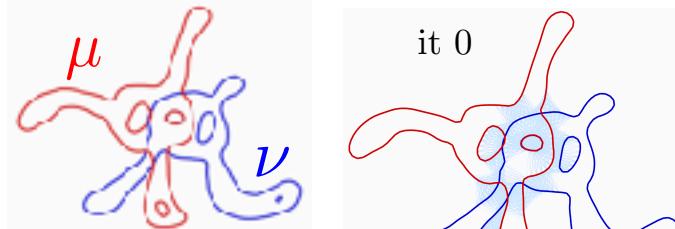
Figure 8.4: Top row: display of ψ such that $\|\alpha - \beta\|_k = \|\psi \star (\alpha - \beta)\|_{L^2(\mathbb{R}^2)}$, formally defined over Fourier as $\hat{\psi}(\omega) = \sqrt{\hat{\varphi}(\omega)}$ where $k^*(x, x') = \varphi(x - x')$. Bottom row: display of $\psi \star (\alpha - \beta)$. (G, σ) stands for Gaussian kernel of variance σ^2 and ED for Energy Distance kernel (in which case $\psi(x) = 1/\sqrt{\|x\|}$).

OT Loss for Diffeomorphic Registration

Joint work with J. Feydy, B. Charier, F-X. Vialard.

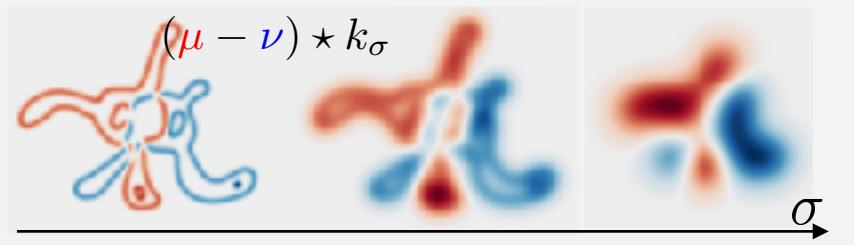
Shape registration: $\min_{\varphi \text{ diffeo}} D(\varphi(\mu), \nu) + R(\varphi)$

loss regularity



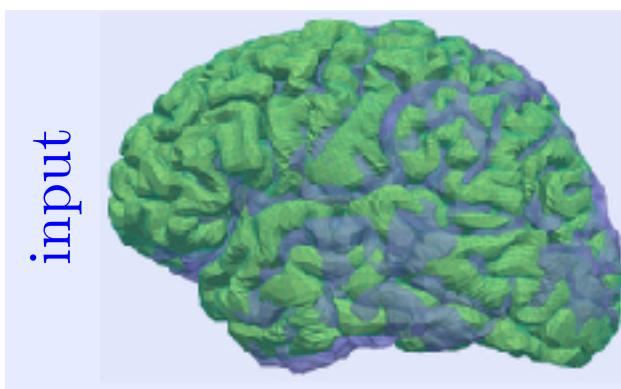
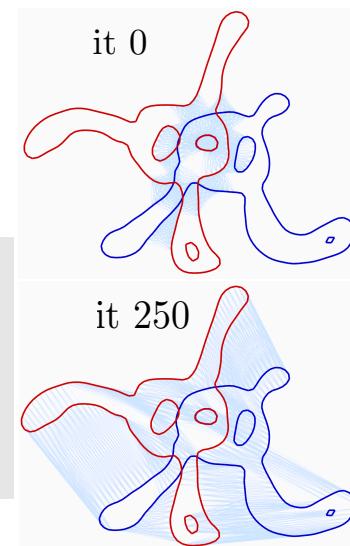
Hilbertian loss (MMD/RKHS):

$$D(\mu, \nu) = \|k_\sigma \star (\mu - \nu)\|_{L^2}^2$$

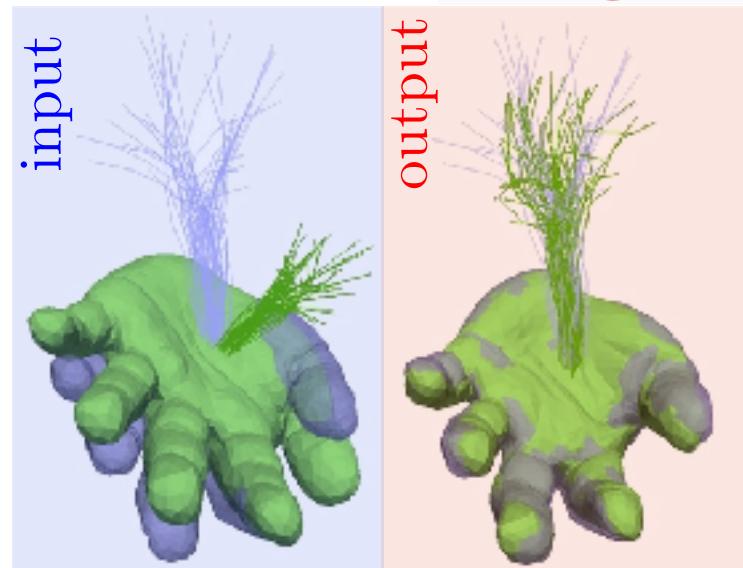
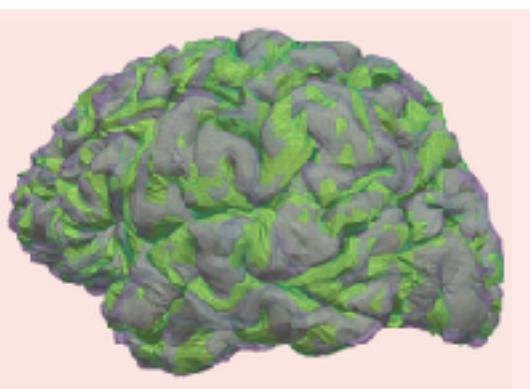


Sinkhorn divergence:

$$D(\mu, \nu) = \bar{W}_\varepsilon(\mu, \nu)$$



output



- Do not use OT for registration ... but as a loss.
- Sinkhorn's iterates “propagate” a small bandwidth kernel.
- Automatic differentiation: game changer for advanced loss and models.

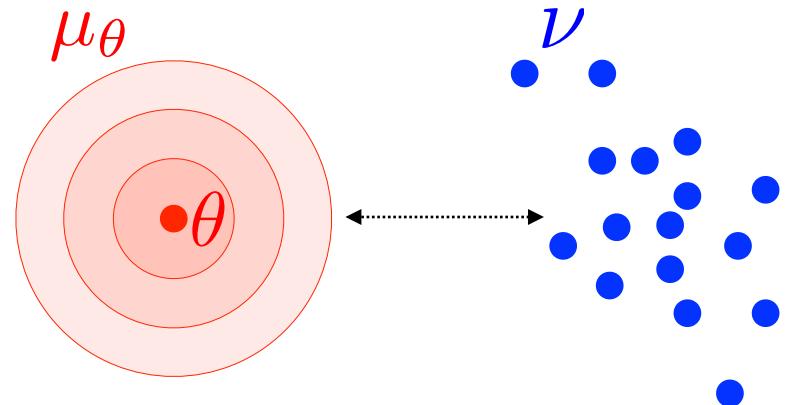
Overview

- Csiszar Divergences
- Dual Norms and MMD
- **Minimum Kantorovitch Estimators**
- Deep Generative Models Fitting

Density Fitting and Generative Models

Observations: $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model: $\theta \mapsto \mu_\theta$



Density fitting: $d\mu_\theta(y) = f_\theta(y)dy$

$$\min_{\theta} \widehat{\text{KL}}(\nu | \mu_\theta) \stackrel{\text{def.}}{=} - \sum_j \log(f_\theta(y_j))$$

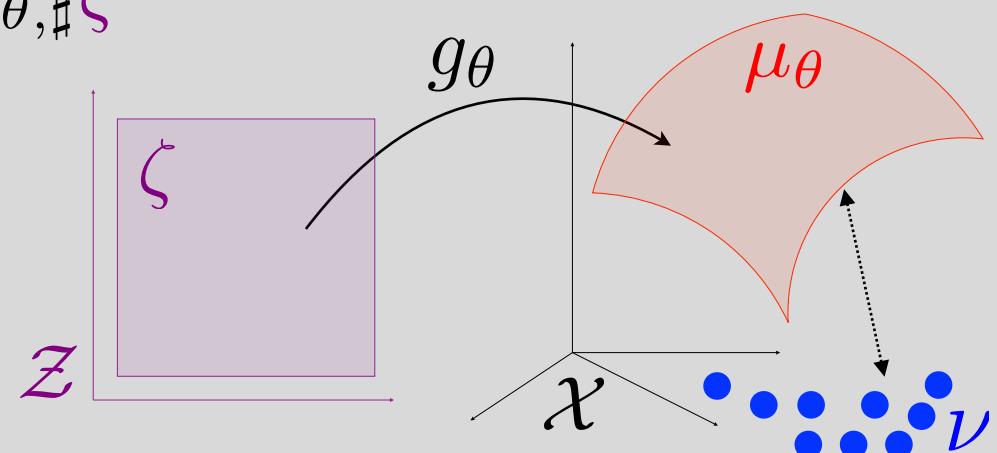
Maximum likelihood (MLE)

Generative model fit: $\mu_\theta = g_{\theta, \sharp} \zeta$

$$\widehat{\text{KL}}(\nu | \mu_\theta) = +\infty$$

→ MLE undefined.

→ Need a weaker metric.



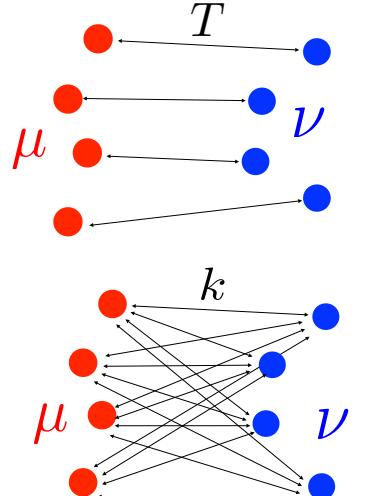
Loss Functions for Measures

Density fitting: $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

Optimal Transport Distances

$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



Maximum Mean Discrepancy (MMD)

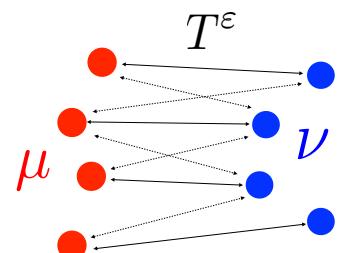
$$\|\mu - \nu\|_k^2 \stackrel{\text{def.}}{=} \frac{1}{N^2} \sum_{i,i'} k(x_i, x_{i'}) + \frac{1}{P^2} \sum_{j,j'} k(y_j, y_{j'}) - \frac{2}{NP} \sum_{i,j} k(x_i, y_j)$$

Gaussian: $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$. Energy distance: $k(x, y) = -\|x - y\|^2$.

Sinkhorn divergences [Cuturi 13]

$$W_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} \sum_{i,j} T_{i,j}^{\varepsilon} \|x_i - y_j\|^p$$

$$\bar{W}_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} W_{\varepsilon}(\mu, \nu)^p - \frac{1}{2} W_{\varepsilon}(\mu, \mu)^p - \frac{1}{2} W_{\varepsilon}(\nu, \nu)^p$$



Theorem: [Ramdas, G.Trillos, Cuturi 17]

$$\bar{W}_{\varepsilon}(\mu, \nu)^p \xrightarrow{\varepsilon \rightarrow 0} W(\mu, \nu)^p \quad \xrightarrow{\varepsilon \rightarrow +\infty} \frac{W(\mu, \nu)^p}{\|\mu - \nu\|_k^2}$$

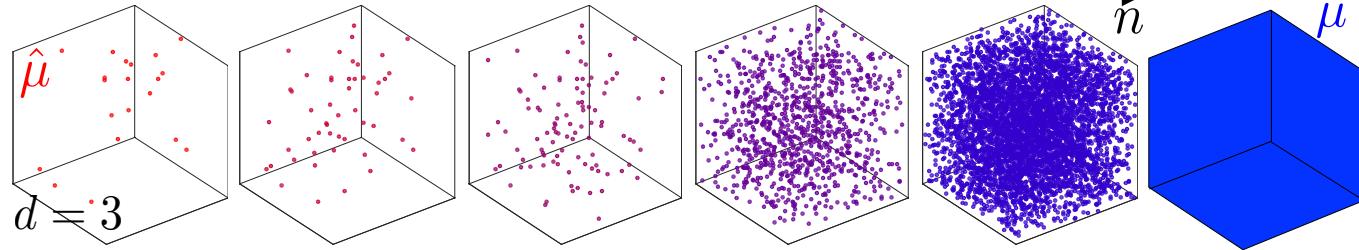
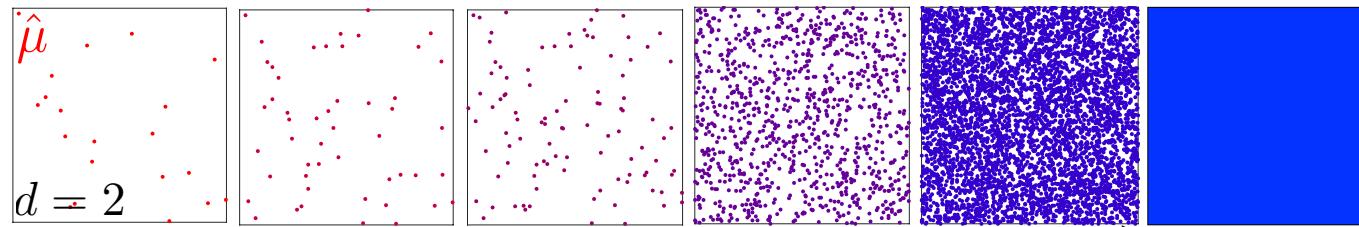
for $k(x, y) = -\|x - y\|^p$

Best of both worlds:

→ cross-validate ε

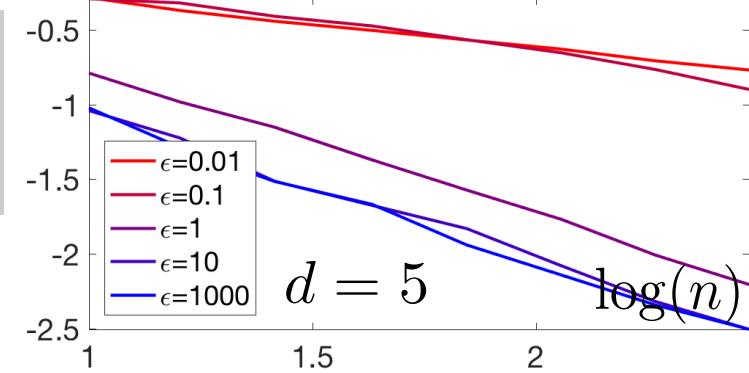
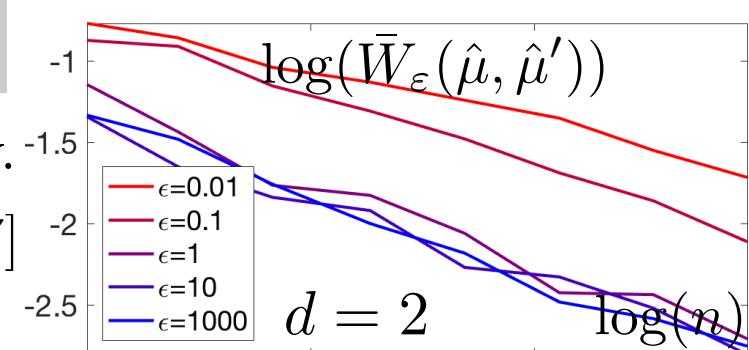
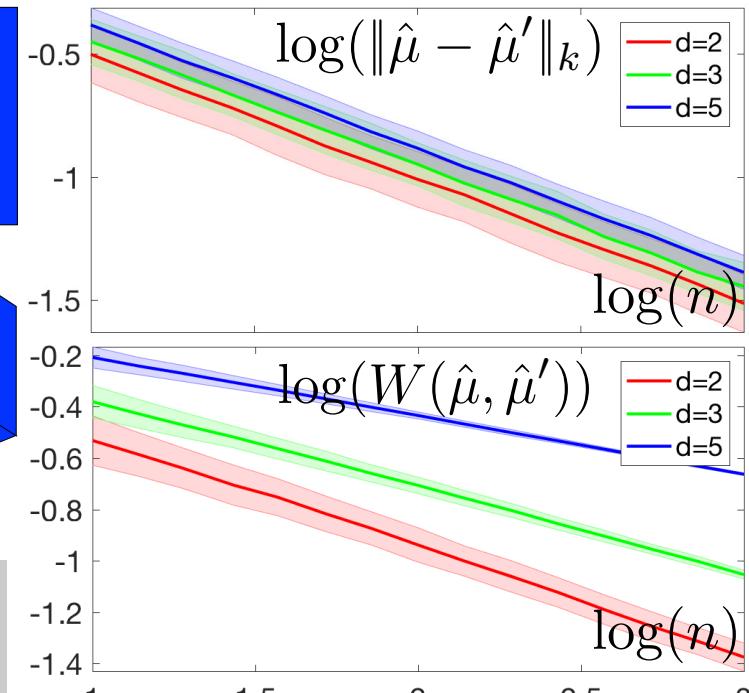
- Scale free (no σ , no heavy tail kernel).
- Non-Euclidean, arbitrary ground distance.
- Less biased gradient.
- No curse of dimension (low sample complexity).

Sample Complexity



$$\text{Theorem: } \mathbb{E}(|W(\hat{\mu}, \hat{\nu}) - W(\mu, \nu)|) = O(n^{-\frac{1}{d}})$$

$$\mathbb{E}(|\|\hat{\mu} - \hat{\nu}\|_k - \|\mu - \nu\|_k|) = O(n^{-\frac{1}{2}})$$



Open problem: sample complexity of \bar{W}_ϵ ?

Overview

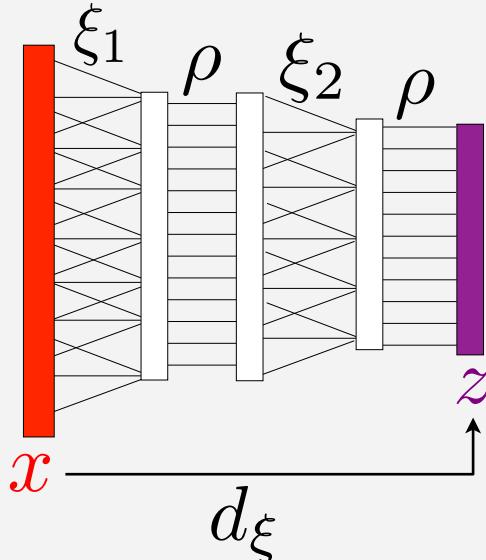
- Csiszar Divergences
- Dual Norms and MMD
- Minimum Kantorovitch Estimators
- Deep Generative Models Fitting

Deep Discriminative vs Generative Models

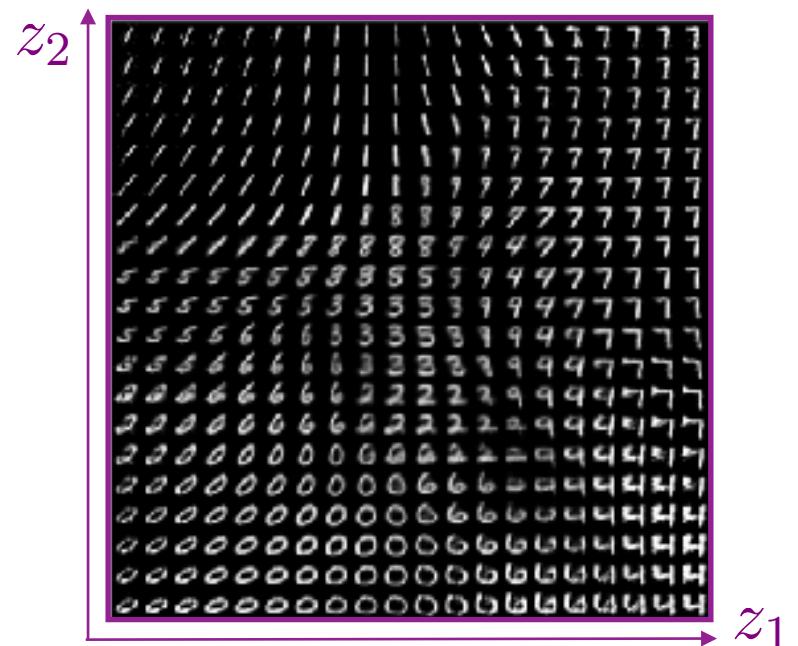
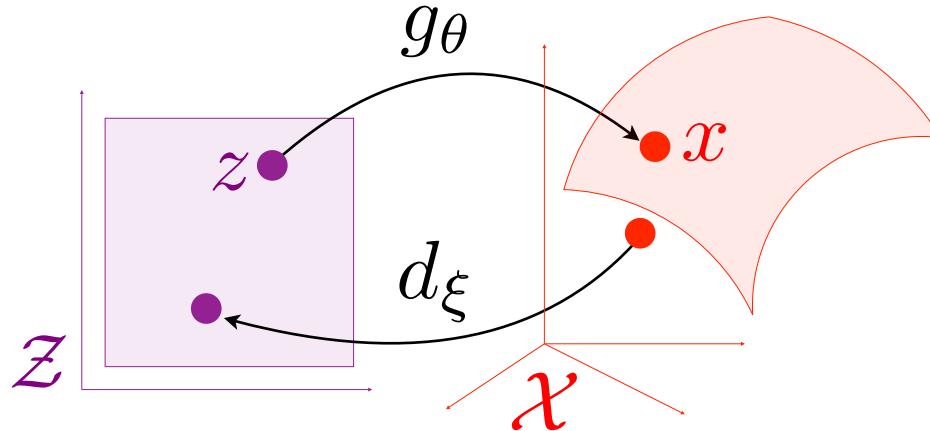
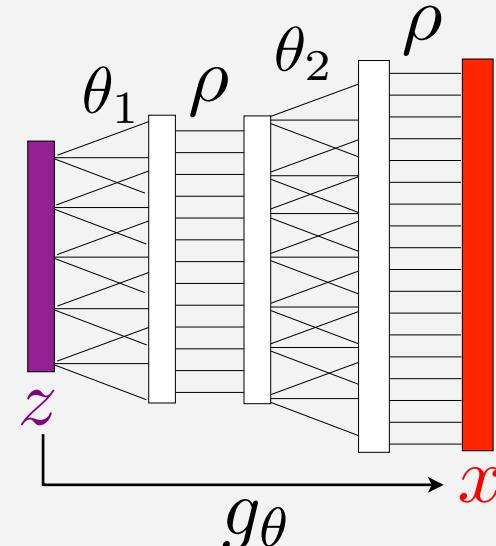
Deep networks:

$$d_\xi(\mathbf{x}) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(\mathbf{x}) \dots)$$
$$g_\theta(\mathbf{z}) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(\mathbf{z}) \dots)$$

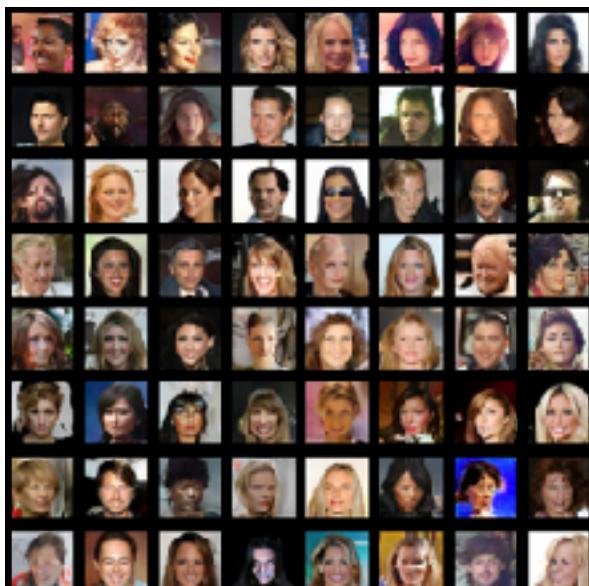
Discriminative



Generative



Examples of Image Generation



Inputs



Small ε



Large ε

→ Need to learn the metric $c(x, y) = \|d_\xi(x) - d_\xi(y)\|^p$ (\sim GANs)

→ Performance evaluation of generative models is an open problem.



Progressive Growing of GANs for Improved Quality, Stability, and Variation
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018