

# Generalized Sliced-Wasserstein Distances

A new family of optimal transport metrics.

Neurips 2019

---

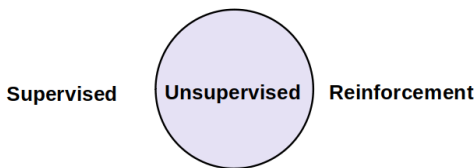
Soheil Kolouri<sup>1</sup>, Kimia Nadjahi<sup>2</sup>, Umut Şimşekli<sup>2</sup>,  
Roland Badeau<sup>2</sup>, Gustavo K. Rohde<sup>3</sup>

<sup>1</sup> HRL Laboratories, <sup>2</sup> Télécom Paris, <sup>3</sup> University of Virginia

# Generative Models

---

# Unsupervised Learning



Learning useful information from a dataset of *unlabeled* samples

*"Humans build a model of the world through predictive unsupervised learning."* – Yann Le Cun

*"One of the main challenges for AI remains unsupervised learning, at which humans are much better than machines."* – Yoshua Bengio

*"What I cannot create, I cannot understand."* – Richard Feynman

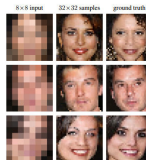
# The Generative Approach

One of the most promising approaches toward unsupervised learning.  
Learns the data distribution in an unsupervised manner.

Why?



Missing data  
NVIDIA, 2018

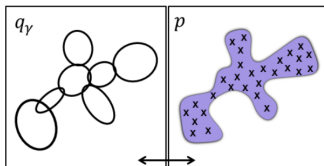


Super resolution  
Dahl et al., CVPR 2017

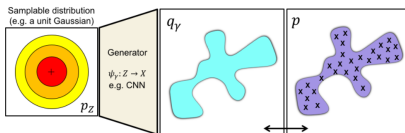


Domain transfer  
Zhu et al., ICCV 2017

# Examples of Generative Models



Gaussian Mixture Models



Deep Generative Modeling (e.g., GANs)

Measuring the *dissimilarity between probability distributions* is at the heart of generative models.

# **Distances between Probability Measures**

---

## Examples of distances

Let  $\mu$  and  $\nu$  be probability measures with finite  $p$ 'th moment, defined on  $X \subset \mathbb{R}^d$ , i.e.,

$$\text{for some } x_0 \in X, \int_X \|x - x_0\|^p d\mu(x) < +\infty$$

with corresponding densities  $l_\mu$  and  $l_\nu$ .

**How “far” is  $\mu$  from  $\nu$ ?**

$$L_p\text{-metrics, } p \geq 1: \quad \left( \int_X |l_\mu(x) - l_\nu(x)|^p dx \right)^{1/p}$$

$$\text{Hellinger distance:} \quad \left( \frac{1}{2} \int_X (\sqrt{l_\mu(x)} - \sqrt{l_\nu(x)})^2 dx \right)^{1/2}$$

$$\text{Kullback-Leibler divergence:} \quad \int_X l_\mu(x) \log \left( \frac{l_\mu(x)}{l_\nu(x)} \right) dx$$

## Examples of distances

Let  $\mu$  and  $\nu$  be probability measures with finite  $p$ 'th moment, defined on  $X \subset \mathbb{R}^d$ . Suppose  $X$  is endowed with a distance  $d$ .

**Wasserstein distance of order  $p$ .**

$$\mathbf{W}_p(\mu, \nu) = \left( \inf_{\gamma \in \Gamma(\mu, \nu)} \int_{X \times X} d(x, y)^p \, d\gamma(x, y) \right)^{1/p},$$

where  $\Gamma(\mu, \nu)$  is the set of probability measures on  $X \times X$  with marginals  $\mu$  (resp.  $\nu$ ) for the first (resp. second) variable.



# Link with Optimal Transport

A solution  $\gamma$  is called the *transport plan* between  $\mu$  and  $\nu$ .

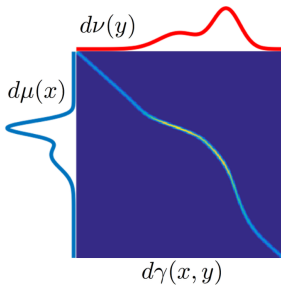
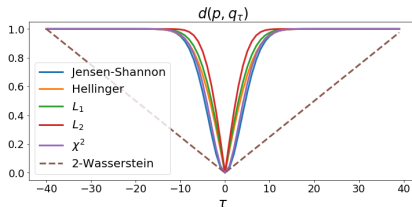
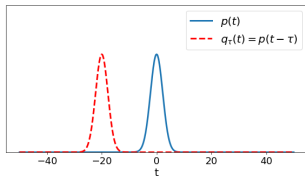


Illustration of optimal transport between two probability measures in the continuous case.

# Comparison of distances



$$\frac{\partial d(p, q_\tau)}{\partial \tau} \approx 0 \text{ whenever } p \text{ and } q_\tau \text{ are non-overlapping.}$$

The Wasserstein distance captures the underlying geometry of the space and is suitable for learning.

# Limitations of Optimal Transport

## **Computing the Wasserstein distance is expensive.**

OT amounts to solving a large-scale linear program.

Computational complexity:  $O(n^3 \log(n))$ , where  $n$  is the number of data samples.

## **Curse of dimensionality.**

The error made when approximating the Wasserstein distance from samples grows *exponentially fast* with the dimension of the space.

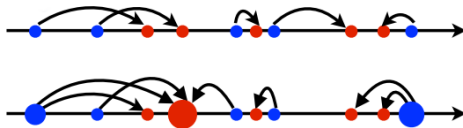
# Sliced-Wasserstein Distance

## Wasserstein distance in 1D.

If  $X \subset \mathbb{R}$ ,

$$W_p(\mu, \nu) = \left( \int_0^1 d(F_\mu^{-1}(t), F_\nu^{-1}(t))^p dt \right)^{1/p},$$

where  $F_\mu^{-1}$  and  $F_\nu^{-1}$  denote the quantile functions of  $\mu$  and  $\nu$  resp.



For empirical distributions: simply compute the average distance between the sorted samples. (Source: Computational Optimal Transport, G. Peyré and M. Cuturi)

# Radon transform

Consider a function  $h \in L^1(\mathbb{R}^d)$  where

$$L^1(\mathbb{R}^d) = \{h : \mathbb{R}^d \rightarrow \mathbb{R} \mid \int_{\mathbb{R}^d} |h(x)| dx < \infty\}.$$

The **standard Radon transform**,  $\mathcal{R}$ , maps  $h$  to the infinite set of its *integrals over the hyperplanes* of  $\mathbb{R}^d$ :

$$\forall (t, \theta) \in \mathbb{R} \times \mathbb{S}^{d-1}, \mathcal{R}h(t, \theta) = \int_{\mathbb{R}^d} h(x) \delta(t - \langle x, \theta \rangle) dx,$$

where,

$\mathbb{S}^{d-1} \subset \mathbb{R}^d$  : the  $d$ -dimensional unit sphere,

$\delta(\cdot)$  : the Dirac delta function in 1D,

$\langle \cdot, \cdot \rangle$  : the Euclidean inner-product.

Each hyperplane can be written as:

$$H(t, \theta) = \{x \in \mathbb{R}^d \mid \langle x, \theta \rangle = t\},$$

(A level set of  $g(x, \theta) = \langle x, \theta \rangle$ )

Fix  $\theta \in \mathbb{S}^{d-1}$ . The integrals over all hyperplanes orthogonal to  $\theta$  define a continuous function  $\mathcal{R}h(\cdot, \theta) : \mathbb{R} \rightarrow \mathbb{R}$  which is a “slice” of  $h$ .

# Sliced-Wasserstein distance

## **Motivation.**

Define an alternative optimal transport distance which does not suffer from a computational burden.

## **Idea.**

Utilize the closed-form formula of the Wasserstein distance in 1D:

1. Obtain representations in 1D for a higher-dimensional distribution through linear projections (via the Radon transform),
2. Calculate the distance between two input distributions as a functional on the Wasserstein distance of their 1D representations.

## Sliced-Wasserstein distance

Let  $\mu$  and  $\nu$  be probability measures with finite  $p$ 'th moment, defined on  $X \subset \mathbb{R}^d$ , with corresponding densities  $l_\mu$  and  $l_\nu$ .

**Sliced-Wasserstein distance of order  $p$ .**

$$\mathbf{SW}_p(\mu, \nu) = \left( \int_{\mathbb{S}^{d-1}} \mathbf{W}_p^p(\mathcal{R}l_\mu(\cdot, \theta), \mathcal{R}l_\nu(\cdot, \theta)) d\theta \right)^{1/p}$$

The SW-distance is, indeed, a distance [1, 3]



## Limitation of Sliced-Wasserstein distance

In practice, we use a Monte-Carlo approximation where samples  $\{\theta_I\}_I$  are uniformly drawn on  $\mathbb{S}^{d-1}$ .

But in high-dimension, there is a high chance that, for any randomly sampled  $\theta$ ,  $\mathbf{W}_p(\mathcal{R}I_\mu(\cdot, \theta), \mathcal{R}I_\nu(\cdot, \theta)) \approx 0$ .

# Maximum Sliced-Wasserstein distance

**Sliced-Wasserstein distance of order  $p$ .**

$$\mathbf{SW}_p(\mu, \nu) = \left( \int_{\mathbb{S}^{d-1}} \mathbf{W}_p^p(\mathcal{R}l_\mu(\cdot, \theta), \mathcal{R}l_\nu(\cdot, \theta)) d\theta \right)^{1/p}$$

**Maximum Sliced-Wasserstein distance of order  $p$ . [2]**

$$\max\text{-}\mathbf{SW}_p(\mu, \nu) = \max_{\theta \in \mathbb{S}^{d-1}} \mathbf{W}_p(\mathcal{R}l_\mu(\cdot, \theta), \mathcal{R}l_\nu(\cdot, \theta))$$

# Maximum Sliced-Wasserstein distance

Sliced-Wasserstein distance of order  $p$ .

$$\mathbf{SW}_p(\mu, \nu) = \left( \int_{\mathbb{S}^{d-1}} \mathbf{W}_p^p(\mathcal{R}l_\mu(\cdot, \theta), \mathcal{R}l_\nu(\cdot, \theta)) d\theta \right)^{1/p}$$

Maximum Sliced-Wasserstein distance of order  $p$ . [2]

$$\max\text{-}\mathbf{SW}_p(\mu, \nu) = \max_{\theta \in \mathbb{S}^{d-1}} \mathbf{W}_p(\mathcal{R}l_\mu(\cdot, \theta), \mathcal{R}l_\nu(\cdot, \theta))$$

What happens if we use *nonlinear* projections?

# **Generalized Sliced-Wasserstein Distances**

---

# Generalized Radon transform

The Generalized Radon transform (GRT,  $\mathcal{G}$ ) integrates the data distribution over **general hyper-surfaces** (instead of hyperplanes)

$H(t, \theta) = \{x \in \mathbb{R}^d \mid g(x, \theta) = t\}$ , i.e. :

$$\forall (t, \theta) \in \mathbb{R} \times \Omega_\theta, \mathcal{G}h(t, \theta) = \int_{\mathbb{R}^d} h(x) \delta(t - g(x, \theta)) dx$$

where  $\Omega_\theta$  is a compact set of feasible parameters for  $\theta \mapsto g(\cdot, \theta)$ .

# Generalized Radon transform

$g$  is the *defining function* and must satisfy:

H1.  $g$  is a real-valued  $C^\infty$  function on  $\mathcal{X} \times (\mathbb{R}^n \setminus \{0\})$

H2.  $g(x, \theta)$  is homogeneous of degree one in  $\theta$ , i.e.,

$$\forall \lambda \in \mathbb{R}, g(x, \lambda\theta) = \lambda g(x, \theta)$$

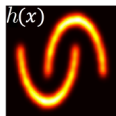
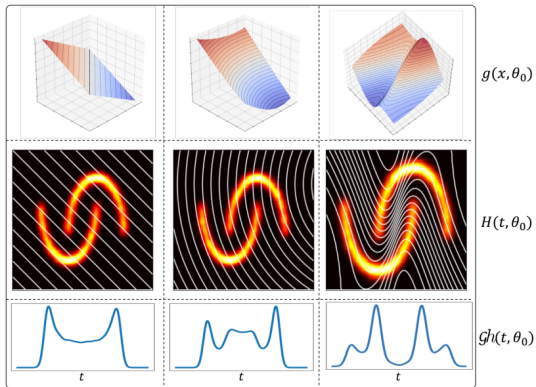
H3.  $g$  is non-degenerate in the sense that,

$$\forall (x, \theta) \in \mathcal{X} \times \mathbb{R}^n \setminus \{0\}, \frac{\partial g}{\partial x}(x, \theta) \neq 0$$

H4. The mixed Hessian of  $g$  is strictly positive, i.e.,

$$\det \left( \left( \frac{\partial^2 g}{\partial x^i \partial \theta^j} \right)_{i,j} \right) > 0$$

# Generalized Radon transform



Input distribution

$Gh(t, \theta)$ : Slices with respect to different  $g(t, \theta)$

$$H(t, \theta) = \{x | g(x, \theta) = t\}$$

# Generalized Sliced-Wasserstein distances

Let  $\mu$  and  $\nu$  be probability measures with finite  $p$ 'th moment, defined on  $X \subset \mathbb{R}^d$ , with corresponding densities  $l_\mu$  and  $l_\nu$ .

**Generalized Sliced-Wasserstein distance of order  $p$ .**

$$\mathbf{GSW}_p(\mu, \nu) = \left( \int_{\Omega_\theta} \mathbf{W}_p^p(\mathcal{G}l_\mu(\cdot, \theta), \mathcal{G}l_\nu(\cdot, \theta)) d\theta \right)^{1/p}$$

**Maximum Generalized Sliced-Wasserstein distance of order  $p$ .**

$$\max\text{-}\mathbf{GSW}_p(\mu, \nu) = \max_{\theta \in \Omega_\theta} \mathbf{W}_p(\mathcal{G}l_\mu(\cdot, \theta), \mathcal{G}l_\nu(\cdot, \theta))$$



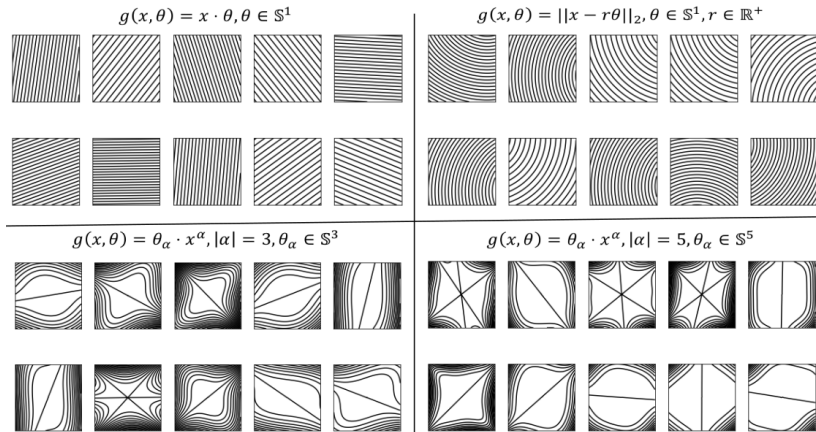
## Proposition

*The Generalized Sliced-Wasserstein and Maximum Generalized Sliced-Wasserstein distances of order  $p$  are, indeed, distances if and only if the Generalized Radon transform is injective.*

If the GRT is not injective, then GSW and max-GSW are *pseudo-metrics*.

# Choice of the defining function $g$

There exists a family of defining functions that provide injective GRT (e.g., homogeneous polynomials with odd degree).



Randomly generated integration hypersurfaces (curves in 2D)

## Choice of the defining function $g$

The memory complexity of polynomial projections grows exponentially with the data dimension and the polynomial degree.

Alternative solution: **define  $g$  as a neural network**. We propose a multi-layer fully connected network with *leaky ReLU* activations.

With a neural network-based defining function, minimizing the max-GSW distance between two distributions is analogical to the adversarial learning.

# Experiments

---

Consider the following problem:

$$\min_{\mu} \mathbf{GSW}_p(\mu, \nu),$$

where  $\nu$  is a target distribution and  $\mu$  is the source distribution, initialized to a normal distribution.

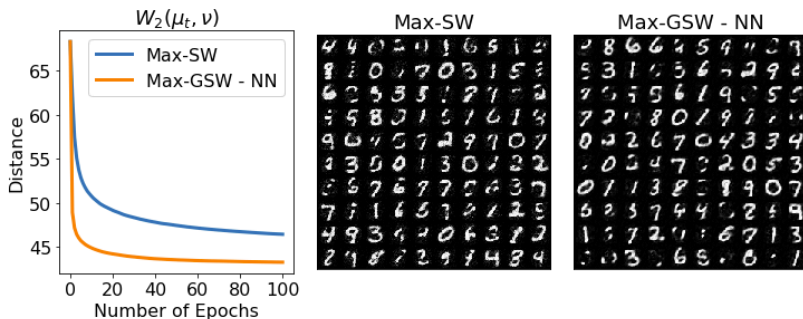
The optimization is then solved iteratively via:

$$\partial_t \mu_t = -\nabla \mathbf{GSW}_p(\mu_t, \nu), \quad \mu_0 = \mathcal{N}(0, 1).$$

## GSW flows: **Synthetical datasets**

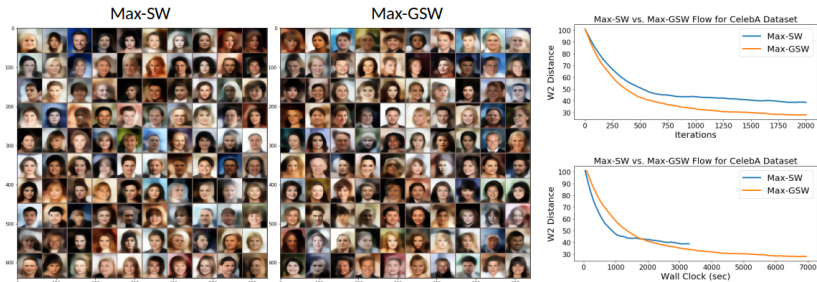
(See videos in the supplementary files.)

# GSW flows: Hand-written digit recognition



Wasserstein distance of order 2 between source and target distributions for the MNIST dataset.

# GSW flows: Face generation



Wasserstein distance of order 2 between source and target distributions for the CelebA dataset.



## Conclusion

---

# Summary

- We introduced a new family of optimal transport metrics for probability measures that generalizes the Sliced-Wasserstein distance.
- We provided theoretical conditions that yield GSW and max-GSW to be, indeed, distances.
- We empirically demonstrated the superior performance of GSW and max-GSW over SW in various generative modeling applications.



N. Bonnotte.

**Unidimensional and evolution methods for optimal transportation.**

PhD thesis, Université Paris 11, France, 2013.



I. Deshpande, Y.-T. Hu, R. Sun, A. Pyrros, N. Siddiqui, S. Koyejo, Z. Zhao, D. Forsyth, and A. Schwing.

**Max-sliced wasserstein distance and its use for gans.**

In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.



S. Kolouri, Y. Zou, and G. K. Rohde.

**Sliced-Wasserstein kernels for probability distributions.**

In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4876–4884, 2016.