

# Deep Generative Models

## Lecture 6

Roman Isachenko



2025, Spring

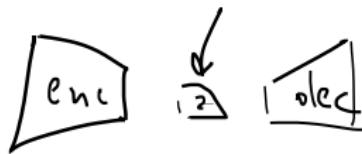
## Recap of previous lecture

### Assumptions

- ▶ Let  $c \sim \text{Categorical}(\pi)$ , where

$$\pi = (\pi_1, \dots, \pi_K), \quad \pi_k = P(c = k), \quad \sum_{k=1}^K \pi_k = 1.$$

- ▶ Let VAE model has discrete latent representation  $c$  with prior  $p(c) = \text{Uniform}\{1, \dots, K\}$ .



### ELBO

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(c|\mathbf{x}, \phi)} \log p(\mathbf{x}|c, \theta) - \underbrace{KL(q(c|\mathbf{x}, \phi)||p(c))}_{\text{KL term}} \rightarrow \max_{\phi, \theta} .$$

$$KL(q(c|\mathbf{x}, \phi)||p(c)) = \underbrace{-H(q(c|\mathbf{x}, \phi))}_{\text{KL term}} + \underbrace{\log K}_{\text{KL term}}$$

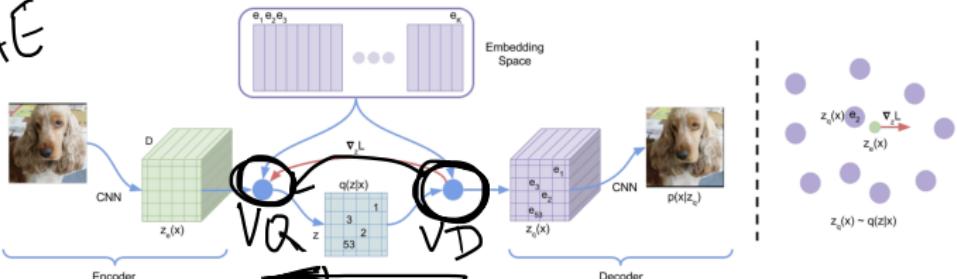
### Vector quantization

Define the dictionary space  $\{\mathbf{e}_k\}_{k=1}^K$ , where  $\mathbf{e}_k \in \mathbb{R}^L$ ,  $K$  is the size of the dictionary.

$$\mathbf{z}_q = \mathbf{q}(\mathbf{z}) = \mathbf{e}_{k^*}, \quad \text{where } k^* = \arg \min_k \|\mathbf{z} - \mathbf{e}_k\|.$$

# Recap of previous lecture

VQVAE



## Deterministic variational posterior

$$q(c_{ij} = k^* | \mathbf{x}, \phi) = \begin{cases} 1, & \text{for } k^* = \arg \min_k \|[\mathbf{z}_e]_{ij} - \mathbf{e}_k\|; \\ 0, & \text{otherwise.} \end{cases}$$

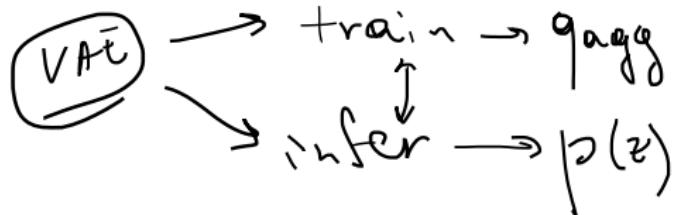
ELBO

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(c|\mathbf{x}, \phi)} \log p(\mathbf{x}|\mathbf{e}_c, \theta) - \log K = \underbrace{\log p(\mathbf{x}|\mathbf{z}_q, \theta)}_{\text{ELBO term}} - \underbrace{\log K}_{\text{KL divergence}}$$

## Straight-through gradient estimation

$$\frac{\partial \log p(\mathbf{x}|\mathbf{z}_q, \theta)}{\partial \phi} = \frac{\partial \log p(\mathbf{x}|\mathbf{z}_q, \theta)}{\partial \mathbf{z}_q} \cdot \frac{\partial \mathbf{z}_q}{\partial \phi} \approx \frac{\partial \log p(\mathbf{x}|\mathbf{z}_q, \theta)}{\partial \mathbf{z}_q} \cdot \frac{\partial \mathbf{z}_e}{\partial \phi}$$

## Recap of previous lecture



## Theorem

$$\frac{1}{n} \sum_{i=1}^n KL(q(\mathbf{z}|\mathbf{x}_i, \phi) || p(\mathbf{z})) = KL(q_{\text{agg}}(\mathbf{z}|\phi) || p(\mathbf{z})) + \mathbb{I}_q[\mathbf{x}, \mathbf{z}].$$

## ELBO surgery

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\phi, \theta}(\mathbf{x}_i) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{q(\mathbf{z}|\mathbf{x}_i, \phi)} \log p(\mathbf{x}_i|\mathbf{z}, \theta)}_{\text{Reconstruction loss}} - \underbrace{\mathbb{I}_q[\mathbf{x}, \mathbf{z}]}_{\text{MI}} - \underbrace{KL(q_{\text{agg}}(\mathbf{z}|\phi) || p(\mathbf{z}))}_{\text{Marginal KL}}$$

## Optimal prior

$$KL(q_{\text{agg}}(\mathbf{z}|\phi) || p(\mathbf{z})) = 0 \Leftrightarrow p(\mathbf{z}) = q_{\text{agg}}(\mathbf{z}) = \frac{1}{n} \sum_{i=1}^n q(\mathbf{z}|\mathbf{x}_i, \phi).$$

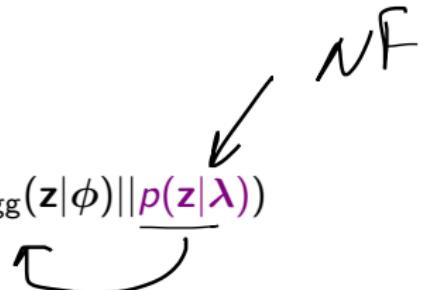
The optimal prior distribution  $p(\mathbf{z})$  is the aggregated variational posterior distribution  $q_{\text{agg}}(\mathbf{z}|\phi)$ .

## Recap of previous lecture

- ▶ Standard Gaussian  $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}) \Rightarrow$  over-regularization;
- ▶  $\underline{p(\mathbf{z})} = q_{\text{agg}}(\mathbf{z}|\phi) = \frac{1}{n} \sum_{i=1}^n q(\mathbf{z}|\mathbf{x}_i, \phi) \Rightarrow$  overfitting and highly expensive.

## ELBO revisiting

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\phi, \theta}(\mathbf{x}_i) = \text{RL} - \text{MI} - KL(q_{\text{agg}}(\mathbf{z}|\phi) || \overbrace{p(\mathbf{z}|\lambda)})$$



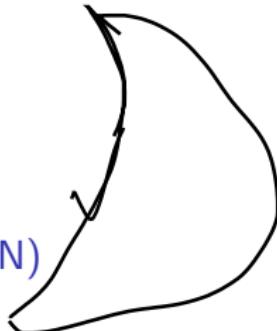
It is Forward KL with respect to  $p(\mathbf{z}|\lambda)$ .

## ELBO with learnable VAE prior

$$\begin{aligned}\mathcal{L}_{\phi, \theta}(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \phi)} [\log p(\mathbf{x}|\mathbf{z}, \theta) + \log p(\mathbf{z}|\lambda) - \log q(\mathbf{z}|\mathbf{x}, \phi)] \\ &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \phi)} \left[ \log p(\mathbf{x}|\mathbf{z}, \theta) + \underbrace{\left( \log p(f_\lambda(\mathbf{z})) + \log |\det(\mathbf{J}_f)| \right)}_{\text{flow-based prior}} - \log q(\mathbf{z}|\mathbf{x}, \phi) \right] \\ \mathbf{z} &= \mathbf{f}_\lambda^{-1}(\mathbf{z}^*) = \mathbf{g}_\lambda(\mathbf{z}^*), \quad \mathbf{z}^* \sim p(\mathbf{z}^*) = \mathcal{N}(0, \mathbf{I})\end{aligned}$$

# Outline

1. Likelihood-free learning ✓
2. Generative adversarial networks (GAN)
3. Wasserstein distance
4. Wasserstein GAN



# Outline

1. Likelihood-free learning
2. Generative adversarial networks (GAN)
3. Wasserstein distance
4. Wasserstein GAN

## Likelihood based models

Poor likelihood  
Great samples

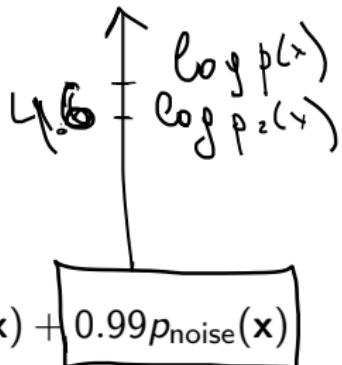
$$p(\mathbf{x}) = \prod p(x_i) \quad \rightarrow$$

$\sum_{i=1}^n \log p(x_i)$

Great likelihood  
Poor samples

$$p_1(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathcal{N}(\mathbf{x} | \mathbf{x}_i, \epsilon \mathbf{I})$$

$$p_2(\mathbf{x}) = 0.01 p(\mathbf{x}) + 0.99 p_{\text{noise}}(\mathbf{x})$$



For small  $\epsilon$  this model will generate samples with great quality, but likelihood of test sample will be very poor.

$$\begin{aligned} \log [0.01p(\mathbf{x}) + 0.99p_{\text{noise}}(\mathbf{x})] &\geq \\ \geq \log [0.01p(\mathbf{x})] &= \boxed{\log p(\mathbf{x}) - \log 100} \end{aligned}$$

Noisy irrelevant samples, but for high dimensions  $\log p(\mathbf{x})$  becomes proportional to  $m$ .

- ▶ Likelihood is not a perfect quality measure for generative model.
- ▶ Likelihood could be intractable.

# Likelihood-free learning

## Where did we start

We would like to approximate true data distribution  $\pi(\mathbf{x})$ . Instead of searching true  $\pi(\mathbf{x})$  over all probability distributions, learn function approximation  $p(\mathbf{x}|\theta) \approx \pi(\mathbf{x})$ .

Imagine we have two sets of samples

- ▶  $\{\mathbf{x}_i\}_{i=1}^{n_1} \sim \pi(\mathbf{x})$  – real samples;
- ▶  $\{\mathbf{x}_i\}_{i=1}^{n_2} \sim p(\mathbf{x}|\theta)$  – generated (or fake) samples.

Let define discriminative model (classifier):

$$p(y = 1|\mathbf{x}) = P(\{\mathbf{x} \sim \pi(\mathbf{x})\}); \quad p(y = 0|\mathbf{x}) = P(\{\mathbf{x} \sim p(\mathbf{x}|\theta)\})$$

## Assumption

Generative distribution  $p(\mathbf{x}|\theta)$  equals to the true distribution  $\pi(\mathbf{x})$  if we can not distinguish them using discriminative model  $p(y|\mathbf{x})$ . It means that  $p(y = 1|\mathbf{x}) = 0.5$  for each sample  $\mathbf{x}$ .

# Generative adversarial networks (GAN)

- ▶ The more powerful discriminative model we will have, the more likely we will get the "best" generative distribution  $p(\mathbf{x}|\theta)$ .
- ▶ The most common way to learn a classifier is to minimize cross entropy loss.

$$p(\mathbf{x}|\theta)$$

Cross entropy for discriminative model

$$\min_{p(y|\mathbf{x})} \left[ -\mathbb{E}_{\pi(\mathbf{x})} \log p(y=1|\mathbf{x}) - \mathbb{E}_{p(\mathbf{x}|\theta)} \log p(y=0|\mathbf{x}) \right]$$

$$\max_{p(y|\mathbf{x})} \left[ \mathbb{E}_{\pi(\mathbf{x})} \log p(y=1|\mathbf{x}) + \mathbb{E}_{p(\mathbf{x}|\theta)} \log p(y=0|\mathbf{x}) \right]$$

Generative model

Assume generative model  $p(\mathbf{x}, \mathbf{z}|\theta) = p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z})$  with the base distribution  $p(\mathbf{z})$  and deterministic map  $p(\mathbf{x}|\mathbf{z}, \theta) = \delta(\mathbf{x} - \mathbf{G}_\theta(\mathbf{z}))$ .

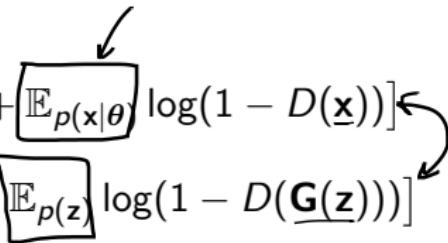
# Generative adversarial networks (GAN)

## Cross entropy for discriminative model

$$\max_{p(y|x)} [\mathbb{E}_{\pi(x)} \log p(y=1|x) + \mathbb{E}_{p(x|\theta)} \log p(y=0|x)]$$

- ▶ **Discriminator:** a classifier  $p(y=1|x, \phi) = D_\phi(x) \in [0, 1]$ , which distinguishes real samples from  $\pi(x)$  and generated samples from  $p(x|\theta)$ . Discriminator tries to **minimize** cross entropy.
- ▶ **Generator:** generative model  $x = G_\theta(z)$  with  $z \sim p(z)$ , which makes the generated sample more realistic. Generator tries to **maximize** cross entropy.

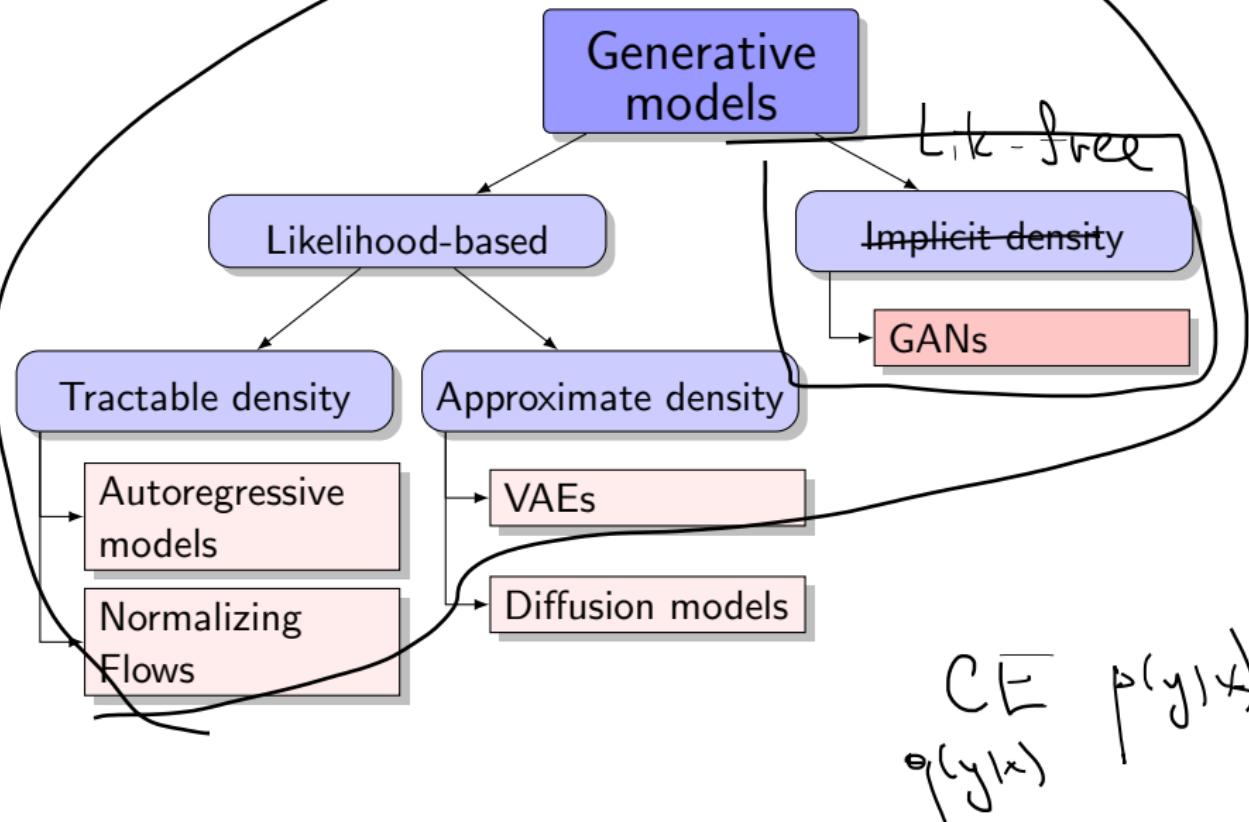
## GAN Objective

$$\begin{aligned} & \min_G \max_D [\mathbb{E}_{\pi(x)} \log D(x) + \boxed{\mathbb{E}_{p(x|\theta)} \log(1 - D(\underline{x}))}] \\ & \min_G \max_D [\mathbb{E}_{\pi(x)} \log D(x) + \boxed{\mathbb{E}_{p(z)} \log(1 - D(\underline{G(z)}))}] \end{aligned}$$


# Outline

1. Likelihood-free learning
2. Generative adversarial networks (GAN)
3. Wasserstein distance
4. Wasserstein GAN

# Generative models zoo



# GAN optimality

## Theorem

The minimax game

$$\min_{\underline{G}} \max_D \left[ \mathbb{E}_{\pi(x)} \log D(x) + \mathbb{E}_{p(z)} \log(1 - D(\mathbf{G}(z))) \right]$$

$V(G, D)$

has the global optimum  $\pi(x) = p(x|\theta)$ , in this case  $D^*(x) = 0.5$ .

Proof (fixed  $G$ )

$$\begin{aligned} V(G, D) &= \mathbb{E}_{\pi(x)} \log D(x) + \mathbb{E}_{p(x|\theta)} \log(1 - D(x)) \\ &= \int \underbrace{[\pi(x) \log D(x) + p(x|\theta) \log(1 - D(x))]}_{y(D)} dx \end{aligned}$$

$$\frac{dy(D)}{dD} = \frac{\pi(x)}{D(x)} - \frac{p(x|\theta)}{1 - D(x)} = 0 \Rightarrow D^*(x) = \boxed{\frac{\pi(x)}{\pi(x) + p(x|\theta)}}$$

# GAN optimality

Proof continued (fixed  $D = D^*$ )

$$\begin{aligned} \overbrace{V(G, D^*)} &= \mathbb{E}_{\pi(x)} \log \left( \frac{\pi(x)}{\pi(x) + p(x|\theta)} \right) + \mathbb{E}_{p(x|\theta)} \log \left( \frac{p(x|\theta)}{\pi(x) + p(x|\theta)} \right) \\ &= KL \left( \pi(x) \parallel \frac{\pi(x) + p(x|\theta)}{2} \right) + KL \left( p(x|\theta) \parallel \frac{\pi(x) + p(x|\theta)}{2} \right) - 2 \log 2 \\ &= \boxed{2 JSD(\pi(x) \parallel p(x|\theta))} - 2 \log 2. \end{aligned}$$

Jensen-Shannon divergence (symmetric KL divergence)

$$\boxed{JSD(\pi(x) \parallel p(x|\theta))} = \frac{1}{2} \left[ KL \left( \pi(x) \parallel \frac{\pi(x) + p(x|\theta)}{2} \right) + KL \left( p(x|\theta) \parallel \frac{\pi(x) + p(x|\theta)}{2} \right) \right] \geq 0$$

Could be used as a distance measure!

$$\boxed{V(G^*, D^*) = -2 \log 2}, \quad \boxed{\pi(x) = p(x|\theta)}, \quad \boxed{D^*(x) = 0.5.} \quad \frac{1)}{11+1}$$

# GAN optimality

$$\mathcal{D} = \mathcal{D}^\pi$$

## Theorem

The minimax game

$$\min_G \max_D \underbrace{\left[ \mathbb{E}_{\pi(x)} \log D(x) + \mathbb{E}_{p(z)} \log(1 - D(\mathbf{G}(z))) \right]}_{V(G,D)}$$

has the global optimum  $\pi(x) = p(x|\theta)$ , in this case  $D^*(x) = 0.5$ .

## Expectations

If the generator could be any function and the discriminator is optimal at every step, then the generator is guaranteed to converge to the data distribution.  $\mathbb{P}(x)$

## Reality

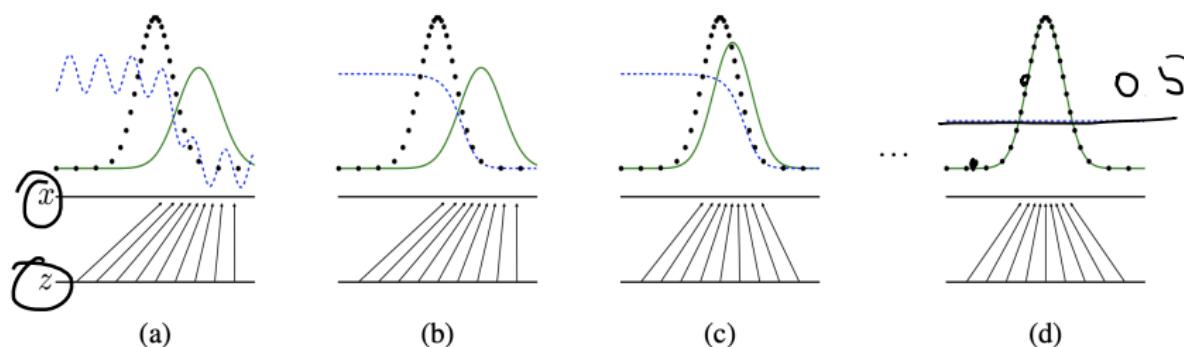
- ▶ Generator updates are made in parameter space, discriminator is not optimal at every step.
- ▶ Generator and discriminator loss keeps oscillating during GAN training.

# GAN training

Let further assume that generator and discriminator are parametric models:  $D_\phi(\mathbf{x})$  and  $\mathbf{G}_\theta(\mathbf{z})$ .

## Objective

$$\min_{\theta} \max_{\phi} [\mathbb{E}_{\pi(\mathbf{x})} \log D_\phi(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D_\phi(\mathbf{G}_\theta(\mathbf{z})))]$$



- ▶  $\mathbf{z} \sim p(\mathbf{z})$  is a latent variable.
- ▶  $p(\mathbf{x}|\mathbf{z}, \theta) = \delta(\mathbf{x} - \mathbf{G}_\theta(\mathbf{z}))$  is deterministic decoder (like NF).
- ▶ We do not have encoder at all.

## Mode collapse

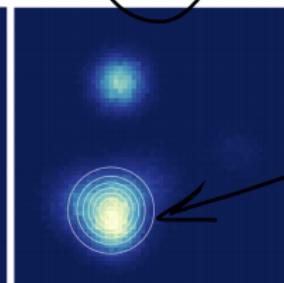
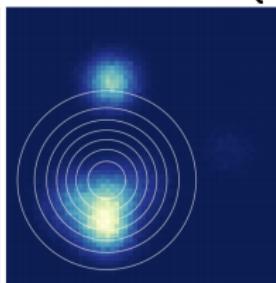
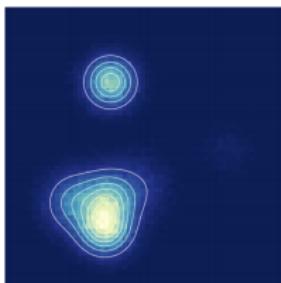
The phenomena where the generator of a GAN collapses to one or few distribution modes.

$$p(l \times | \theta) = \mathcal{N}(\mu, \sigma^2)$$

$$KL(\pi \| p) \rightarrow \min_{\theta}$$

$$JSD(\pi \| p) \rightarrow \min_{\theta}$$

Data  $\pi(l)$



JSD



Alternate architectures, adding regularization terms, injecting small noise perturbations and other millions bags and tricks are used to avoid the mode collapse.

Goodfellow I. J. et al. Generative Adversarial Networks, 2014

Metz L. et al. Unrolled Generative Adversarial Networks, 2016

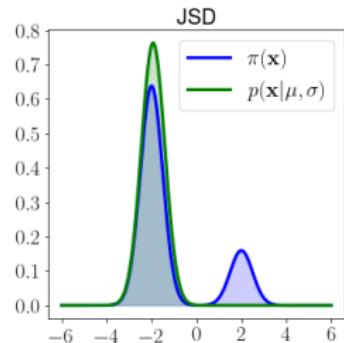
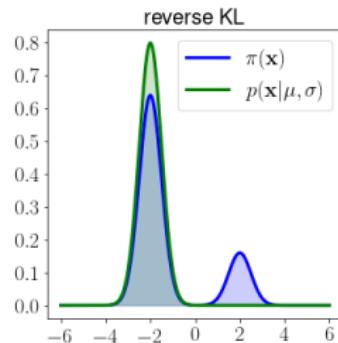
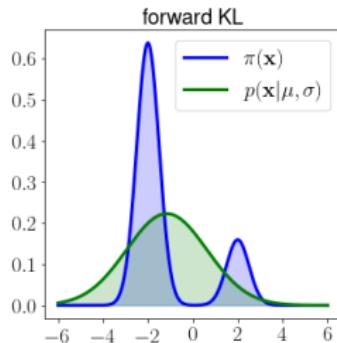
# Jensen-Shannon vs Kullback-Leibler

- ▶  $\pi(\mathbf{x})$  is a fixed mixture of 2 gaussians.
- ▶  $p(\mathbf{x}|\mu, \sigma) = \mathcal{N}(\mu, \sigma^2)$ .

Mode covering vs mode seeking

$$KL(\pi||p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}, \quad KL(p||\pi) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{\pi(\mathbf{x})} d\mathbf{x}$$

$$JSD(\pi||p) = \frac{1}{2} \left[ KL \left( \pi(\mathbf{x}) || \frac{\pi(\mathbf{x}) + p(\mathbf{x})}{2} \right) + KL \left( p(\mathbf{x}) || \frac{\pi(\mathbf{x}) + p(\mathbf{x})}{2} \right) \right]$$

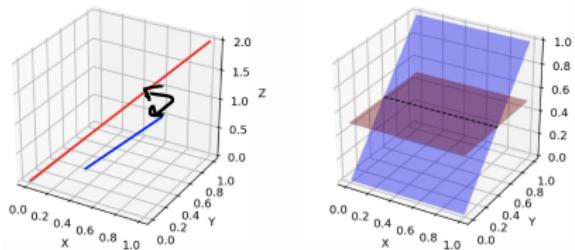


# Outline

1. Likelihood-free learning
2. Generative adversarial networks (GAN)
3. Wasserstein distance
4. Wasserstein GAN

## Informal theoretical results

- The dimensionality of  $\mathbf{z}$  is lower than the dimensionality of  $\mathbf{x}$ .  
Hence, support of  $p(\mathbf{x}|\theta)$  with  $\mathbf{x} = \mathbf{G}_\theta(\mathbf{z})$  lies on low-dimensional manifold.
- Distribution of real images  $\pi(\mathbf{x})$  is also concentrated on a low dimensional manifold.



$x \in \mathbb{R}^3$   
1024x1024

$x_{ijk}$

- If  $\pi(\mathbf{x})$  and  $p(\mathbf{x}|\theta)$  have disjoint supports, then there is a smooth optimal discriminator.
- For such low-dimensional disjoint manifolds

$$KL(\pi||p) = KL(p||\pi) = \infty, \quad JSD(\pi||p) = \log 2$$

$G_\theta$   
 $O$

Weng L. From GAN to WGAN, 2019

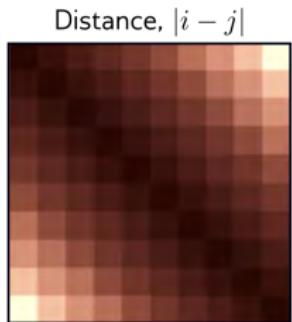
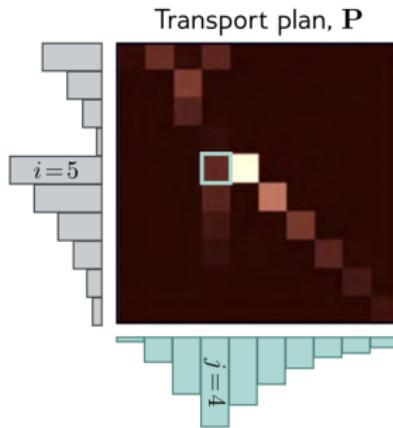
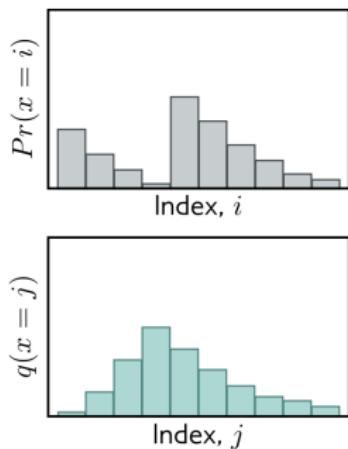
Arjovsky M., Bottou L. Towards Principled Methods for Training Generative Adversarial Networks, 2017

# Wasserstein distance (discrete)

A.k.a. **Earth Mover's distance**.

## Optimal transport formulation

The minimum cost of moving and transforming a pile of dirt in the shape of one probability distribution to the shape of the other distribution.



$$\text{Wasserstein distance} = \sum \mathbf{P} \cdot |i - j|$$

## Wasserstein distance (continuous)

$$W(\pi, p) = \inf_{\gamma \in \Gamma(\pi, p)} \mathbb{E}_{(x,y) \sim \gamma} \|x - y\| = \inf_{\gamma \in \Gamma(\pi, p)} \|x - y\| \gamma(x, y) dx dy$$

- ▶  $\gamma(x, y)$  – transportation plan (the amount of "dirt" that should be transported from point  $x$  to point  $y$ )

$$\int \gamma(x, y) dx = p(y); \quad \int \gamma(x, y) dy = \pi(x).$$

- ▶  $\Gamma(\pi, p)$  – the set of all joint distributions  $\gamma(x, y)$  with marginals  $\pi$  and  $p$ .
- ▶  $\gamma(x, y)$  – the amount,  $\|x - y\|$  – the distance.

## Wasserstein metric

$$W_s(\pi, p) = \inf_{\gamma \in \Gamma(\pi, p)} (\mathbb{E}_{(x,y) \sim \gamma} \|x - y\|^s)^{1/s}$$

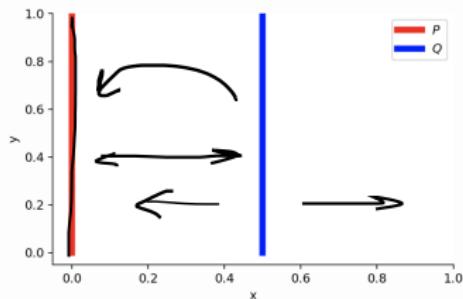
Here we will use  $W(\pi, p) = W_1(\pi, p)$  that corresponds to the optimal transport formulation.

# Wasserstein distance vs KL vs JSD

Consider 2d distributions

$$\boxed{\pi(x, y) = (0, U[0, 1])}$$

$$p(x, y|\theta) = (\theta, U[0, 1])$$



- $\theta = 0$ . Distributions are the same

$$KL(\pi||p) = KL(p||\pi) = JSD(p||\pi) = W(\pi, p) = 0$$

- $\theta \neq 0$

$$KL(\pi||p) = \int_{U[0,1]} 1 \log \frac{1}{0} dy = \infty = KL(p||\pi)$$

$$JSD(\pi||p) = \frac{1}{2} \left( \int_{U[0,1]} 1 \log \frac{1}{1/2} dy + \int_{U[0,1]} 1 \log \frac{1}{1/2} dy \right) = \log 2$$

$$\boxed{W(\pi, p) = |\theta|}$$

# Wasserstein distance vs KL vs JSD

## Theorem 1

Let  $\mathbf{G}_\theta(\mathbf{z})$  be (almost) any feedforward neural network, and  $p(\mathbf{z})$  a prior over  $\mathbf{z}$  such that  $\mathbb{E}_{p(\mathbf{z})} \|\mathbf{z}\| < \infty$ . Then therefore  $W(\pi, p)$  is continuous everywhere and differentiable almost everywhere.

## Theorem 2

Let  $\pi$  be a distribution on a compact space  $\mathcal{X}$  and  $\{p_t\}_{t=1}^\infty$  be a sequence of distributions on  $\mathcal{X}$ .

$$KL(\pi || p_t) \rightarrow 0 \text{ (or } KL(p_t || \pi) \rightarrow 0) \quad (1)$$

$$JSD(\pi || p_t) \rightarrow 0 \quad (2)$$

$$W(\pi || p_t) \rightarrow 0 \quad (3)$$

Then, considering limits as  $t \rightarrow \infty$ , (1) implies (2), (2) implies (3).

# Outline

1. Likelihood-free learning
2. Generative adversarial networks (GAN)

3. Wasserstein distance

4. Wasserstein GAN

$$W(\pi \parallel \rho)$$

# Wasserstein GAN

## Wasserstein distance

$$W(\pi||p) = \inf_{\gamma \in \Gamma(\pi, p)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \|\mathbf{x} - \mathbf{y}\| = \underbrace{\inf_{\gamma \in \Gamma(\pi, p)}}_{\text{circled}} \int \|\mathbf{x} - \mathbf{y}\| \underbrace{\gamma(\mathbf{x}, \mathbf{y})}_{\text{circled}} d\mathbf{x} d\mathbf{y}$$

The infimum across all possible joint distributions in  $\Gamma(\pi, p)$  is intractable.

## Theorem (Kantorovich-Rubinstein duality)

$$W(\pi||p) = \frac{1}{K} \max_{\substack{\text{circled} \\ \|f\|_L \leq K}} [\mathbb{E}_{\pi(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{p(\mathbf{x})} f(\mathbf{x})],$$

→ min

where  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ ,  $\|f\|_L \leq K$  are  $K$ -Lipschitz continuous functions ( $f : \mathcal{X} \rightarrow \mathbb{R}$ )

$$|f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq K \|\mathbf{x}_1 - \mathbf{x}_2\|, \quad \text{for all } \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}.$$

Now we need only samples to get Monte Carlo estimate for  $W(\pi||p)$ .

# Wasserstein GAN

Theorem (Kantorovich-Rubinstein duality)

$$W(\pi || p) = \frac{1}{K} \max_{\|f\|_L \leq K} [\mathbb{E}_{\pi(x)} f(x) - \mathbb{E}_{p(x)} f(x)],$$

- ▶ Now we have to ensure that  $f$  is  $K$ -Lipschitz continuous.
- ▶ Let  $f_\phi(x)$  be a feedforward neural network parametrized by  $\phi$ .
- ▶ If parameters  $\phi$  lie in a compact set  $\Phi$  then  $f_\phi(x)$  will be  $|K|$ -Lipschitz continuous function.
- ▶ Let the parameters be clamped to a fixed box  $\boxed{\Phi \in [-c, c]^d}$  (e.g.  $c = 0.01$ ) after each gradient update.

$$\begin{aligned} K \cdot W(\pi || p) &= \max_{\|f\|_L \leq K} [\mathbb{E}_{\pi(x)} f(x) - \mathbb{E}_{p(x)} f(x)] \\ &\geq \max_{\phi \in \Phi} [\mathbb{E}_{\pi(x)} f_\phi(x) - \mathbb{E}_{p(x)} f_\phi(x)] \end{aligned}$$

# Wasserstein GAN

## Standard GAN objective

$$\min_{\theta} \max_{\phi} \mathbb{E}_{\pi(x)} \log D_{\phi}(x) + \mathbb{E}_{p(z)} \log(1 - D_{\phi}(G_{\theta}(z)))$$

## WGAN objective

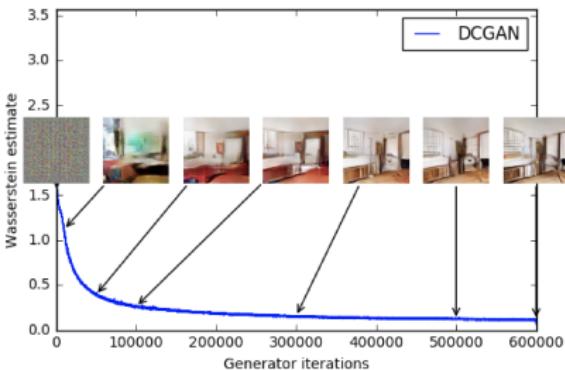
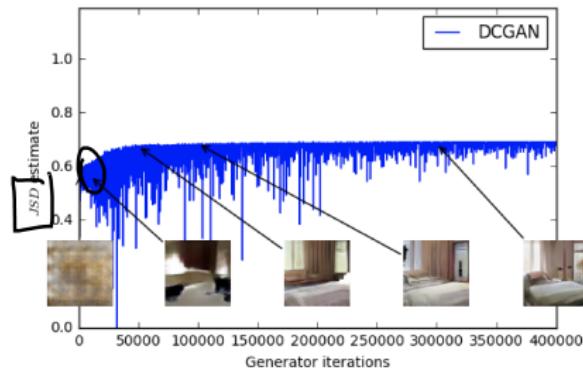
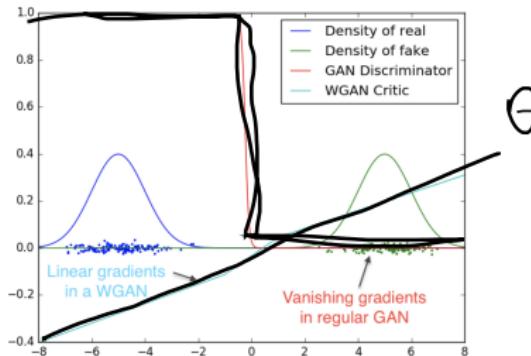
$$\min_{\theta} W(\pi || p) \approx \min_{\theta} \max_{\phi \in \Phi} [\mathbb{E}_{\pi(x)} f_{\phi}(x) - \mathbb{E}_{p(z)} f_{\phi}(G_{\theta}(z))].$$

$f \notin [0, 1] \subset (-\infty, \infty)$

- ▶ Discriminator  $D$  is similar to the function  $f$ , but it is not a classifier anymore. In the WGAN model, function  $f$  is usually called **critic**.
- ▶ "Weight clipping is a clearly terrible way to enforce a Lipschitz constraint".
  - ▶ If the clipping parameter  $c$  is too large, it is hard to train the critic till optimality.
  - ▶ If the clipping parameter  $c$  is too small, it could lead to vanishing gradients.

# Wasserstein GAN

- ▶ WGAN has non-zero gradients for disjoint supports.
- ▶  $JSD(\pi||p)$  correlates poorly with the sample quality. Stays constant nearly maximum value  $\log 2 \approx 0.69$ .
- ▶  $W(\pi||p)$  is highly correlated with the sample quality.



## Summary

- ▶ Likelihood is not a perfect criteria to measure quality of generative model.
- ▶ Adversarial learning suggests to solve minimax problem to match the distributions.
- ▶ GAN tries to optimize Jensen-Shannon divergence (in theory).
- ▶ KL and JS divergences work poorly as model objective in the case of disjoint supports.
- ▶ Earth-Mover distance is a more appropriate objective function for distribution matching problem.
- ▶ Kantorovich-Rubinstein duality gives the way to calculate the EM distance using only samples.
- ▶ Wasserstein GAN uses the weight clipping to ensure the Lipschitness of the critic.