

# Deep Generative Models

## Lecture 7

Roman Isachenko



2025, Spring

## Recap of previous lecture

### Likelihood-free learning

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

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.

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

### 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  $\mathbf{x}$ .

- ▶ **Generator:** generative model  $\mathbf{x} = \mathbf{G}(\mathbf{z})$ , which makes generated sample more realistic.
- ▶ **Discriminator:** a classifier  $D(\mathbf{x}) \in [0, 1]$ , which distinguishes real samples from generated samples.

## Recap of previous lecture

### GAN optimality theorem

The minimax game

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

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

$$\min_G V(G, D^*) = \min_G [2JSD(\pi || p) - \log 4] = -\log 4, \quad \pi(\mathbf{x}) = p(\mathbf{x}|\theta).$$

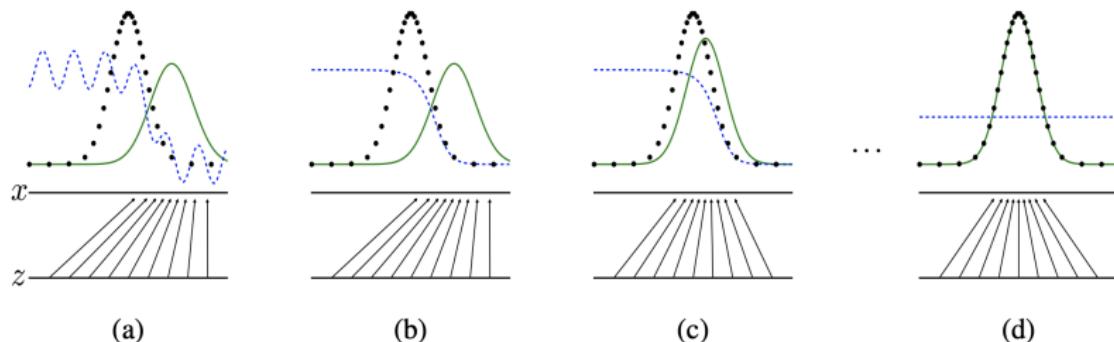
If the generator could be **any** function and the discriminator is **optimal** at every step, then the generator is **guaranteed to converge** to the data distribution.

## Recap of previous lecture

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

## Objective

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



## Recap of previous lecture

### Main problems of standard GAN

- ▶ Vanishing gradients (solution: non-saturating GAN);
- ▶ Mode collapse (caused by Jensen-Shannon divergence).

### Standard GAN

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

### Informal theoretical results

The real images distribution  $\pi(\mathbf{x})$  and the generated images distribution  $p(\mathbf{x}|\theta)$  are low-dimensional and have disjoint supports.  
In this case

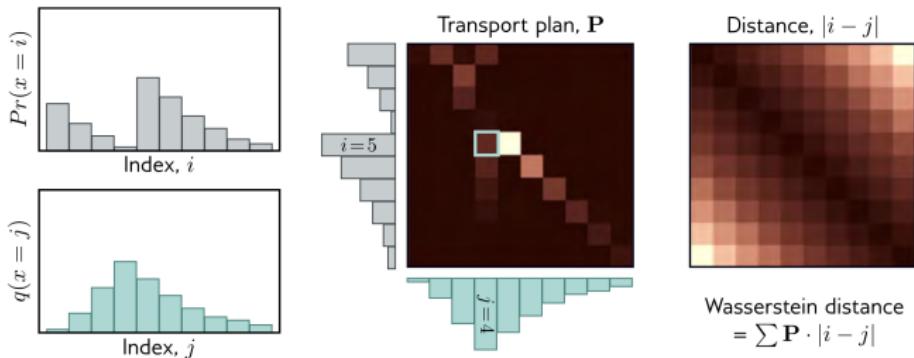
$$KL(\pi||p) = KL(p||\pi) = \infty, \quad JSD(\pi||p) = \log 2.$$

---

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

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

# Recap of previous lecture



## Wasserstein distance

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

- ▶  $\gamma(\mathbf{x}, \mathbf{y})$  – transportation plan (the amount of "dirt" that should be transported from point  $\mathbf{x}$  to point  $\mathbf{y}$ ).
- ▶  $\Gamma(\pi, p)$  – the set of all joint distributions  $\gamma(\mathbf{x}, \mathbf{y})$  with marginals  $\pi$  and  $p$  ( $\int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = p(\mathbf{y})$ ,  $\int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \pi(\mathbf{x})$ ).
- ▶  $\gamma(\mathbf{x}, \mathbf{y})$  – the amount,  $\|\mathbf{x} - \mathbf{y}\|$  – the distance.

## Recap of previous lecture

Theorem (Kantorovich-Rubinstein duality)

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

where  $\|f\|_L \leq K$  are  $K$ -Lipschitz continuous functions.

## WGAN objective

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

- ▶ Function  $f$  in WGAN is usually called *critic*.
- ▶ If parameters  $\phi$  lie in a compact set  $\Phi \in [-c, c]^d$  then  $f(x, \phi)$  will be  $K$ -Lipschitz continuous function.

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

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

# Evaluation of likelihood-free models

## Likelihood-based models

- ▶ **train part:** fit the model.
- ▶ **validation part:** tune the hyperparameters.
- ▶ **test part:** evaluate generalization by reporting the likelihood.

Not all models have tractable likelihood  
(VAE: compare ELBO values; GAN: ???).

## What do we want from samples?

- ▶ Sharpness



- ▶ Diversity



# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

## Wasserstein metric

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

## Wassestein GAN (optimal transport)

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

### Theorem

If  $\pi(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\pi, \boldsymbol{\Sigma}_\pi)$ ,  $p(\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$ , then

$$W_2^2(\pi, p) = \|\boldsymbol{\mu}_\pi - \boldsymbol{\mu}_p\|^2 + \text{tr} \left[ \boldsymbol{\Sigma}_\pi + \boldsymbol{\Sigma}_p - 2 \left( \boldsymbol{\Sigma}_\pi^{1/2} \boldsymbol{\Sigma}_p \boldsymbol{\Sigma}_\pi^{1/2} \right)^{1/2} \right]$$

## Frechet Inception Distance

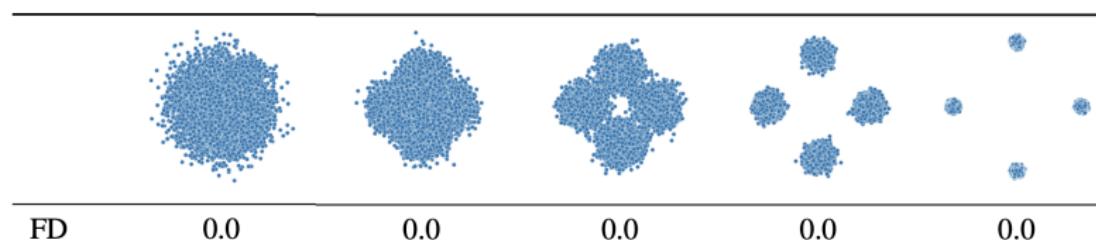
$$\text{FID}(\pi, p) = W_2^2(\pi, p)$$

## Frechet Inception Distance (FID)

$$\text{FID}(\pi, p) = \|\boldsymbol{\mu}_\pi - \boldsymbol{\mu}_p\|^2 + \text{tr} \left[ \boldsymbol{\Sigma}_\pi + \boldsymbol{\Sigma}_p - 2 \left( \boldsymbol{\Sigma}_\pi^{1/2} \boldsymbol{\Sigma}_p \boldsymbol{\Sigma}_\pi^{1/2} \right)^{1/2} \right]$$

- ▶ FID is calculated in the latent space  $\mathbf{z}$ .
- ▶ We take pretrained image embedder to get the latent representations  $\mathbf{z} = \mathbf{f}(\mathbf{x})$ .
- ▶  $\boldsymbol{\mu}_\pi$ ,  $\boldsymbol{\Sigma}_\pi$  and  $\boldsymbol{\mu}_p$ ,  $\boldsymbol{\Sigma}_p$  are the statistics of the latent representations  $\mathbf{z}$  for the samples from  $\pi(\mathbf{x})$  and  $p(\mathbf{x}|\theta)$ .

$$FID(p(\mathbf{x}), \mathcal{N}(0, \mathbf{I}))$$



## Frechet Inception Distance (FID)

$$\text{FID}(\pi, p) = \|\boldsymbol{\mu}_\pi - \boldsymbol{\mu}_p\|^2 + \text{tr} \left[ \boldsymbol{\Sigma}_\pi + \boldsymbol{\Sigma}_p - 2 \left( \boldsymbol{\Sigma}_\pi^{1/2} \boldsymbol{\Sigma}_p \boldsymbol{\Sigma}_\pi^{1/2} \right)^{1/2} \right]$$

### Drawbacks

- ▶ Dependence on the pretrained classification model.
- ▶ Usage of the normality assumption.
- ▶ May not correspond to human evaluation.

Model	Model-A	Model-B
FID	21.40	18.42
$\text{FID}_\infty$	20.16	17.19
Human rater preference	92.5%	6.9%

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

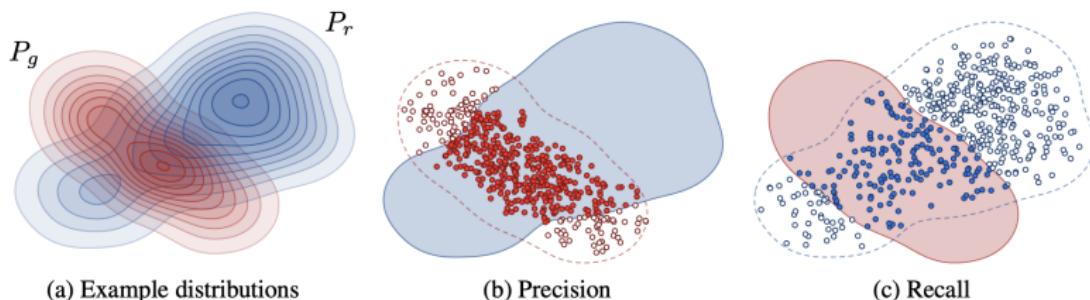
## 3. Score matching

## 4. Denoising score matching

# Precision-Recall

What do we want from samples

- ▶ **Sharpness:** generated samples should be of high quality.
- ▶ **Diversity:** their variation should match that observed in the training set.



- ▶ **Precision** denotes the fraction of generated images that are realistic.
- ▶ **Recall** measures the fraction of the training data manifold covered by the generator.

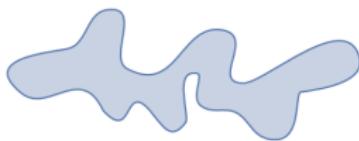
## Precision-Recall

- ▶  $\mathcal{S}_\pi = \{\mathbf{x}_i\}_{i=1}^n \sim \pi(\mathbf{x})$  – real samples;
- ▶  $\mathcal{S}_p = \{\mathbf{x}_i\}_{i=1}^n \sim p(\mathbf{x}|\theta)$  – generated samples.

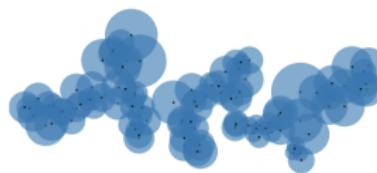
Define binary function:

$$\mathbb{I}(\mathbf{x}, \mathcal{S}) = \begin{cases} 1, & \text{if exists } \mathbf{x}' \in \mathcal{S} : \|\mathbf{x} - \mathbf{x}'\|_2 \leq \|\mathbf{x}' - \text{NN}_k(\mathbf{x}', \mathcal{S})\|_2; \\ 0, & \text{otherwise.} \end{cases}$$

$$\text{Precision}(\mathcal{S}_\pi, \mathcal{S}_p) = \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{S}_p} \mathbb{I}(\mathbf{x}, \mathcal{S}_\pi); \quad \text{Recall}(\mathcal{S}_\pi, \mathcal{S}_p) = \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{S}_\pi} \mathbb{I}(\mathbf{x}, \mathcal{S}_p).$$



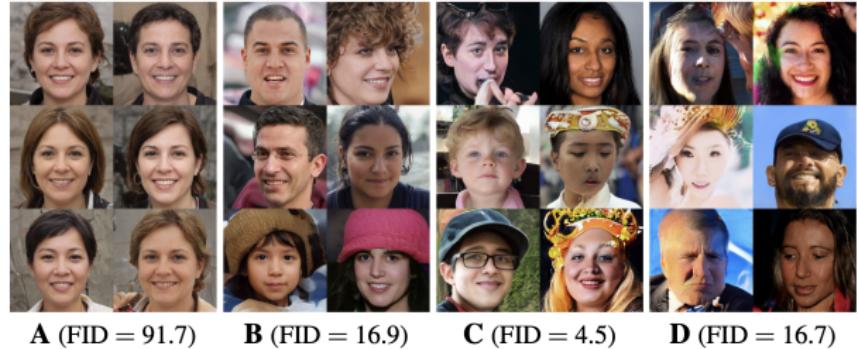
(a) True manifold



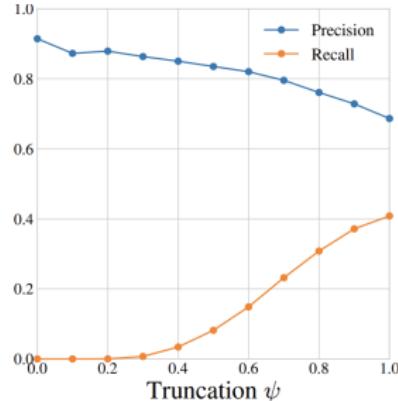
(b) Approx. manifold

Embed the samples using the pretrained network (as for FID).

# Precision-Recall



$\psi = 0.0 \quad \psi = 0.3 \quad \psi = 0.7 \quad \psi = 1.0$



## Truncation trick

BigGAN: truncated normal sampling

$$p(\mathbf{z}|\psi) = \frac{\mathcal{N}(\mathbf{z}|0, \mathbf{I})}{\int_{-\infty}^{\psi} \mathcal{N}(\mathbf{z}'|0, \mathbf{I}) d\mathbf{z}'}$$

Elements of  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$  which fall outside a predefined range are resampled.

StyleGAN

$$\mathbf{z}' = \hat{\mathbf{z}} + \psi \cdot (\mathbf{z} - \hat{\mathbf{z}}), \quad \hat{\mathbf{z}} = \mathbb{E}_{\mathbf{z}} \mathbf{z}$$

- ▶ Constant  $\psi$  is a tradeoff between diversity and fidelity.
- ▶  $\psi = 0.7$  is used for most of the results.

---

Brock A., Donahue J., Simonyan K. Large Scale GAN Training for High Fidelity Natural Image Synthesis, 2018

Karras T., Laine S., Aila T. A Style-Based Generator Architecture for Generative Adversarial Networks, 2018

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

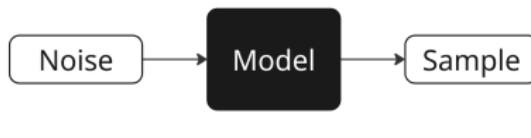
## 2. Langevin dynamic

## 3. Score matching

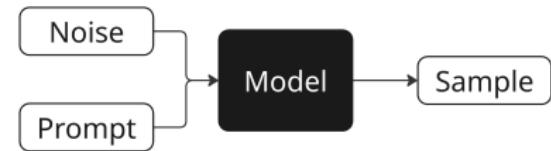
## 4. Denoising score matching

# CLIP score

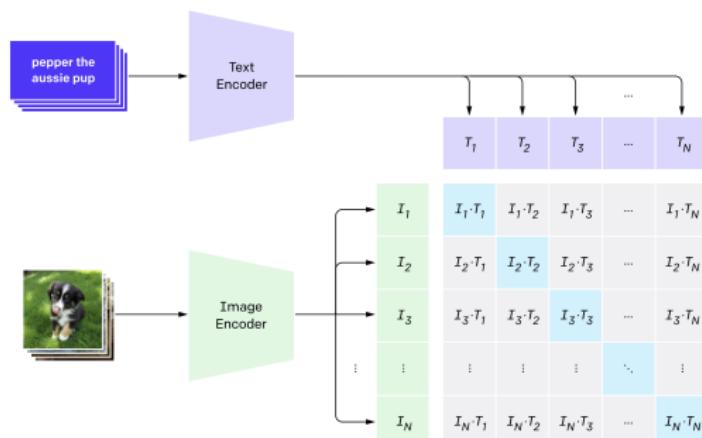
## Unconditional model



## Conditional model



We need the way to measure not only generated image quality, but also its relevance to the prompt.



# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

**Human evaluation**

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

# Human Evaluation

- ▶ There is no perfect automated metric.
- ▶ The best way to evaluate the generative model is to make human evaluation.
- ▶ It is essential to evaluate different aspects.

Аспект	Yandex ART 2.0	Mj 6.1	Mj 6	Ideogram	Recraft	Google Imagen3	Dall-E 3	FLUX	SBER Kandi3.1
Релевантность	<b>0,59</b>	<b>0,58</b>	<b>0,63</b>	<b>0,45</b>	<b>0,51</b>	<b>0,50</b>	<b>0,50</b>	<b>0,54</b>	<b>0,75</b>
Эстетика	<b>0,49</b>	<b>0,55</b>	<b>0,55</b>	<b>0,51</b>	<b>0,51</b>	<b>0,61</b>	<b>0,61</b>	<b>0,54</b>	<b>0,59</b>
Комплексность	<b>0,44</b>	<b>0,73</b>	<b>0,70</b>	<b>0,68</b>	<b>0,76</b>	<b>0,75</b>	<b>0,75</b>	<b>0,71</b>	<b>0,74</b>
Дефектность	<b>0,69</b>	<b>0,57</b>	<b>0,68</b>	<b>0,55</b>	<b>0,59</b>	<b>0,63</b>	<b>0,63</b>	<b>0,50</b>	<b>0,75</b>
Предпочтение	<b>0,66</b>	<b>0,60</b>	<b>0,69</b>	<b>0,49</b>	<b>0,54</b>	<b>0,63</b>	<b>0,63</b>	<b>0,51</b>	<b>0,84</b>

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

# Energy-based models

## Unnormalized density

$$p(\mathbf{x}|\theta) = \frac{\hat{p}(\mathbf{x}|\theta)}{Z_\theta}, \quad \text{where } Z_\theta = \int \hat{p}(\mathbf{x}|\theta) d\mathbf{x}$$

- ▶  $\hat{p}(\mathbf{x}|\theta)$  is any non-negative function.
- ▶ If we use the reparametrization  $\hat{p}(\mathbf{x}|\theta) = \exp(-f_\theta(\mathbf{x}))$ , we remove the non-negativite constraint.

## Unnormalized density

The gradient of the normalized log-density equals to the gradient of the unnormalized log-density:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta) = \nabla_{\mathbf{x}} \log \hat{p}(\mathbf{x}|\theta) - \nabla_{\mathbf{x}} \log Z_\theta = \nabla_{\mathbf{x}} \log \hat{p}(\mathbf{x}|\theta)$$

- ▶ Let assume that we already have the density (normalized or unnormalized)  $p(\mathbf{x}|\theta)$ .
- ▶ How to sample from the model?

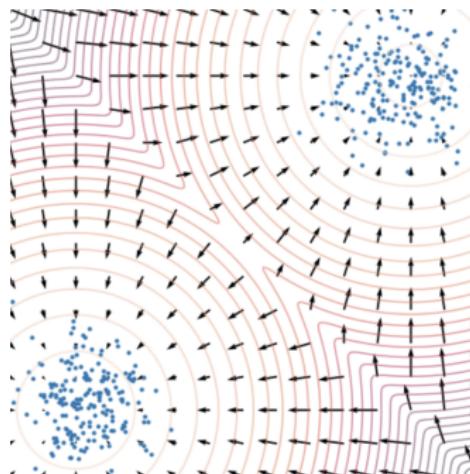
# Langevin dynamic

## Theorem (informal)

Let  $\mathbf{x}_0$  be a random vector. Under some mild regularity conditions samples from the following dynamics will come from  $p(\mathbf{x}|\theta)$  (for small enough  $\eta$  and large enough  $I$ )

$$\mathbf{x}_{I+1} = \mathbf{x}_I + \frac{\eta}{2} \cdot \nabla_{\mathbf{x}_I} \log p(\mathbf{x}_I | \theta) + \sqrt{\eta} \cdot \boldsymbol{\epsilon}_I, \quad \boldsymbol{\epsilon}_I \sim \mathcal{N}(0, \mathbf{I}).$$

- ▶ What do we get if  $\boldsymbol{\epsilon}_I = \mathbf{0}$ ?
- ▶ The density  $p(\mathbf{x}|\theta)$  is a **stationary** distribution for the Markov chain.
- ▶ We take the gradient w.r.t. to  $\mathbf{x}$  (not  $\theta$ ).
- ▶  $\nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta)$  defines the vector field.



# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

## Score matching

### Score function

$$\mathbf{s}_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta)$$

### Langevin dynamic

If we find the score function  $\mathbf{s}_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta)$  we will be able to sample from the model using Langevin dynamic.

$$\mathbf{x}_{I+1} = \mathbf{x}_I + \frac{\eta}{2} \cdot \nabla_{\mathbf{x}_I} \log p(\mathbf{x}_I|\theta) + \sqrt{\eta} \cdot \boldsymbol{\epsilon}_I = \mathbf{x}_I + \frac{\eta}{2} \cdot \mathbf{s}_\theta(\mathbf{x}_I) + \sqrt{\eta} \cdot \boldsymbol{\epsilon}_I.$$

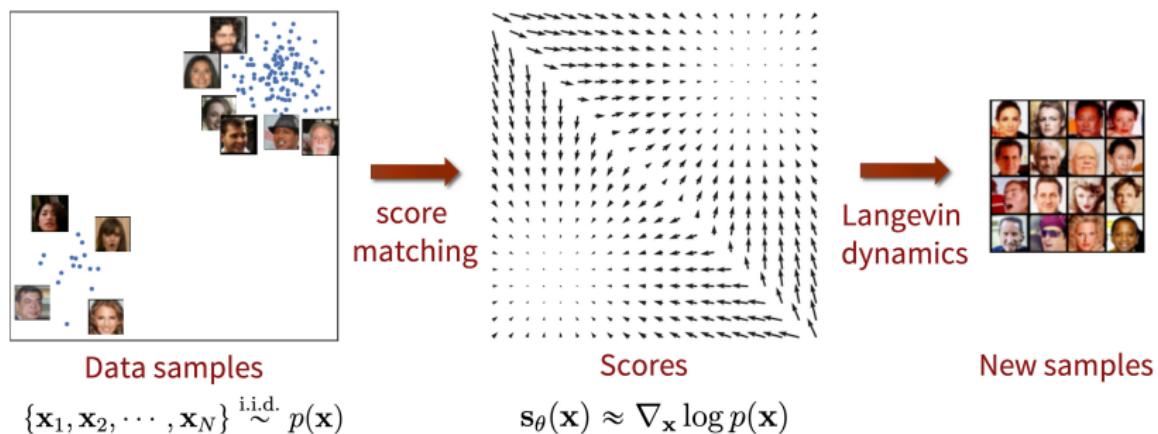
### Fisher divergence

$$\begin{aligned} D_F(\pi, p) &= \frac{1}{2} \mathbb{E}_\pi \left\| \nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta) - \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) \right\|_2^2 = \\ &= \frac{1}{2} \mathbb{E}_\pi \left\| \mathbf{s}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) \right\|_2^2 \rightarrow \min_\theta \end{aligned}$$

# Score matching

## Fisher divergence

$$D_F(\pi, p) = \frac{1}{2} \mathbb{E}_\pi \| \mathbf{s}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) \|_2^2 \rightarrow \min_{\theta}$$



**Problem:** We do not know  $\nabla_{\mathbf{x}} \log \pi(\mathbf{x})$ .

# Outline

## 1. Evaluation of likelihood-free models

Frechet Inception Distance (FID)

Precision-Recall

CLIP score

Human evaluation

## 2. Langevin dynamic

## 3. Score matching

## 4. Denoising score matching

## Denoising score matching

Let perturb original data  $\mathbf{x} \sim \pi(\mathbf{x})$  by random normal noise

$$\mathbf{x}_\sigma = \mathbf{x} + \sigma \cdot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), \quad q(\mathbf{x}_\sigma | \mathbf{x}) = \mathcal{N}(\mathbf{x}, \sigma^2 \cdot \mathbf{I})$$

$$q(\mathbf{x}_\sigma) = \int q(\mathbf{x}_\sigma | \mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}.$$

### Assumption

The solution of

$$\frac{1}{2} \mathbb{E}_{q(\mathbf{x}_\sigma)} \| \mathbf{s}_{\theta, \sigma}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma) \|_2^2 \rightarrow \min_{\theta}$$

satisfies  $\mathbf{s}_{\theta, \sigma}(\mathbf{x}_\sigma) \approx \mathbf{s}_{\theta, 0}(\mathbf{x}_0) = \mathbf{s}_\theta(\mathbf{x})$  if  $\sigma$  is small enough.

- ▶ The score function of the noised data is almost the same as the score function of the original data.
- ▶ Score function  $\mathbf{s}_{\theta, \sigma}(\mathbf{x}_\sigma)$  parametrized by  $\sigma$ .
- ▶ **Note:** We don't know  $q(\mathbf{x}_\sigma)$ , just like  $\pi(\mathbf{x})$ .

# Denoising score matching

## Theorem

$$\begin{aligned}\mathbb{E}_{q(\mathbf{x}_\sigma)} \|\mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma)\|_2^2 &= \\ &= \mathbb{E}_{\pi(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_\sigma|\mathbf{x})} \|\mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma|\mathbf{x})\|_2^2 + \text{const}(\theta)\end{aligned}$$

## Gradient of the noise kernel

$$\mathbf{x}_\sigma = \mathbf{x} + \sigma \cdot \boldsymbol{\epsilon}, \quad q(\mathbf{x}_\sigma|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \sigma^2 \cdot \mathbf{I})$$

$$\nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma|\mathbf{x}) = -\frac{\mathbf{x}_\sigma - \mathbf{x}}{\sigma^2} = -\frac{\boldsymbol{\epsilon}}{\sigma}$$

- ▶ The RHS does not need to compute  $\nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma)$  and even  $\nabla_{\mathbf{x}_\sigma} \log \pi(\mathbf{x}_\sigma)$ .
- ▶  $\mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma)$  tries to **denoise** the noised samples  $\mathbf{x}_\sigma$ .

# Denoising score matching

Initial objective:

$$\mathbb{E}_{\pi(\mathbf{x})} \left\| \mathbf{s}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) \right\|_2^2 \rightarrow \min_{\theta}$$

Noised objective:

$$\mathbb{E}_{q(\mathbf{x}_\sigma)} \left\| \mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}} \log q(\mathbf{x}_\sigma) \right\|_2^2 \rightarrow \min_{\theta}$$

This is equivalent to denoising task

$$\mathbb{E}_{\pi(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_\sigma|\mathbf{x})} \left\| \mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma|\mathbf{x}) \right\|_2^2 \rightarrow \min_{\theta}$$

$$\mathbb{E}_{\pi(\mathbf{x})} \mathbb{E}_{\mathcal{N}(0, \mathbf{I})} \left\| \mathbf{s}_{\theta,\sigma}(\mathbf{x} + \sigma \cdot \boldsymbol{\epsilon}) + \frac{\boldsymbol{\epsilon}}{\sigma} \right\|_2^2 \rightarrow \min_{\theta}$$

## Langevin dynamic

$$\mathbf{x}_{I+1} = \mathbf{x}_I + \frac{\eta}{2} \cdot \mathbf{s}_{\theta,\sigma}(\mathbf{x}_I) + \sqrt{\eta} \cdot \boldsymbol{\epsilon}_I, \quad \boldsymbol{\epsilon}_I \sim \mathcal{N}(0, \mathbf{I}).$$

## Summary

- ▶ Frechet Inception Distance is the most popular metric for the implicit models evaluation.
- ▶ Precision-recall allow to select model that compromises the sample quality and the sample diversity.
- ▶ CLIP score is frequently used to measure text-to-image relevance.
- ▶ The best way to measure the generated image quality is to make human evaluation.
- ▶ Langevin dynamics allows to sample from the generative model using the gradient of the log-likelihood.
- ▶ Score matching proposes to minimize the Fisher divergence to get the score function.
- ▶ Denoising score matching minimizes the Fisher divergence on noisy samples. It allows to estimate the Fisher divergence using samples.