

# Deep Generative Models

## Lecture 7

Roman Isachenko

Moscow Institute of Physics and Technology  
Yandex School of Data Analysis

2025, Autumn

# Recap of Previous Lecture

## Likelihood-Free Learning

- ▶ Likelihood isn't a perfect metric for generative models.
- ▶ Likelihood may 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 (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

The generative distribution  $p(\mathbf{x}|\theta)$  matches the true distribution  $\pi(\mathbf{x})$  if we can't distinguish between them using a discriminative model  $p(y|\mathbf{x})$ .

- ▶ **Generator:** a generative model  $\mathbf{x} = \mathbf{G}(\mathbf{z})$  that produces more realistic samples.
- ▶ **Discriminator:** a classifier  $D(\mathbf{x}) \in [0, 1]$  distinguishing real 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 a global optimum at  $\pi(\mathbf{x}) = p(\mathbf{x}|\theta)$ , and then  $D^*(\mathbf{x}) = 0.5$ .

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

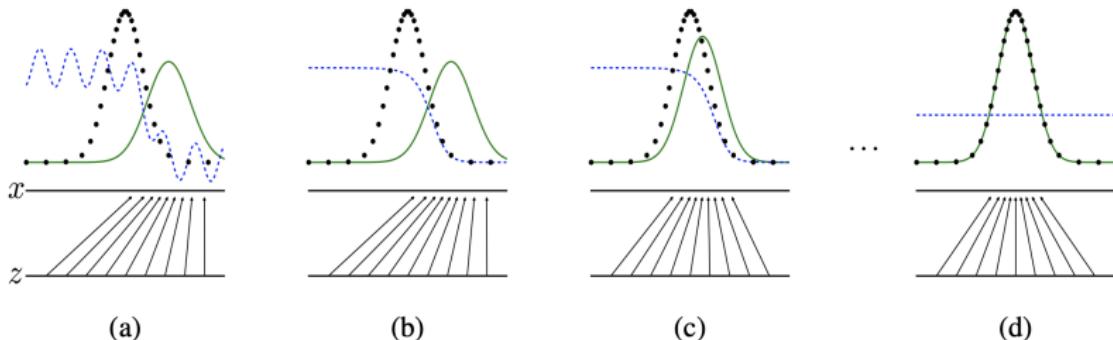
If the generator can be **any** function and the discriminator is **optimal** at each step, then the generator is **guaranteed to converge** to the data distribution.

## Recap of Previous Lecture

- ▶ The generator is updated in the parameter space; the discriminator isn't optimal at every iteration.
- ▶ Both generator and discriminator loss typically oscillate 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 Issues With Standard GANs

- ▶ Vanishing gradients (solution: non-saturating GAN).
- ▶ Mode collapse (arises from Jensen-Shannon divergence).

## Standard GAN

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

## Informal Theoretical Results

Both the data distribution  $\pi(x)$  and the generative distribution  $p(x|\theta)$  are low-dimensional with disjoint supports. In such cases,

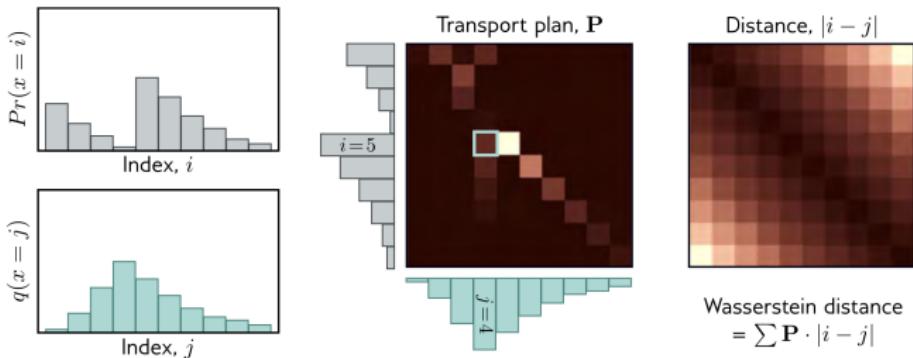
$$\text{KL}(\pi\|p) = \text{KL}(p\|\pi) = \infty, \quad \text{JSD}(\pi\|p) = \log 2.$$

---

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 (amount of "dirt" to transport from  $\mathbf{x}$  to  $\mathbf{y}$ ).
- ▶  $\Gamma(\pi, p)$  – 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  $\|f\|_L \leq K$  denotes  $K$ -Lipschitz continuous functions.

## WGAN Objective

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

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

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

# Outline

## 1. Evaluation of Likelihood-Free Models

Frechet Inception Distance (FID)

Precision-Recall

CLIP Score

Human Evaluation

## 2. Langevin Dynamics

## 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 Dynamics

## 3. Score Matching

## 4. Denoising Score Matching

# Evaluation of Likelihood-Free Models

## Likelihood-Based Models

- ▶ **Train:** fit the model.
- ▶ **Validation:** tune hyperparameters.
- ▶ **Test:** assess generalization by reporting likelihood.

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

## Desirable Properties for Samples

- ▶ Sharpness



- ▶ Diversity



# Outline

## 1. Evaluation of Likelihood-Free Models

Frechet Inception Distance (FID)

Precision-Recall

CLIP Score

Human Evaluation

## 2. Langevin Dynamics

## 3. Score Matching

## 4. Denoising Score Matching

# Wasserstein Metric

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

## Wasserstein 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

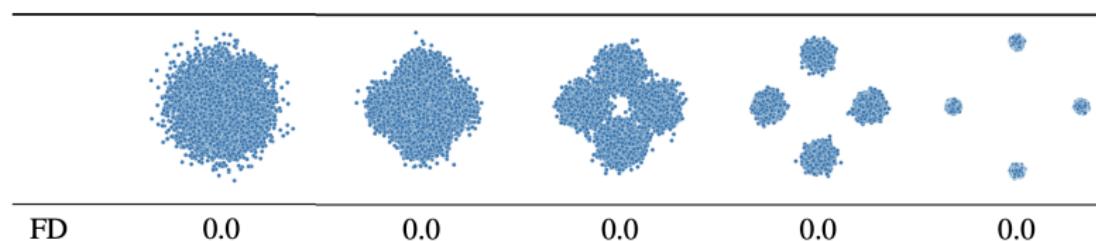
$$\text{FID}(\pi, p) = W_2^2(\pi, p)$$

## 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 computed in the latent space  $\mathbf{z}$ .
- ▶ We use a pretrained image embedder to get latent representations  $\mathbf{z} = \mathbf{f}(\mathbf{x})$ .
- ▶  $\boldsymbol{\mu}_\pi$ ,  $\boldsymbol{\Sigma}_\pi$  and  $\boldsymbol{\mu}_p$ ,  $\boldsymbol{\Sigma}_p$  are statistics of latent representations for 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

- ▶ Depends on the pretrained classification network.
- ▶ Uses the normality assumption.
- ▶ May not correlate with 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 Dynamics

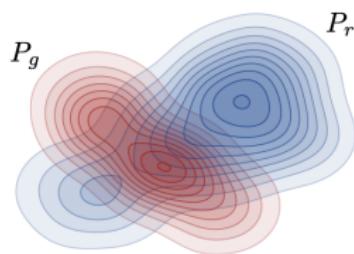
## 3. Score Matching

## 4. Denoising Score Matching

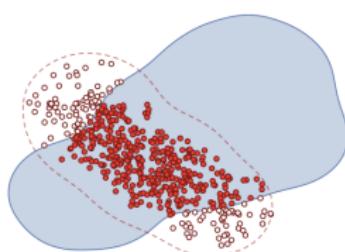
# Precision-Recall

## Desirable Properties for Generated Samples

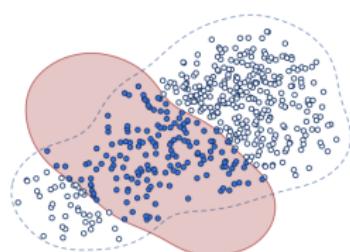
- ▶ **Sharpness:** generated samples should possess high visual quality.
- ▶ **Diversity:** their variation should match that in the training data.



(a) Example distributions



(b) Precision



(c) Recall

- ▶ **Precision** denotes the fraction of generated images that look realistic.
- ▶ **Recall** measures how well the generator covers the training data manifold.

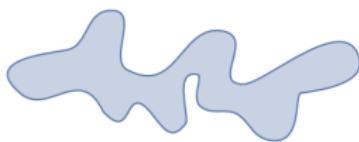
## Precision-Recall

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

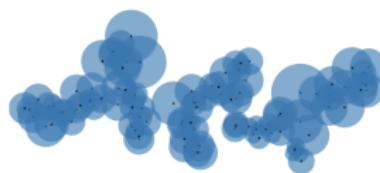
Define a binary function:

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

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



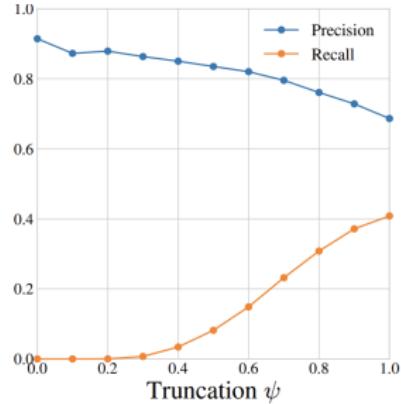
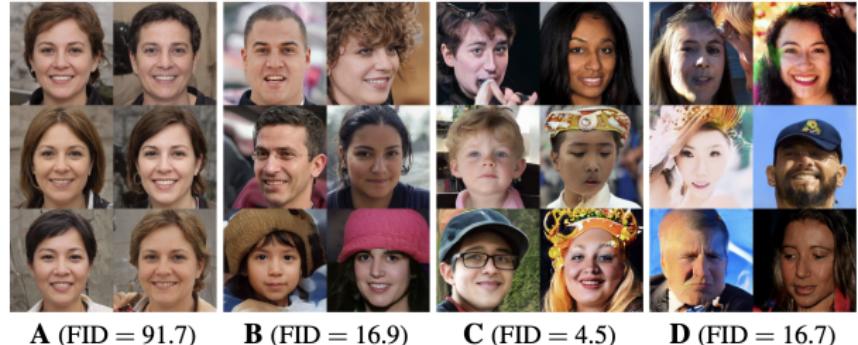
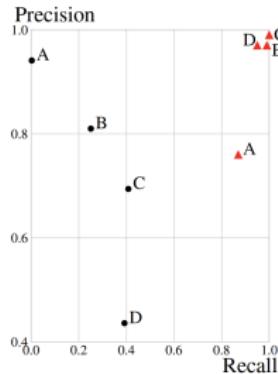
(a) True manifold



(b) Approx. manifold

Embed the samples using a pretrained network (as in FID).

# Precision-Recall



Kynkäanniemi T. et al. Improved precision and recall metric for assessing generative models, 2019

# Outline

## 1. Evaluation of Likelihood-Free Models

Frechet Inception Distance (FID)

Precision-Recall

**CLIP Score**

Human Evaluation

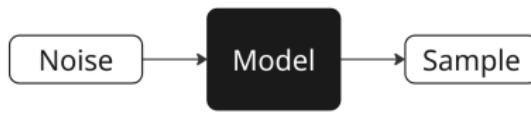
## 2. Langevin Dynamics

## 3. Score Matching

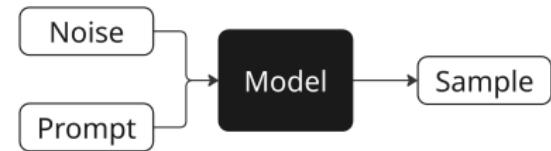
## 4. Denoising Score Matching

# CLIP Score

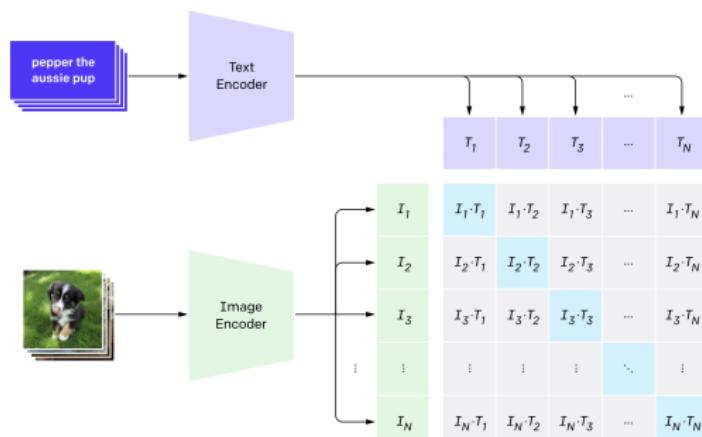
## Unconditional Model



## Conditional Model



We need a way to measure not only the quality of the generated image, but also how well it's aligned with the prompt.



# Outline

## 1. Evaluation of Likelihood-Free Models

Frechet Inception Distance (FID)

Precision-Recall

CLIP Score

Human Evaluation

## 2. Langevin Dynamics

## 3. Score Matching

## 4. Denoising Score Matching

# Human Evaluation

- ▶ No automated metric is perfect.
- ▶ The best way to evaluate generative models is by human assessment.
- ▶ It's important to assess various properties.

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 Dynamics

## 3. Score Matching

## 4. Denoising Score Matching

# Energy-Based Models

## Unnormalized Density

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

- ▶  $\hat{p}(\mathbf{x}|\theta)$  can be any non-negative function.
- ▶ If we reparameterize as  $\hat{p}(\mathbf{x}|\theta) = \exp(-f_\theta(\mathbf{x}))$ , we eliminate the non-negativity constraint.

## Unnormalized Density

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

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

- ▶ Suppose we already have this density (normalized or not)  $p(\mathbf{x}|\theta)$ .
- ▶ How can we sample from the model?

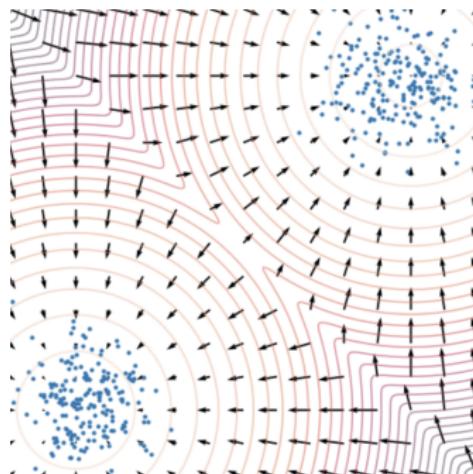
# Langevin Dynamics

## Theorem (Informal)

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

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

- ▶ What if  $\boldsymbol{\epsilon}_I = \mathbf{0}$ ?
- ▶ The density  $p(\mathbf{x}|\theta)$  is the **stationary** distribution of the Markov chain.
- ▶ The gradient is taken with respect to  $\mathbf{x}$ , not  $\theta$ .
- ▶  $\nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta)$  defines a vector field.



# Outline

## 1. Evaluation of Likelihood-Free Models

Frechet Inception Distance (FID)

Precision-Recall

CLIP Score

Human Evaluation

## 2. Langevin Dynamics

## 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 Dynamics

If we know the score function  $\mathbf{s}_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta)$ , we can generate samples from the model using Langevin dynamics:

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

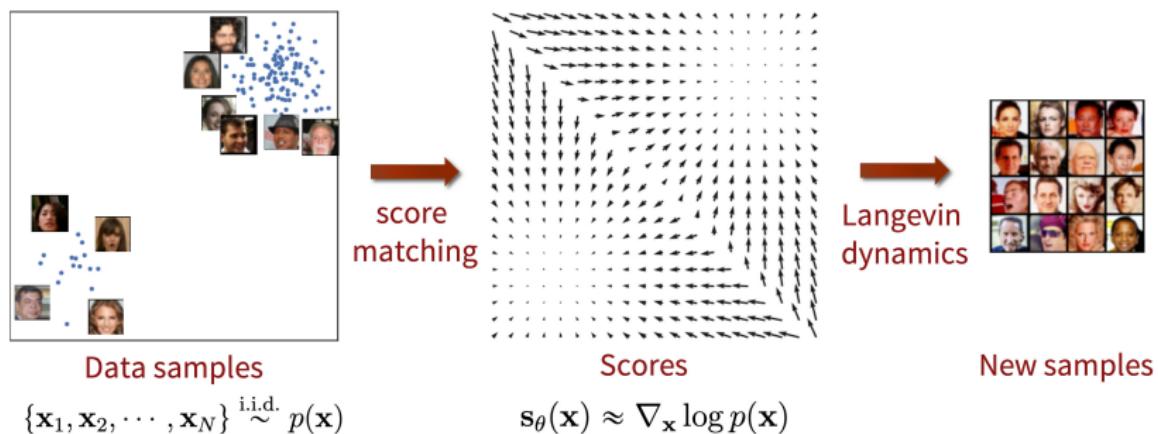
## Fisher Divergence

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

# 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 don't 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 Dynamics

## 3. Score Matching

## 4. Denoising Score Matching

## Denoising Score Matching

Let us perturb the original data  $\mathbf{x} \sim \pi(\mathbf{x})$  with Gaussian 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 to

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

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

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

# 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 right-hand side doesn't require computing  $\nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma)$  or even  $\nabla_{\mathbf{x}_\sigma} \log \pi(\mathbf{x}_\sigma)$ .
- ▶  $\mathbf{s}_{\theta,\sigma}(\mathbf{x}_\sigma)$  is trained 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 a 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 Dynamics

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

## Summary

- ▶ Frechet Inception Distance is the most popular metric for evaluating implicit generative models.
- ▶ Precision-recall allows for choosing a model that balances sample quality and diversity.
- ▶ The CLIP score is widely used to measure text-to-image alignment.
- ▶ The gold standard for evaluating generated image quality is human assessment.
- ▶ Langevin dynamics enable sampling from generative models using gradients of the log-likelihood.
- ▶ Score matching proposes minimizing Fisher divergence to estimate the score function.
- ▶ Denoising score matching optimizes Fisher divergence on noisy data, making it estimable with samples.