

COMP547

DEEP UNSUPERVISED LEARNING

Lecture #9 – Energy and Score-Based Models



**KOÇ
UNIVERSITY**

Aykut Erdem // Koç University // Spring 2026

Previously on COMP547

- Motivation
- Original GAN (Goodfellow et al, 2014)
- Evaluation: Parzen, Inception, Fréchet
- Theory of GANs
- GAN Progression
- Conditional GANs, Cycle-Consistent Adversarial Networks
- Applications



Lecture overview

- Energy-based models
- Score-based Models

Disclaimer: Much of the material and slides for this lecture were borrowed from
—Stefano Ermon and Aditya Grover's Stanford CS236 class
—Yang Song and Stefano Ermon's talk titled "Generative Modeling by Estimating Gradients of the Data Distribution"

Lecture overview

- Motivation
- Energy-based models
- Score-based Models

Inspired by their work, we further investigated the relationship between diffusion models and score-based generative models in an ICLR 2021 paper [20]. We found that the sampling method of diffusion probabilistic models can be integrated with annealed Langevin dynamics of score-based models to create a unified and more powerful sampler (the Predictor-Corrector sampler). By generalizing the number of noise scales to infinity, we further proved that score-based generative models and diffusion probabilistic models can both be viewed as discretizations to stochastic differential equations determined by score functions. This work bridges both score-based generative modeling and diffusion probabilistic modeling into a unified framework.

Collectively, these latest developments seem to indicate that both score-based generative modeling with multiple noise perturbations and diffusion probabilistic models are different perspectives of the same model family, much like how wave mechanics and matrix mechanics are equivalent formulations of quantum mechanics in the history of physics⁵. The perspective of score matching and score-based models allows one to calculate log-likelihoods exactly, solve inverse problems naturally, and is directly connected to energy-based models, Schrödinger bridges and optimal transport [47]. The perspective of diffusion models is naturally connected to VAEs, lossy compression, and can be directly incorporated with variational probabilistic inference. This blog post focuses on the first perspective, but I highly recommend interested readers to learn about the alternative perspective of diffusion models as well (see a great blog by Lilian Weng).

Disclaimer: Much of the material and slides for this lecture were borrowed from

—Stefano Ermon and Aditya Grover's Stanford CS236 class

<https://yang-song.github.io/blog/2021/score/>

—Yang Song and Stefano Ermon's talk titled "Generative Modeling by Estimating Gradients of the Data Distribution"

Lecture overview

- Energy-based models
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - Ising Model, Product of Experts, Restricted Boltzmann machine, Deep Boltzman Machines
 - Training and sampling from EBMs
- Score-based Models

Parameterizing probability distributions

- Probability distributions $p(x)$ are a key building block in generative modeling.
 1. Non-negative: $p(x) \geq 0$
 2. Sum-to-one: $\sum_x p(x) = 1$ (or $\int p(x)dx = 1$ for continuous variables)
- Coming up with a non-negative function $p_\theta(\mathbf{x})$ is not hard. Given any function $f_\theta(\mathbf{x})$, we can choose
 - $g_\theta(\mathbf{x}) = f_\theta(\mathbf{x})^2$
 - $g_\theta(\mathbf{x}) = \exp(f_\theta(\mathbf{x}))$
 - $g_\theta(\mathbf{x}) = |f_\theta(\mathbf{x})|$
 - $g_\theta(\mathbf{x}) = \log(1 + \exp(f_\theta(\mathbf{x})))$

Parameterizing probability distributions

- Probability distributions $p(x)$ are a key building block in generative modeling.
 1. Non-negative: $p(x) \geq 0$
 2. Sum-to-one: $\sum_x p(x) = 1$ (or $\int p(x)dx = 1$ for continuous variables)
Sum-to-one is key!
- Total “volume” is fixed: increasing $p(x_{train})$ guarantees that x_{train} becomes relatively more likely (compared to the rest)
- Problem:
 - $g_\theta(\mathbf{x}) \geq 0$ is easy, but $g_\theta(\mathbf{x})$ might not sum-to-one.
 - $\sum_x g_\theta(\mathbf{x}) = Z(\theta) \neq 1$ in general, so $g_\theta(\mathbf{x})$ is not a valid probability mass function or density

Parameterizing probability distributions

- **Problem:** $g_{\theta}(\mathbf{x}) \geq 0$ is easy, but $g_{\theta}(\mathbf{x})$ might not be normalized
- **Solution:** $p_{\theta}(\mathbf{x}) = \frac{1}{\text{Volume}(g_{\theta})} g_{\theta}(\mathbf{x}) = \frac{1}{\int g_{\theta}(\mathbf{x}) d\mathbf{x}} g_{\theta}(\mathbf{x})$

Then by definition $\int p_{\theta}(\mathbf{x}) d\mathbf{x} = 1$

- **Example:** Choose $g_{\theta}(\mathbf{x})$ so that the volume is **analytically** as a function of θ .
 1. $g_{(\mu, \sigma)}(x) = e^{-\frac{(x-\mu)^2}{2\sigma^2}}$. Volume is: $\int e^{-\frac{x-\mu}{2\sigma^2}} dx = \sqrt{2\pi\sigma^2} \rightarrow$ **Gaussian**
 2. $g_{\lambda}(x) = e^{-\lambda x}$. Volume is: $\int_0^{+\infty} e^{-\lambda x} dx = \frac{1}{\lambda} \rightarrow$ **Exponential**
 3. $g_{\theta}(x) = h(x) \exp\{\theta \cdot T(x)\}$. Volume is $\exp\{A(\theta)\}$, where $A(\theta) = \log \int h(x) \exp\{\theta \cdot T(x)\} d\mathbf{x} \rightarrow$ **Exponential family**
 - Normal, Poisson, exponential, Bernoulli
 - beta, gamma, Dirichlet, Wishart, etc.
- Function forms $g_{\theta}(\mathbf{x})$ need to allow **analytical** integration. Despite being restrictive, they are very useful as building blocks for more complex distributions.

Lecture overview

- **Energy-based models**
 - Parametrizing probability distributions
 - **Energy-based generative modeling**
 - Ising Model, Product of Experts, Restricted Boltzmann machine, Deep Boltzman Machines
 - Training and sampling from EBMs
- Score-based Models

Likelihood based learning

- **Problem:** $g_{\theta}(\mathbf{x}) \geq 0$ is easy, but $g_{\theta}(\mathbf{x})$ might not be normalized
- **Solution:** $p_{\theta}(\mathbf{x}) = \frac{1}{\text{Volume}(g_{\theta})} g_{\theta}(\mathbf{x}) = \frac{1}{\int g_{\theta}(\mathbf{x}) d\mathbf{x}} g_{\theta}(\mathbf{x})$

Typically, choose $g_{\theta}(\mathbf{x})$ so that we know the volume analytically. More complex models can be obtained by combining these building blocks.

1. **Autoregressive:** Products of normalized objects $p_{\theta}(\mathbf{x})p_{\theta'(\mathbf{x})}(\mathbf{y})$

$$\int_{\mathbf{x}} \int_{\mathbf{y}} p_{\theta}(\mathbf{x}) p_{\theta'(\mathbf{x})}(\mathbf{y}) d\mathbf{x} d\mathbf{y} = \int_{\mathbf{x}} p_{\theta}(\mathbf{x}) \underbrace{\int_{\mathbf{y}} p_{\theta'(\mathbf{x})}(\mathbf{y}) d\mathbf{y}}_{=1} d\mathbf{x} = \int_{\mathbf{x}} p_{\theta}(\mathbf{x}) d\mathbf{x} = 1$$

2. **Latent variables:** Mixtures of normalized objects $\alpha p_{\theta}(\mathbf{x}) + (1 - \alpha) p_{\theta'}(\mathbf{x})$
 $\int_{\mathbf{x}} \alpha p_{\theta}(\mathbf{x}) + (1 - \alpha) p_{\theta'}(\mathbf{x}) d\mathbf{x} = \alpha + (1 - \alpha) = 1$

Likelihood based learning

- **Problem:** $g_{\theta}(\mathbf{x}) \geq 0$ is easy, but $g_{\theta}(\mathbf{x})$ might not be normalized
- **Solution:** $p_{\theta}(\mathbf{x}) = \frac{1}{\text{Volume}(g_{\theta})} g_{\theta}(\mathbf{x}) = \frac{1}{\int g_{\theta}(\mathbf{x}) d\mathbf{x}} g_{\theta}(\mathbf{x})$

Typically, choose $g_{\theta}(\mathbf{x})$ so that we know the volume analytically. More complex models can be obtained by combining these building blocks.

1. **Autoregressive:** Products of normalized objects $p_{\theta}(\mathbf{x})p_{\theta'(\mathbf{x})}(\mathbf{y})$

$$\int_{\mathbf{x}} \int_{\mathbf{y}} p_{\theta}(\mathbf{x}) p_{\theta'(\mathbf{x})}(\mathbf{y}) d\mathbf{x} d\mathbf{y} = \int_{\mathbf{x}} p_{\theta}(\mathbf{x}) \underbrace{\int_{\mathbf{y}} p_{\theta'(\mathbf{x})}(\mathbf{y}) d\mathbf{y}}_{=1} d\mathbf{x} = \int_{\mathbf{x}} p_{\theta}(\mathbf{x}) d\mathbf{x} = 1$$

2. **Latent variables:** Mixtures of normalized objects $\alpha p_{\theta}(\mathbf{x}) + (1 - \alpha) p_{\theta'}(\mathbf{x})$
 $\int_{\mathbf{x}} \alpha p_{\theta}(\mathbf{x}) + (1 - \alpha) p_{\theta'}(\mathbf{x}) d\mathbf{x} = \alpha + (1 - \alpha) = 1$

How about using models where the “volume”/normalization constant of \tilde{p} is not easy to compute analytically?

Energy-based model

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

The output of f_{θ} is a scalar value between $-\infty$ and ∞ .

- The volume/normalization constant

$$Z(\theta) = \int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}$$

is also called the **partition** function. Why exponential (and not e.g. $f_{\theta}(\mathbf{x})^2$)?

1. Want to capture very large variations in probability. log-probability is the natural scale we want to work with. Otherwise need highly non-smooth f_{θ} .
2. Exponential families. Many common distributions can be written in this form.
3. These distributions arise under fairly general assumptions in statistical physics (maximum entropy, second law of thermodynamics).
 - $-f_{\theta}(\mathbf{x})$ is called the **energy**, hence the name.
 - Intuitively, configurations \mathbf{x} with low energy (high $f_{\theta}(\mathbf{x})$) are more likely.

Energy-based model

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- Pros:
 - extreme flexibility: can use pretty much any function $f_{\theta}(\mathbf{x})$ you want
- Cons:
 - Sampling from $p_{\theta}(\mathbf{x})$ is hard
 - Evaluating and optimizing likelihood $p_{\theta}(\mathbf{x})$ is hard (learning is hard)
 - No feature learning (but can add latent variables)
- **Curse of dimensionality:** The fundamental issue is that computing $Z(\theta)$ numerically (when no analytic solution is available) scales exponentially in the number of dimensions of \mathbf{x} .
- Nevertheless, some tasks do not require knowing $Z(\theta)$

Applications of Energy-based models

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

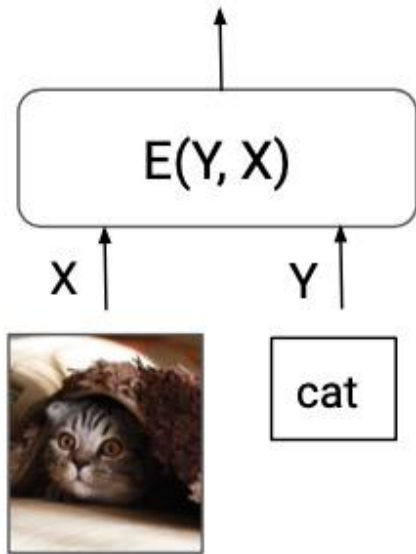
- Given \mathbf{x} , \mathbf{x}' evaluating $p_{\theta}(\mathbf{x})$ or $p_{\theta}(\mathbf{x}')$ requires $Z(\theta)$.
- However, their ratio

$$\frac{p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x}')} = \exp(f_{\theta}(\mathbf{x}) - f_{\theta}(\mathbf{x}'))$$

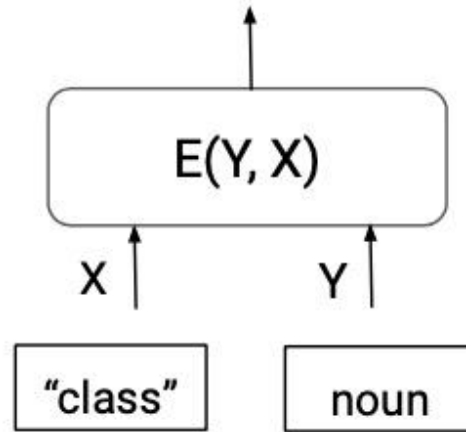
does not involve $Z(\theta)$.

- This means we can easily check which one is more likely. Applications:
 - Anomaly detection
 - Denoising

Applications of Energy-based models



object recognition



sequence labeling

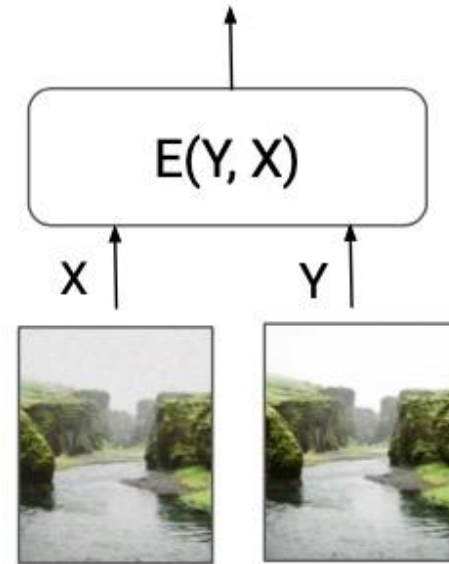


image restoration

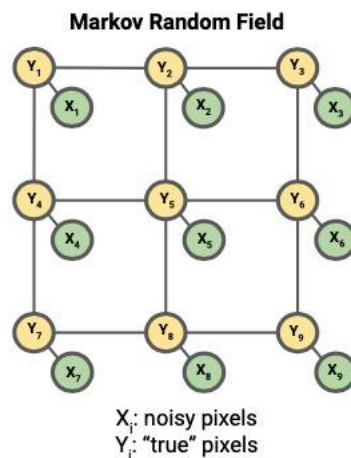
- Given a trained model, many applications require relative comparisons. Hence $Z(\theta)$ is not needed.

Lecture overview

- **Energy-based models**
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - **Ising Model**, Product of Experts, Restricted Boltzmann machine, Deep Boltzman Machines
 - Training and sampling from EBMs
- Score-based Models

Ising Model

- There is a true image $\mathbf{y} \in \{0, 1\}^{3 \times 3}$, and a corrupted image $\mathbf{x} \in \{0, 1\}^{3 \times 3}$. We know \mathbf{x} , and want to somehow recover \mathbf{y} .



- We model the joint probability distribution $p(\mathbf{y}, \mathbf{x})$ as

$$p(\mathbf{y}, \mathbf{x}) = \frac{1}{Z} \exp \left(\sum_i \psi_i (x_i, y_i) + \sum_{(i,j) \in E} \psi_{ij} (y_i, y_j) \right)$$

- $\psi_i (x_i, y_i)$: the i -th corrupted pixel depends on the i -th original pixel
- $\psi_{ij} (y_i, y_j)$: neighboring pixels tend to have the same value
- How did the original image \mathbf{y} look like? Solution: maximize $p(\mathbf{y}|\mathbf{x})$. Or equivalently, maximize $p(\mathbf{y}, \mathbf{x})$.

Lecture overview

- **Energy-based models**
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - Ising Model, **Product of Experts**, Restricted Boltzmann machine, Deep Boltzman Machines
 - Training and sampling from EBMs
- Score-Based Models

Product of Experts

- Suppose you have trained several models $q_{\theta_1}(\mathbf{x}), r_{\theta_2}(\mathbf{x}), t_{\theta_3}(\mathbf{x})$. They can be different models (PixelCNN, Flow, etc.)
- Each one is like an expert that can be used to score how likely an input \mathbf{x} is.
- Assuming the experts make their judgments independently, it is tempting to ensemble them as

$$p_{\theta_1}(\mathbf{x})q_{\theta_2}(\mathbf{x})r_{\theta_3}(\mathbf{x})$$

- To get a valid probability distribution, we need to normalize

$$p_{\theta_1, \theta_2, \theta_3}(\mathbf{x}) = \frac{1}{Z(\theta_1, \theta_2, \theta_3)} q_{\theta_1}(\mathbf{x}) r_{\theta_2}(\mathbf{x}) t_{\theta_3}(\mathbf{x})$$

- Note: similar to an AND operation (e.g., probability is zero as long as one model gives zero probability), unlike mixture models which behave more like OR

Product of Experts

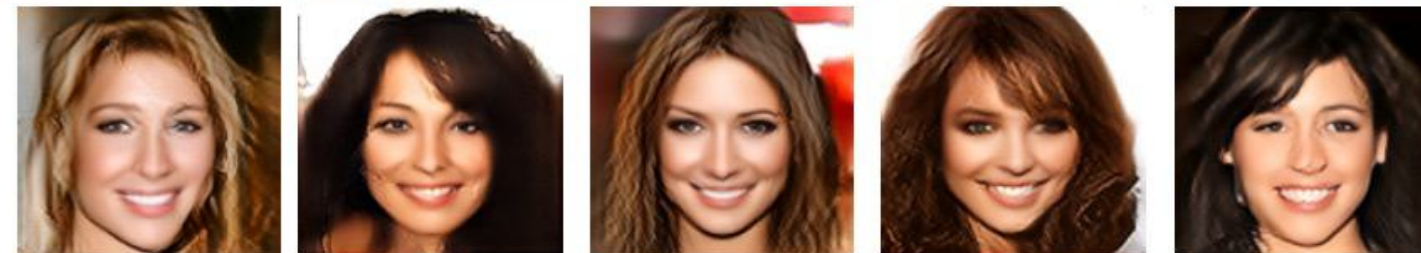
Young
(EBM)



Young
AND
Female
(EBM)



Young
AND Female
AND Smiling
(EBM)



Young
AND Female
AND Smiling
AND Wavy Hair
(EBM)



Lecture overview

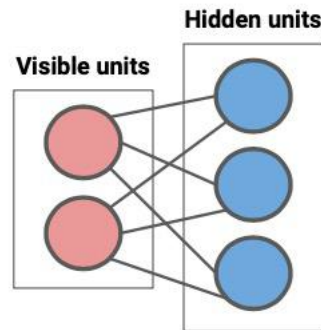
- **Energy-based models**
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - Ising Model, Product of Experts, **Restricted Boltzmann machine**, Deep Boltzman Machines
 - Training and sampling from EBMs
- Score-based Models

Restricted Boltzmann machine (RBM)

- RBM: energy-based model with latent variables
- Two types of variables:
 - $\mathbf{x} \in \{0, 1\}^n$ are visible variables (e.g., pixel values)
 - $\mathbf{z} \in \{0, 1\}^m$ are latent ones

- The joint distribution is

$$p_{W,b,c}(\mathbf{x}, \mathbf{z}) = \frac{1}{Z} \exp(\mathbf{x}^T W \mathbf{z} + b\mathbf{x} + c\mathbf{z}) = \frac{1}{Z} \exp\left(\sum^n \sum^m x_i z_j w_{ij} + b\mathbf{x} + c\mathbf{z}\right)$$



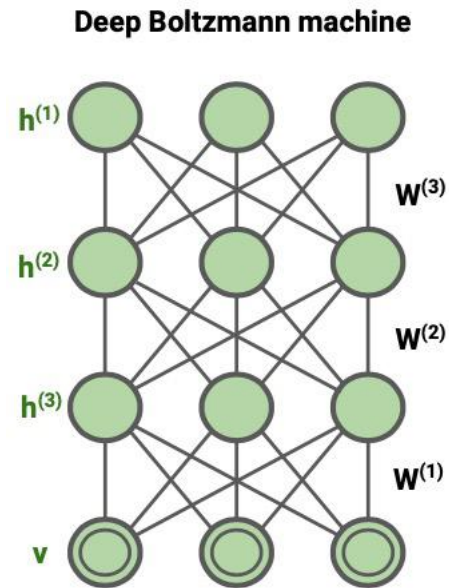
- Restricted because there are no visible-visible and hidden-hidden connections, i.e., $x_i x_j$ or $z_i z_j$ terms in the objective

Lecture overview

- **Energy-based models**
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - Ising Model, Product of Experts, Restricted Boltzmann machine, **Deep Boltzman Machines**
 - Training and sampling from EBMs
- Score-based Models

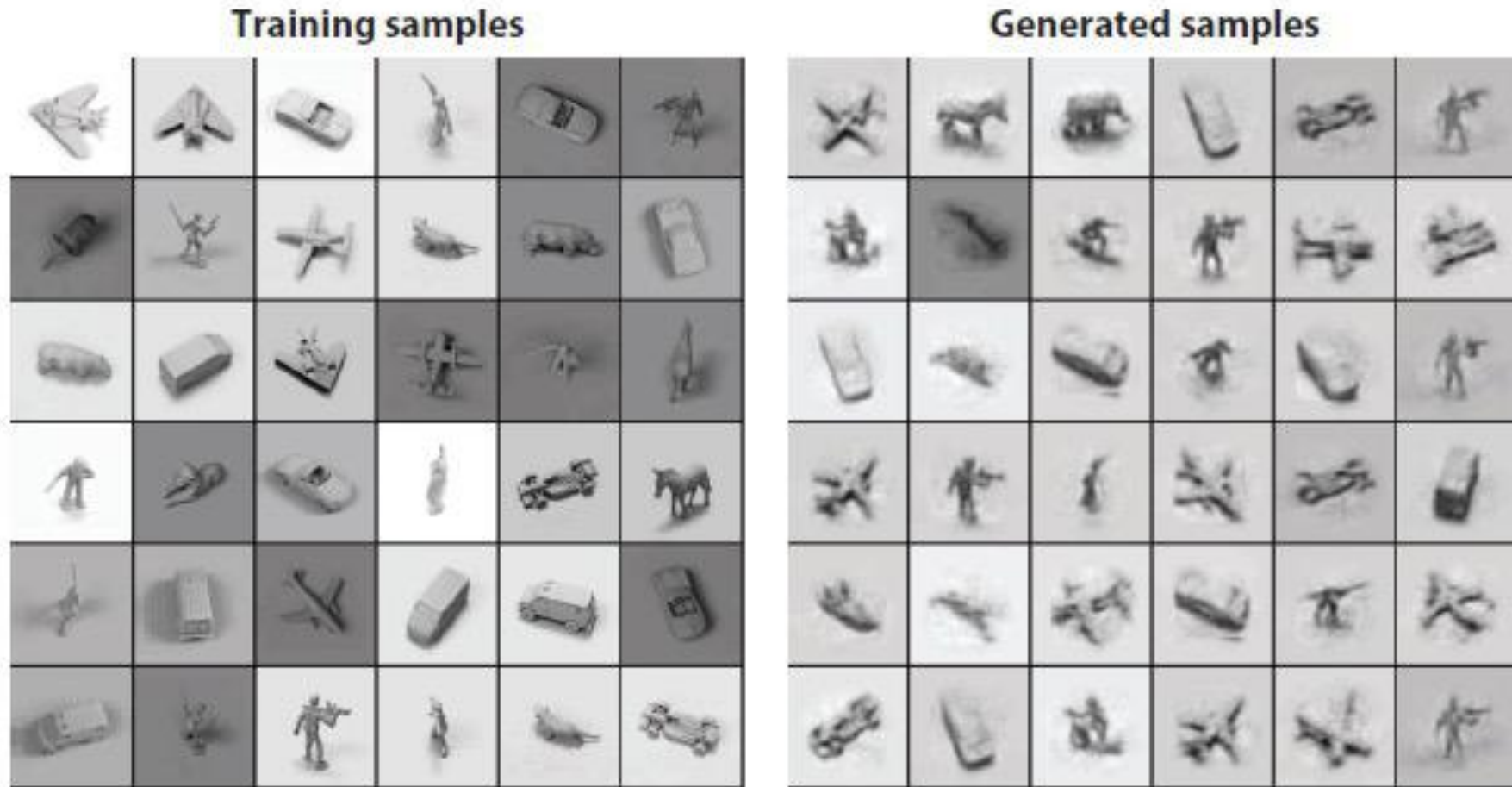
Deep Boltzmann Machines

- Stacked RBMs are one of the first deep generative models:



- Bottom layer variables v are pixel values. Layers above (h) represent “higher-level” features (corners, edges, etc.).
- Early deep neural networks for supervised learning had to be pre-trained like this to make them work.

Deep Boltzmann Machines: Samples



Lecture overview

- **Energy-based models**
 - Parametrizing probability distributions
 - Energy-based generative modeling
 - Ising Model, Product of Experts, Restricted Boltzmann machine, Deep Boltzman Machines
 - **Training and sampling from EBMs**
- Score-based Models

Recap: Energy-based models

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- Pros:

- can plug in pretty much any function $f_{\theta}(\mathbf{x})$ you want

- Cons:

- Sampling is hard

- Evaluating likelihood (learning) is hard

- No feature learning

- **Curse of dimensionality:** The fundamental issue is that computing $Z(\theta)$ numerically (when no analytic solution is available) scales exponentially in the number of dimensions of \mathbf{x}

Computing the normalization constant is hard

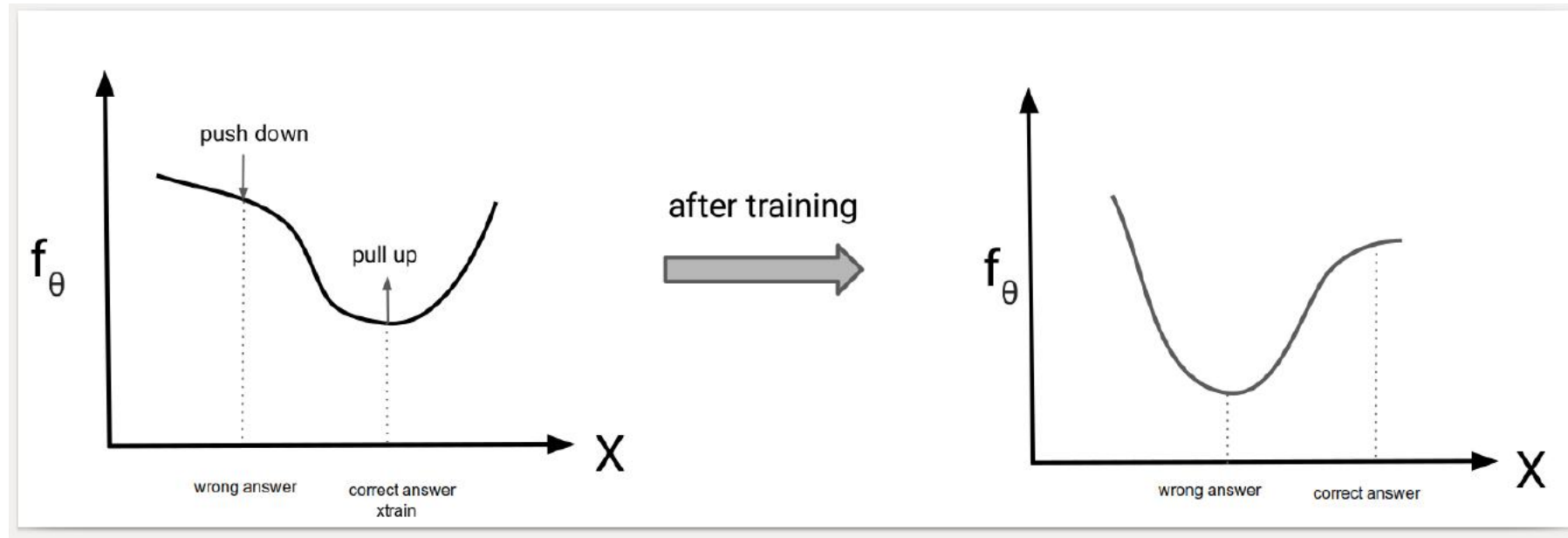
- As an example, the RBM joint distribution is

$$p_{W,b,c}(\mathbf{x}, \mathbf{z}) = \frac{1}{Z} \exp(\mathbf{x}^T W \mathbf{z} + b\mathbf{x} + c\mathbf{z})$$

where

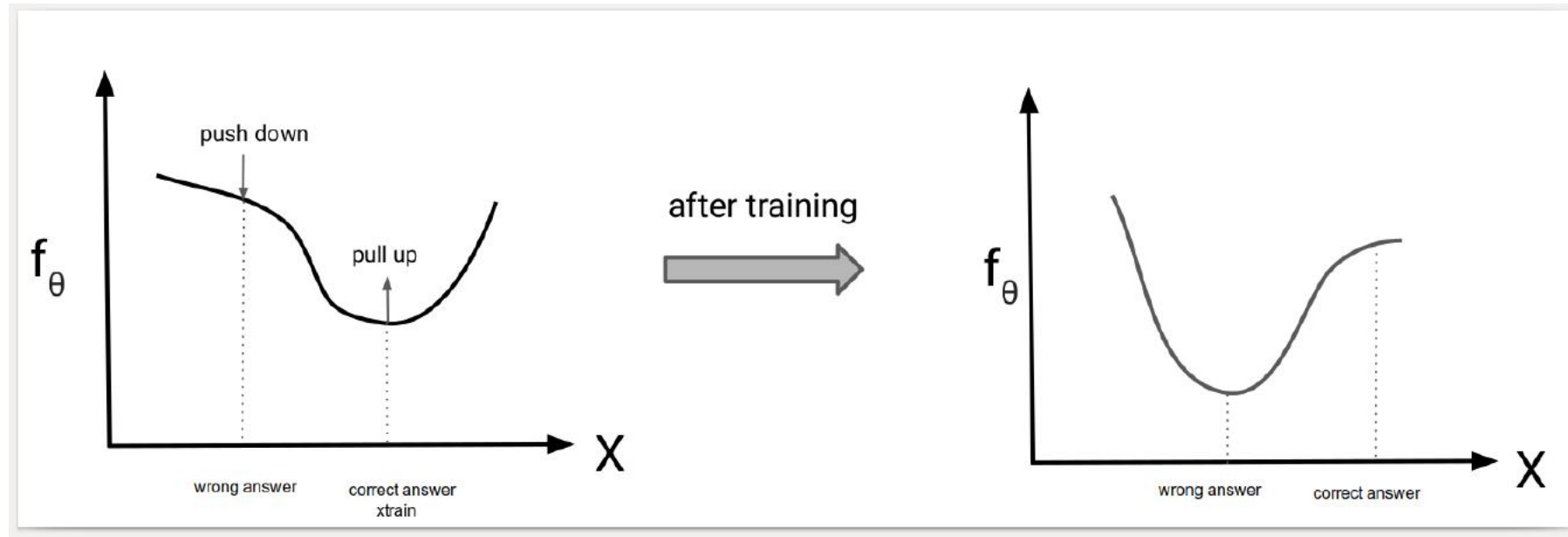
- $\mathbf{x} \in \{0, 1\}^n$ are visible variables (e.g., pixel values)
- $\mathbf{z} \in \{0, 1\}^m$ are latent ones
- The normalization constant (the “volume”) is
$$Z(W, b, c) = \sum_{\mathbf{x} \in \{0,1\}^n} \sum_{\mathbf{z} \in \{0,1\}^m} \exp(\mathbf{x}^T W \mathbf{z} + b\mathbf{x} + c\mathbf{z})$$
- **Note:** it is a well-defined function of the parameters W, b, c , but no simple closed-form. Takes time exponential in n, m to compute. This means that evaluating the objective function $p_{W,b,c}(\mathbf{x}, \mathbf{z})$ for likelihood-based learning is hard.
- **Observation:** Optimizing the likelihood $p_{W,b,c}(\mathbf{x}, \mathbf{z})$ is difficult, but optimizing the un-normalized probability $\exp(\mathbf{x}^T W \mathbf{z} + b\mathbf{x} + c\mathbf{z})$ (wrt trainable parameters W, b, c) is easy.

Training intuition



- **Goal:** maximize $\frac{\exp\{f_\theta(\mathbf{x}_{\text{train}})\}}{Z(\theta)}$. Increase numerator, decrease denominator.
- **Intuition:** because the model is not normalized, increasing the un-normalized log-probability $f_\theta(\mathbf{x}_{\text{train}})$ by changing θ does **not** guarantee that $\mathbf{x}_{\text{train}}$ becomes relatively more likely (compared to the rest).
- We also need to take into account the effect on other “wrong points” and try to “push them down” to also make $Z(\theta)$ small.

Contrastive Divergence



- Goal: maximize $\frac{\exp\{f_\theta(\mathbf{x}_{\text{train}})\}}{Z(\theta)}$
- Idea: Instead of evaluating $Z(\theta)$ exactly, use a Monte Carlo estimate.
- Contrastive divergence algorithm: sample $\mathbf{x}_{\text{sample}} \sim p_\theta$ take step on $\nabla_\theta (f_\theta(\mathbf{x}_{\text{train}}) - f_\theta(\mathbf{x}_{\text{sample}}))$. Make training data more likely than typical sample from the model.

Contrastive Divergence

- Maximize log-likelihood: $\max_{\theta} f_{\theta}(x_{\text{train}}) - \log Z(\theta)$
- Gradient of log-likelihood:

$$\begin{aligned} & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \nabla_{\theta} \log Z(\theta) \\ = & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \frac{\nabla_{\theta} Z(\theta)}{Z(\theta)} \\ = & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \frac{1}{Z(\theta)} \int \nabla_{\theta} \exp\{f_{\theta}(x)\} dx \\ = & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \frac{1}{Z(\theta)} \int \exp\{f_{\theta}(x)\} \nabla_{\theta} f_{\theta}(x) dx \\ = & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \int \frac{\exp\{f_{\theta}(x)\}}{Z(\theta)} \nabla_{\theta} f_{\theta}(x) dx \\ = & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - E_{x_{\text{sample}}} [\nabla_{\theta} f_{\theta}(x_{\text{sample}})] \\ \approx & \nabla_{\theta} f_{\theta}(x_{\text{train}}) - \nabla_{\theta} f_{\theta}(x_{\text{sample}}), \end{aligned}$$

where $x_{\text{sample}} \sim \exp\{f_{\theta}(x_{\text{sample}})\} / Z(\theta)$

- How to sample?

Sampling from energy-based models

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- No direct way to sample like in autoregressive or flow models. Main issue: cannot easily compute how likely each possible sample is
- However, we can easily compare two samples x, x' .
- Use an iterative approach called Markov Chain Monte Carlo:
 1. Initialize x^0 randomly, $t = 0$
 2. Let $x' = x^t + \text{noise}$
 - If $f_{\theta}(x') > f_{\theta}(x^t)$, let $x^{t+1} = x'$
 - Else let $x^{t+1} = x'$ with probability $\exp(f_{\theta}(x') - f_{\theta}(x^t))$
 3. Goto step 2
- Works in theory, but can take a very long time to converge

Sampling from energy-based models

- For any continuous distribution $p_\theta(\mathbf{x})$, suppose we can compute its gradient (the **score function**) $\nabla_{\mathbf{x}} \log p_\theta(\mathbf{x})$
- Let $\pi(\mathbf{x})$ be a prior distribution that is easy to sample from.
- Langevin MCMC.
 - $\mathbf{x}^0 \sim \pi(\mathbf{x})$
 - Repeat $\mathbf{x}^{t+1} \sim \mathbf{x}^t + \epsilon \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}^t) + \sqrt{2\epsilon} \mathbf{z}^t$ for $t = 0, 1, 2, \dots, T - 1$, where $\mathbf{z}^t \sim \mathcal{N}(0, I)$.
 - If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, we have $\mathbf{x}_T \sim p_\theta(\mathbf{x})$.
- Note that for energy-based models

$$\begin{aligned} \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}) &= \nabla_{\mathbf{x}} f_\theta(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z(\theta)}_{=0} \\ &= \nabla_{\mathbf{x}} f_\theta(\mathbf{x}) \end{aligned}$$

Modern energy-based models



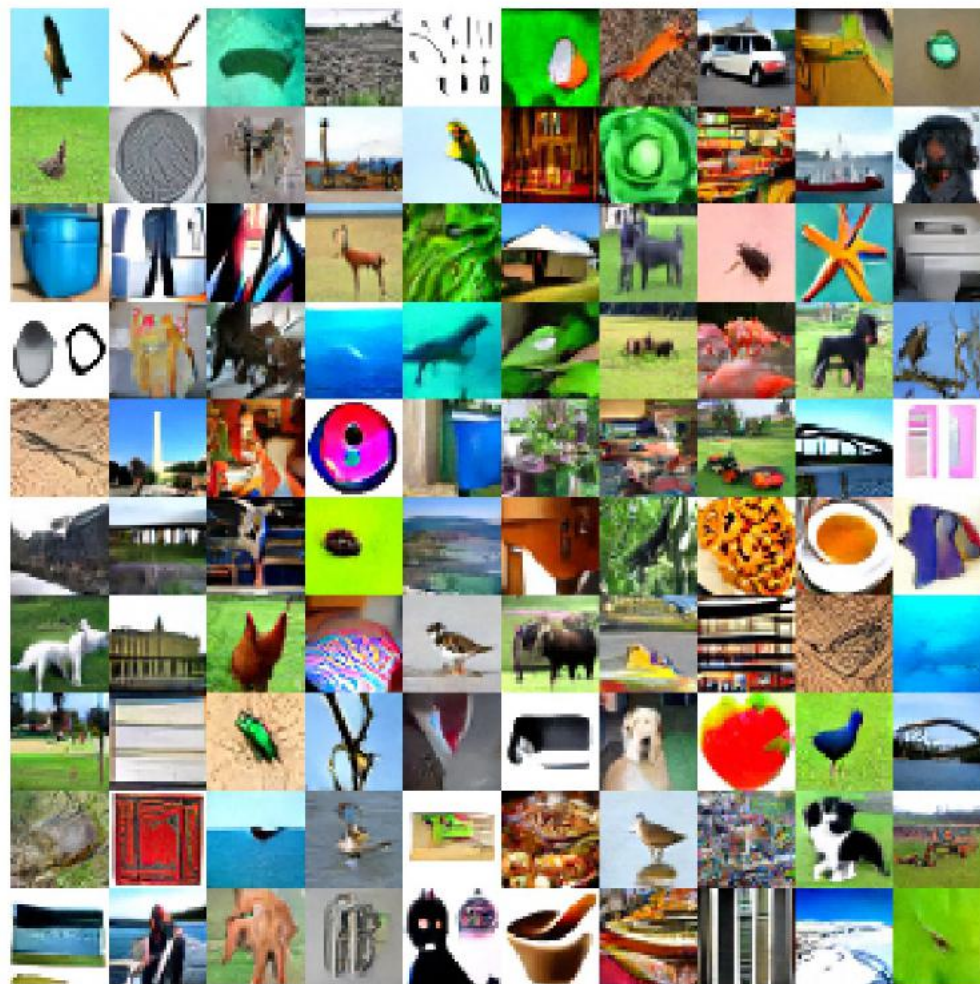
Langevin sampling



Face samples

Image source: Nijkamp et al. 2019

Modern energy-based models



ImageNet samples

Lecture overview

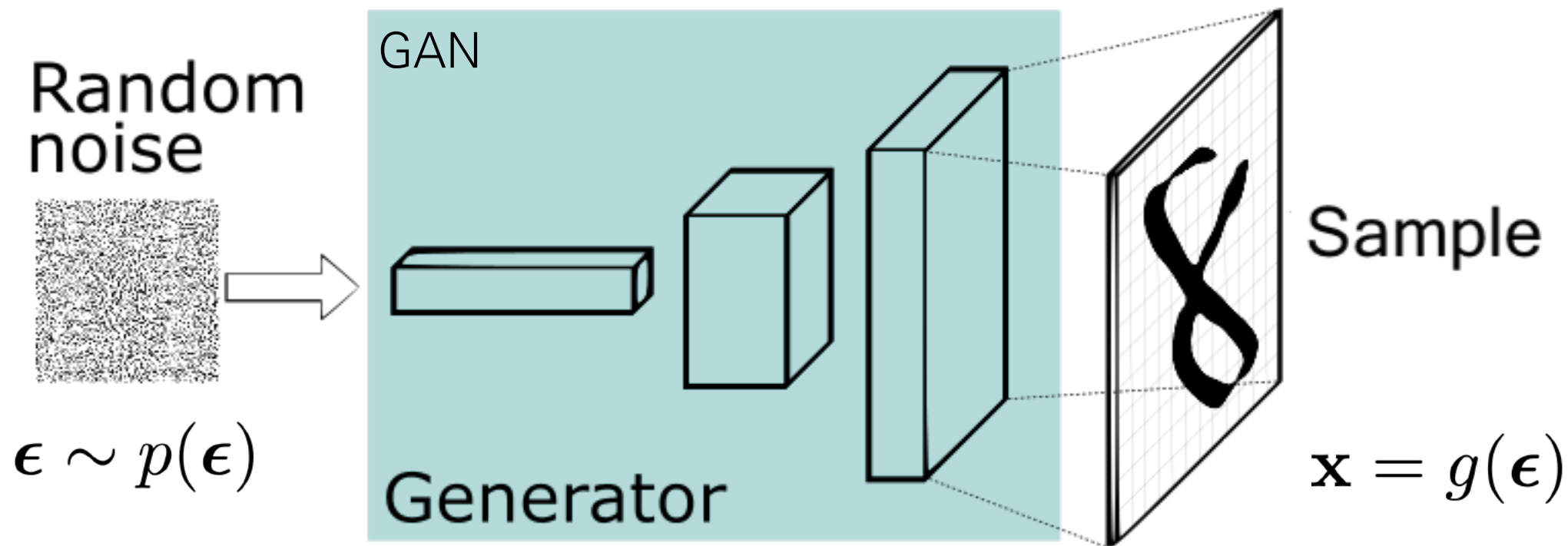
- Energy-based models
- **Score-based Models**
 - Probability distributions and Score functions
 - Score-based Generative Models
 - Sampling from a Score-based Model
 - Latent Score-based Generative Models

Lecture overview

- Energy-based models
- **Score-based Models**
 - Probability Distribution and Score functions
 - Score-based Generative Models
 - Sampling from a Score-based Model
 - Latent Score-based Generative Models

Representations of Probability Distributions

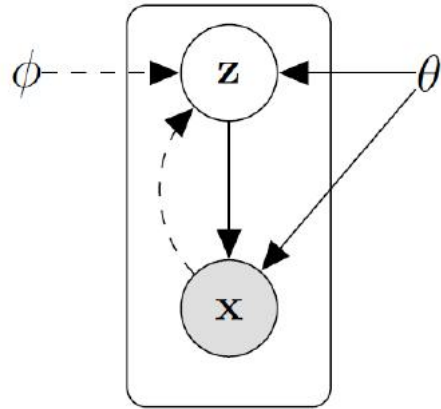
Implicit models: directly represent the sampling process



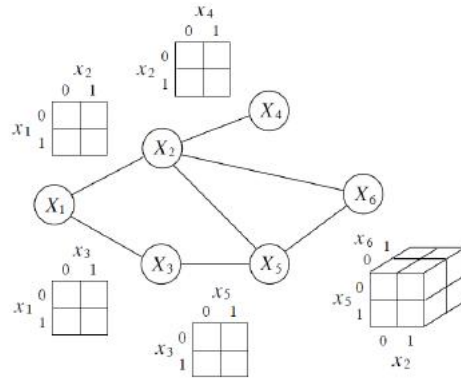
- **Cons:** hard to train, no likelihood, no principled model comparisons

Representations of Probability Distributions

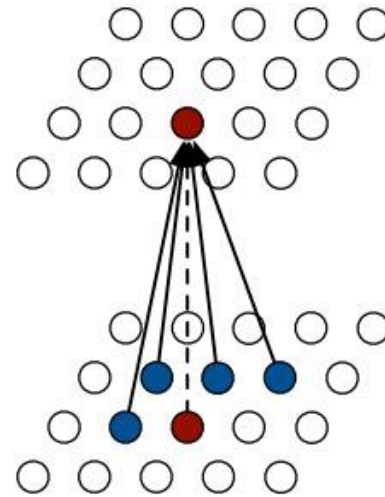
Explicit models: represent a probability density/mass function $p(\mathbf{x})$



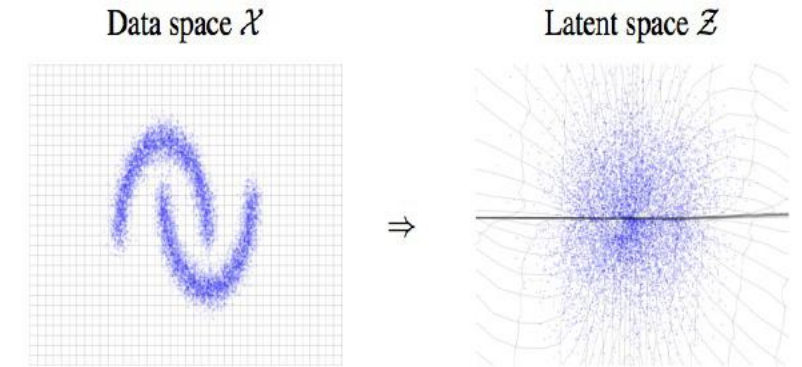
Bayesian networks
(e.g., VAEs)



MRF



Autoregressive
models



Flow models

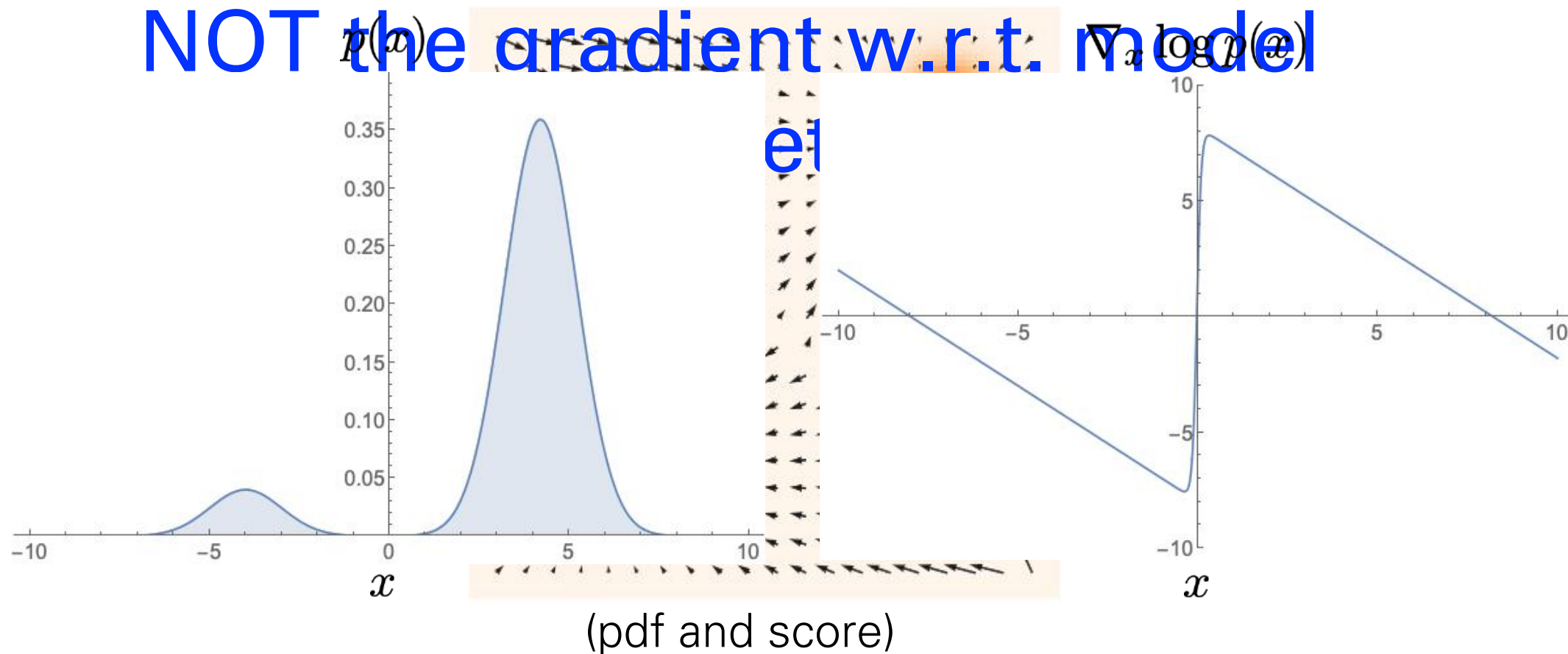
- **Cons:** need to be normalized \rightarrow balance expressivity and tractability

Representation of Probability Distributions

Alternative: The gradient of a probability density wrt the input dimensions

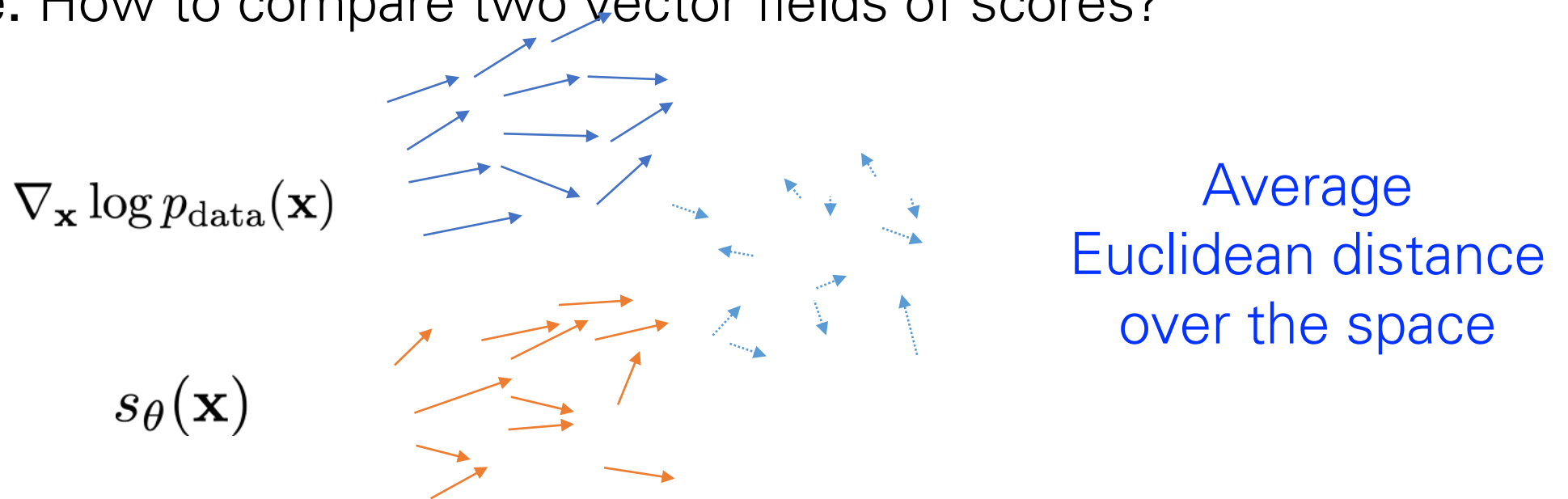
$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) \quad \text{Score}$$

NOT the gradient w.r.t. model



Score Estimation

- **Given:** i.i.d. samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(\mathbf{x})$
- **Task:** Estimating the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- **Score Model:** A trainable vector-valued function $s_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^D$
- **Objective:** How to compare two vector fields of scores?



Score Estimation

- **Given:** i.i.d. samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(\mathbf{x})$
- **Task:** Estimating the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- **Score Model:** A trainable vector-valued function $s_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^D$
- **Objective:** How to compare two vector fields of scores?

$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}} [\|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - s_{\theta}(\mathbf{x})\|_2^2]$$

(Fisher divergence)

- Integration by parts

$$\mathbb{E}_{p_{\text{data}}} \left[\frac{1}{2} \|s_{\theta}(\mathbf{x})\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x})) \right]$$

Score Matching
Hyvärinen (2005)

$$\approx \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{2} \|s_{\theta}(\mathbf{x}_i)\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}_i)) \right]$$

COMP547

DEEP UNSUPERVISED LEARNING

Lecture #9 – Energy and Score-Based Models



**KOÇ
UNIVERSITY**

Aykut Erdem // Koç University // Spring 2026

Recap: Energy-based models

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

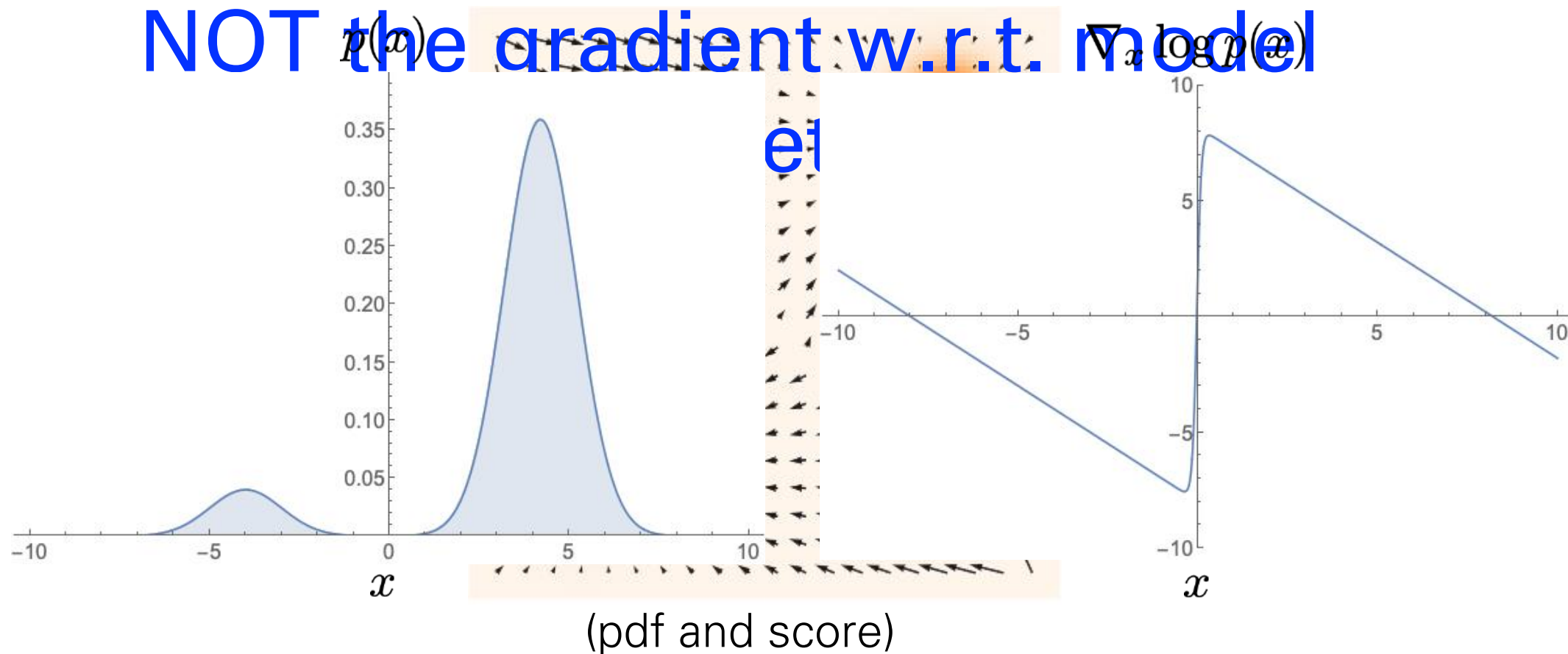
- Pros:
 - can plug in pretty much any function $f_{\theta}(\mathbf{x})$ you want
- Cons:
 - Sampling is hard
 - Evaluating likelihood (learning) is hard
 - No feature learning

Recap: Repr. of Probability Distributions

Alternative: The gradient of a probability density wrt the input dimensions

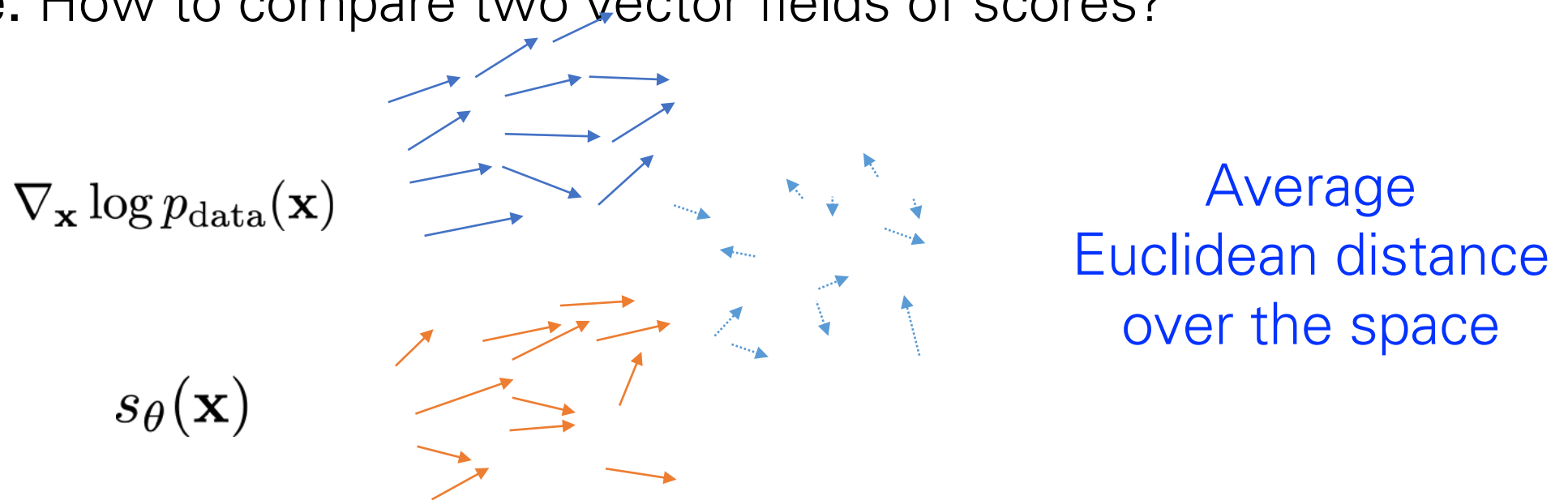
$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) \quad \text{Score}$$

NOT the gradient w.r.t. model



Recap: Score Estimation

- **Given:** i.i.d. samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(\mathbf{x})$
- **Task:** Estimating the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- **Score Model:** A trainable vector-valued function $s_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^D$
- **Objective:** How to compare two vector fields of scores?



Recap: Score Estimation

- **Given:** i.i.d. samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(\mathbf{x})$
- **Task:** Estimating the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$
- **Score Model:** A trainable vector-valued function $s_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^D$
- **Objective:** How to compare two vector fields of scores?

$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}} [\|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - s_{\theta}(\mathbf{x})\|_2^2]$$

(Fisher divergence)

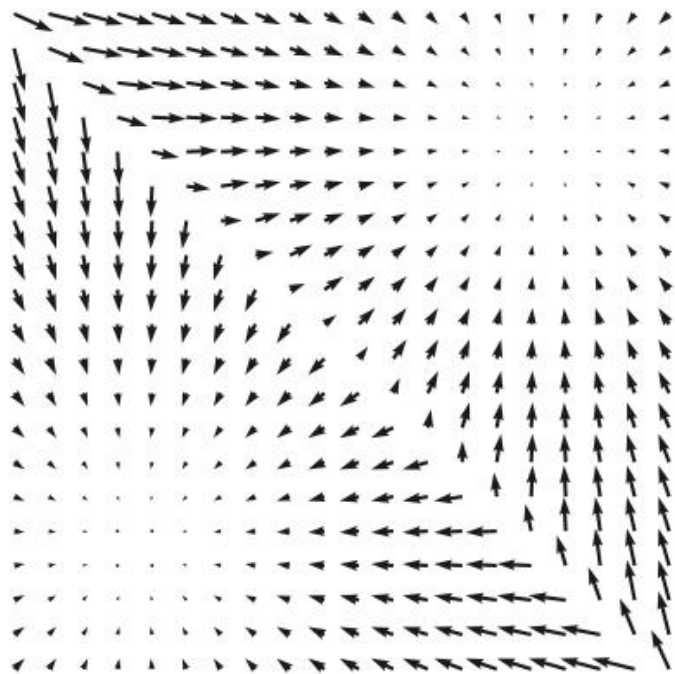
- Integration by parts

$$\mathbb{E}_{p_{\text{data}}} \left[\frac{1}{2} \|s_{\theta}(\mathbf{x})\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x})) \right]$$

Score Matching
Hyvärinen (2005)

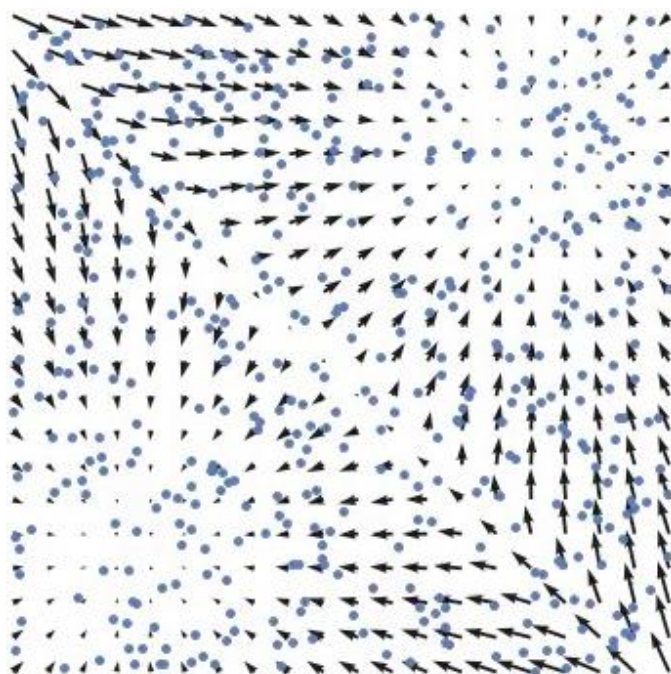
$$\approx \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{2} \|s_{\theta}(\mathbf{x}_i)\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}_i)) \right]$$

From Scores to Samples: Langevin Dynamics



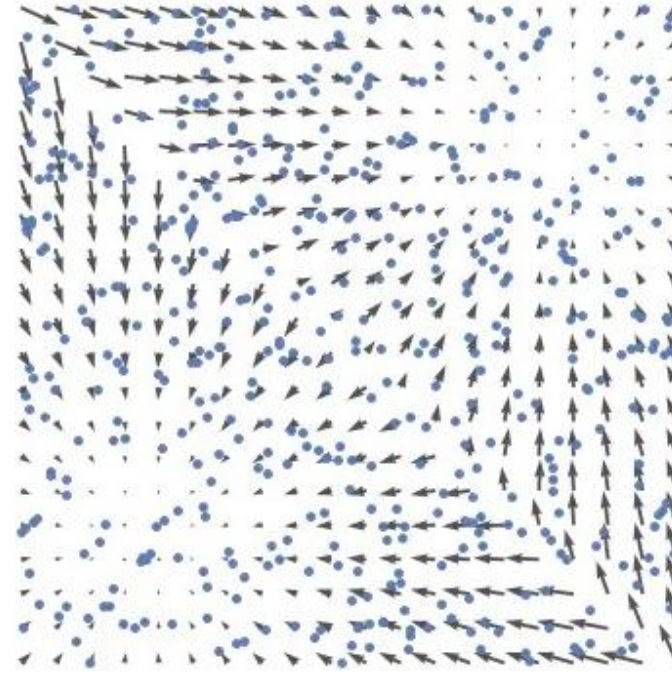
Scores

$$s_{\theta}(\mathbf{x})$$



Follow the scores

$$\tilde{\mathbf{x}}_{t+1} \leftarrow \tilde{\mathbf{x}}_t + \frac{\epsilon}{2} s_{\theta}(\tilde{\mathbf{x}}_t)$$



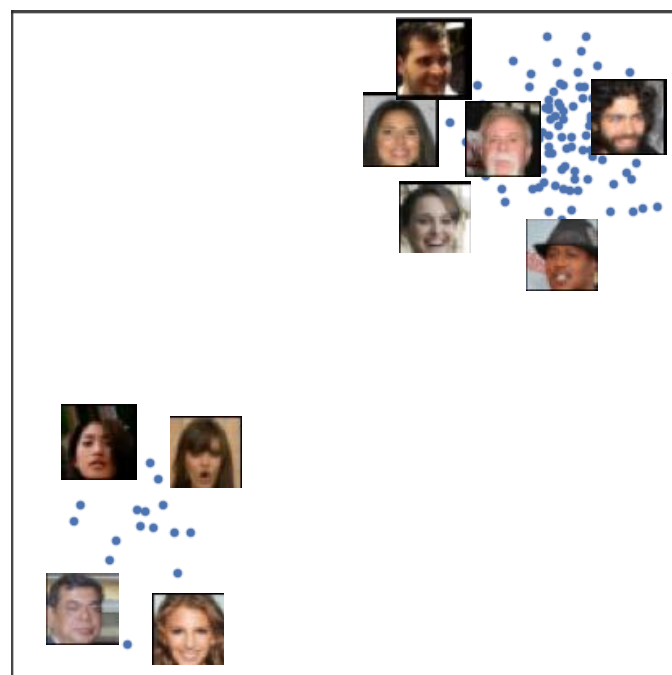
Follow noisy scores:
Langevin dynamics

$$\begin{aligned} \mathbf{z}_t &\sim \mathcal{N}(0, I) \\ \tilde{\mathbf{x}}_{t+1} &\leftarrow \tilde{\mathbf{x}}_t + \frac{\epsilon}{2} s_{\theta}(\tilde{\mathbf{x}}_t) + \sqrt{\epsilon} \mathbf{z}_t \end{aligned}$$

Lecture overview

- Energy-based models
- **Score-based Models**
 - Probability Distribution and Score functions
 - **Score-based Generative Models**
 - Sampling from a Score-based Model
 - Latent Score-based Generative Models

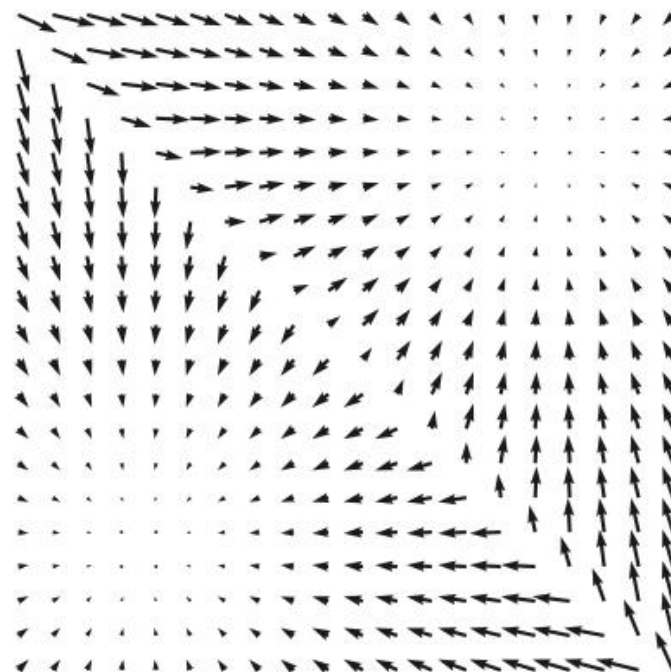
Score-Based Generative Modeling



Data samples

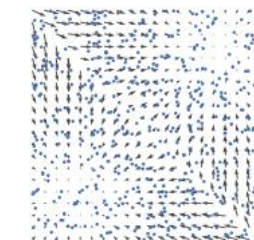
$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(\mathbf{x})$$

Score Matching



Scores

$$s_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$$



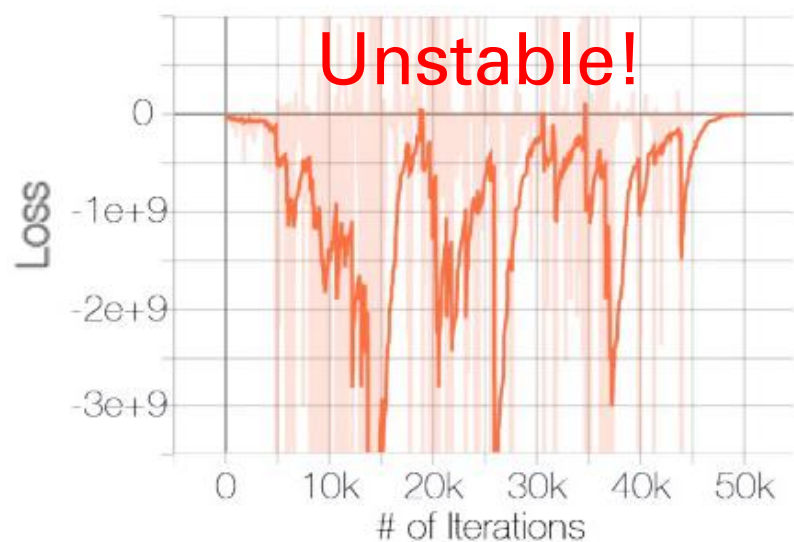
Langevin dynamics



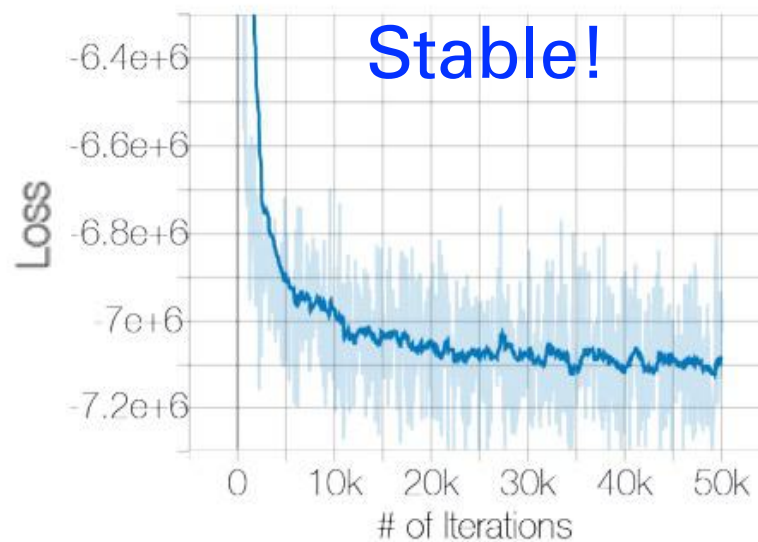
New samples

Adding Noise to Data for Well-Defined Scores

- Scores can be undefined when
 - The support of data distribution is on a low-dimensional manifold
 - The data distribution is discrete
- Solution: adding noise

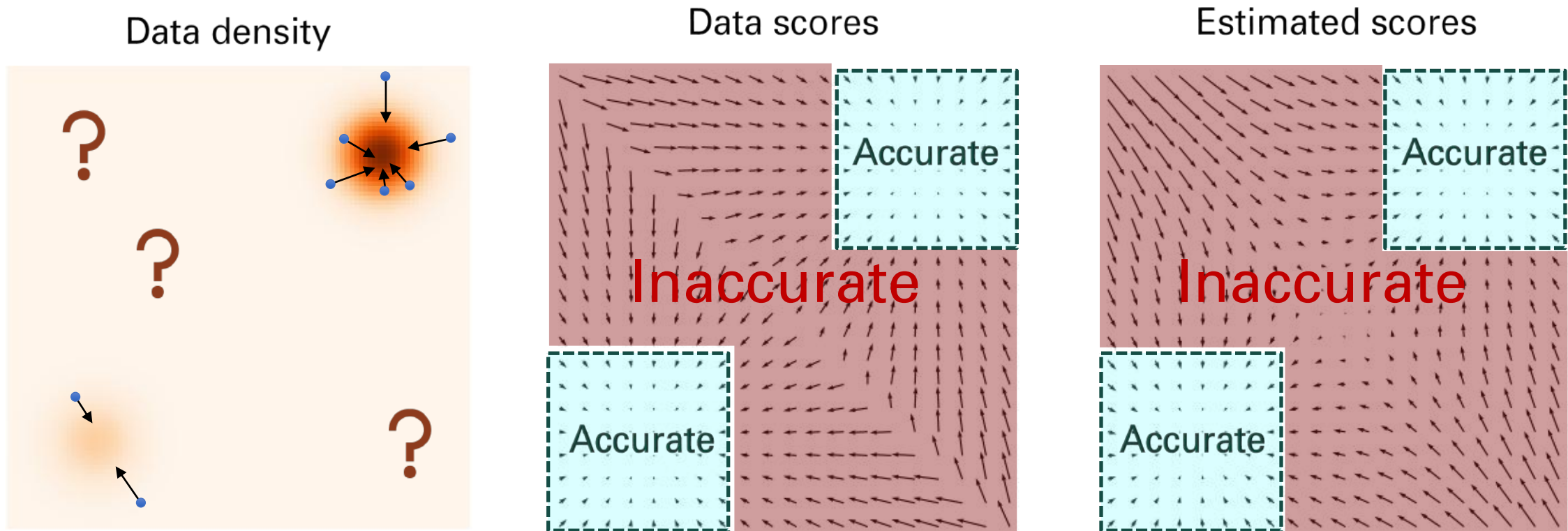


Data unperturbed



Data perturbed with $\mathcal{N}(0; 0.0001)$

Challenge in Low Data Density Regions

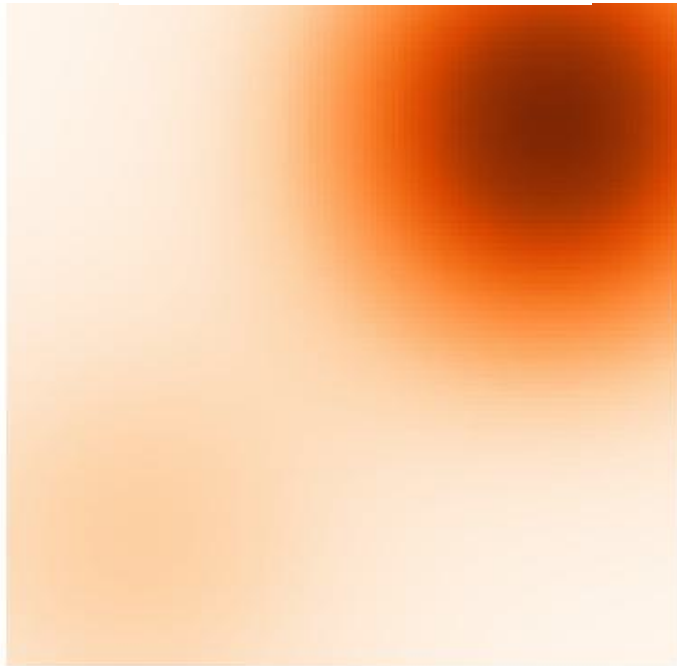


$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}} [\|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - s_{\theta}(\mathbf{x})\|_2^2] \approx \frac{1}{2N} \sum_{i=1}^N \|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}_i) - s_{\theta}(\mathbf{x}_i)\|_2^2$$

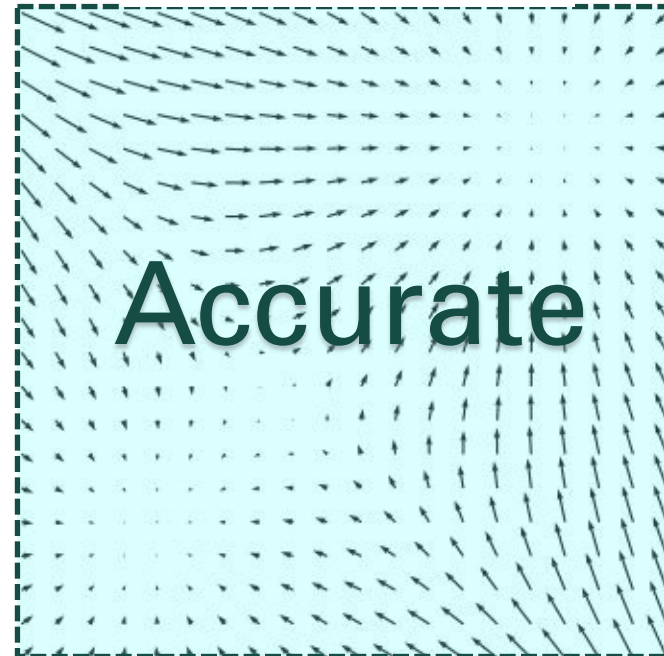
Adding Noise to Data for Better Score Estimation

- Random noise provides samples in low data density regions.

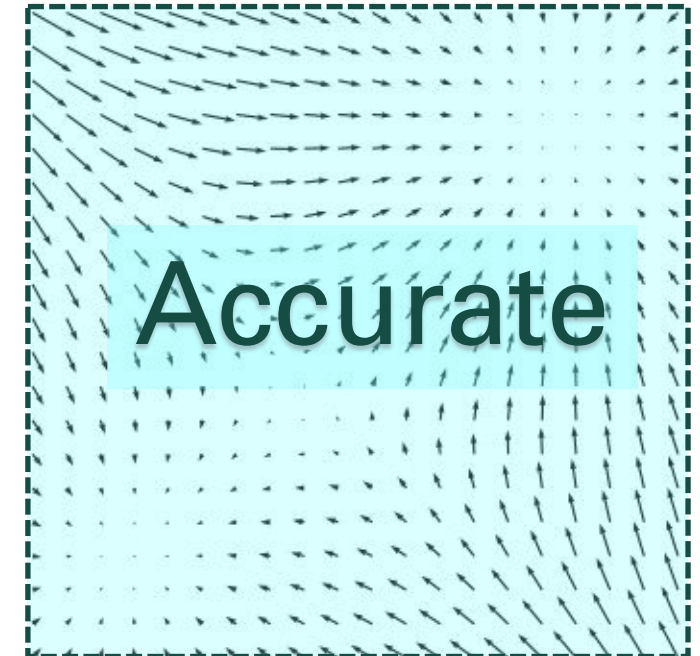
Perturbed density



Perturbed scores

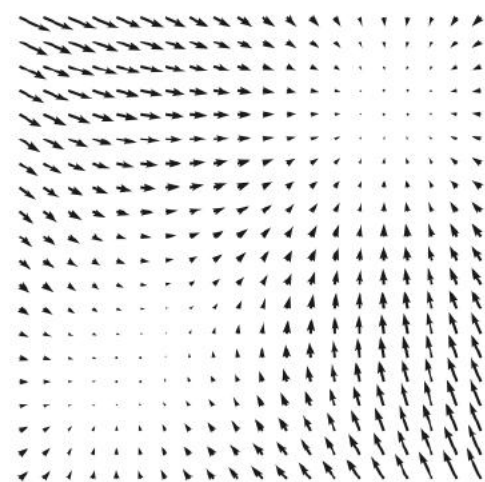
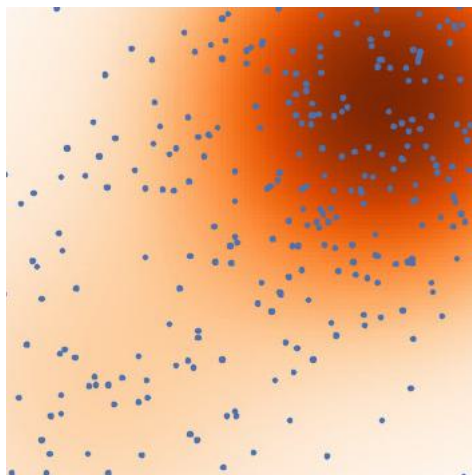


Estimated scores

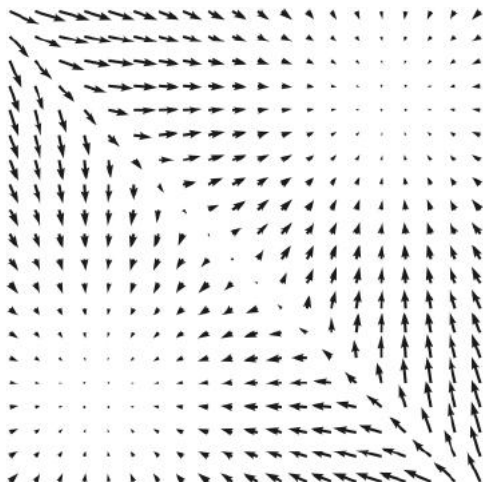
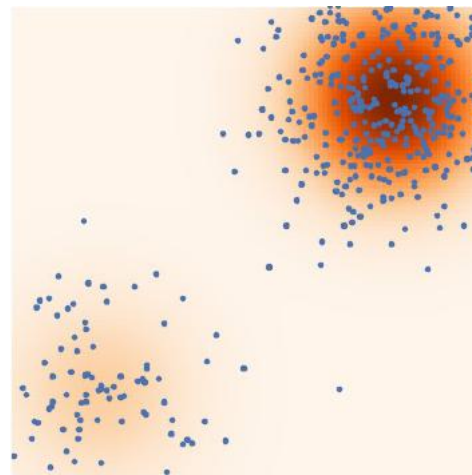


Joint Score Estimation via Noise Conditional Score Networks

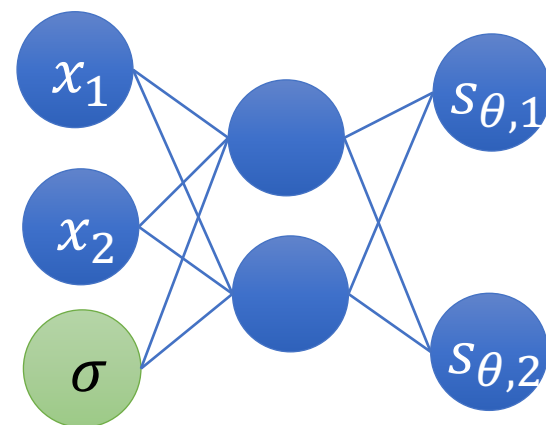
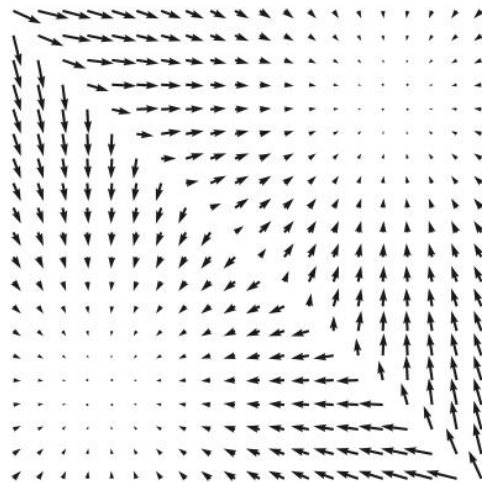
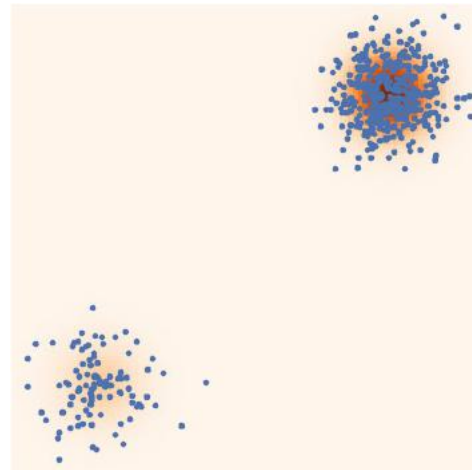
σ_1



σ_2

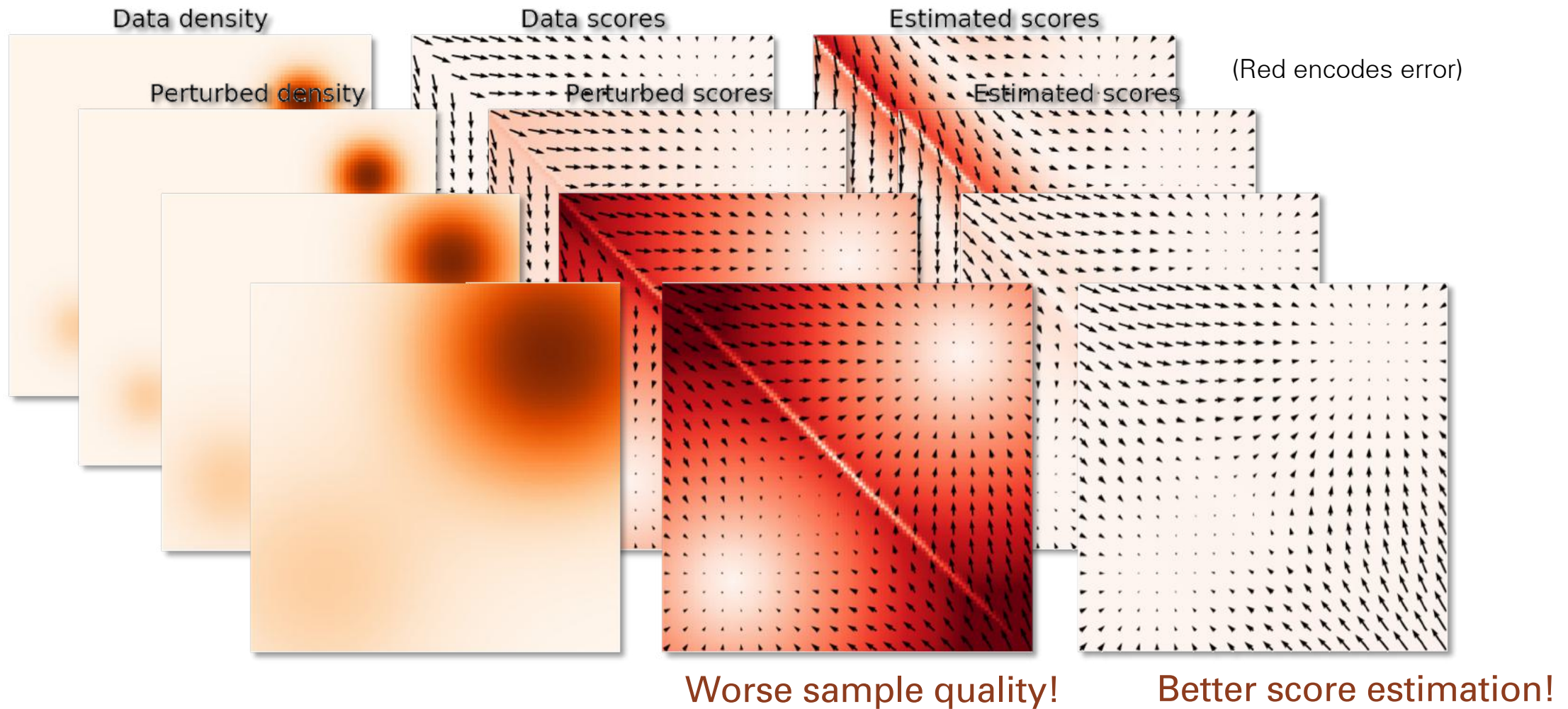


σ_3



Noise Conditional
Score Network
(NCSN)

Sample Quality vs. Estimation Accuracy



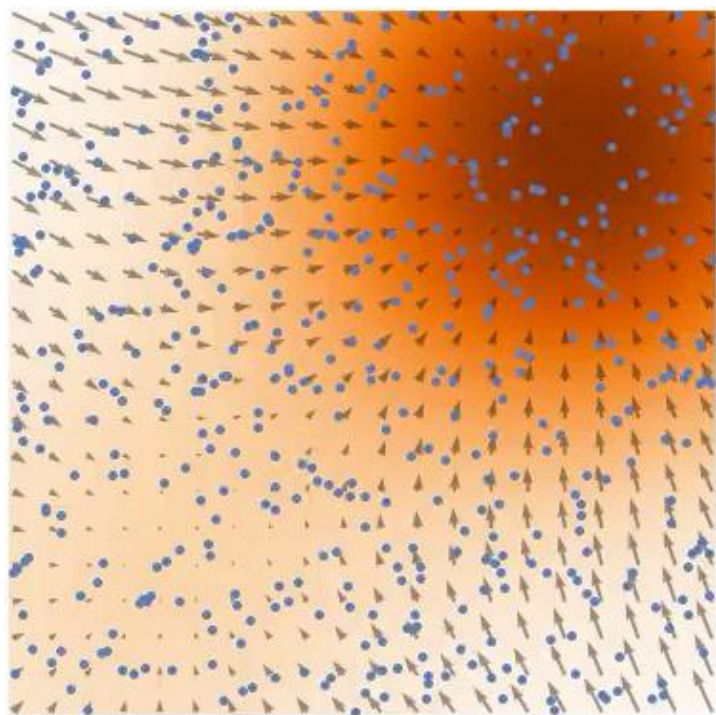
Lecture overview

- Energy-based models
- **Score-based Models**
 - Probability Distribution and Score functions
 - Score-based Generative Models
 - **Sampling from a Score-based Model**
 - Latent Score-based Generative Models

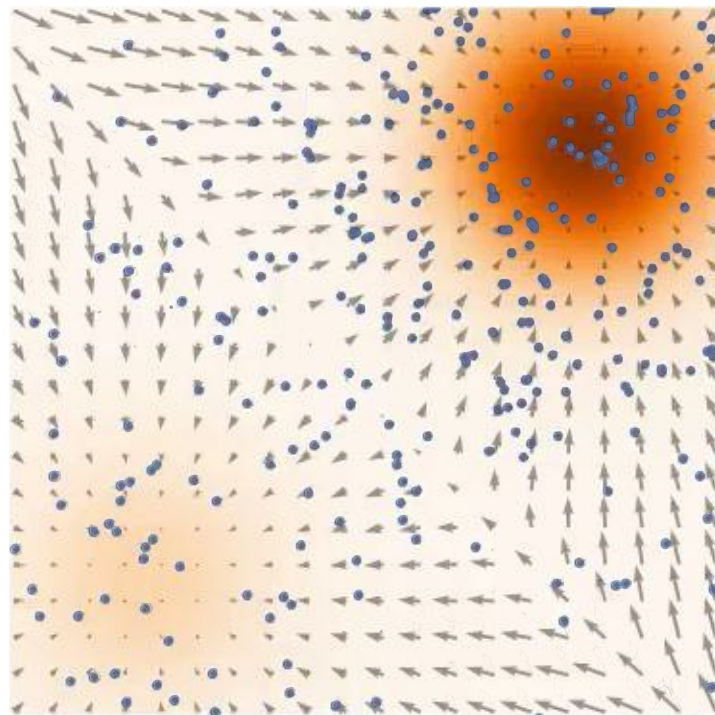
Annealed Langevin Dynamics

Joint Scores to Samples

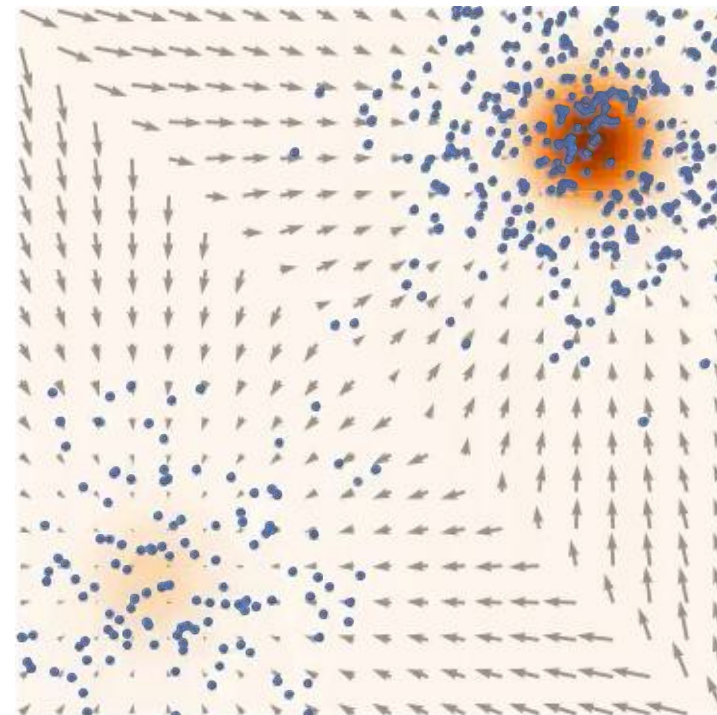
- Sample using $\sigma_1, \sigma_2, \dots, \sigma_L$ sequentially with Langevin dynamics.
- Anneal down the noise level.
- Samples used as initialization for the next level.



σ_1

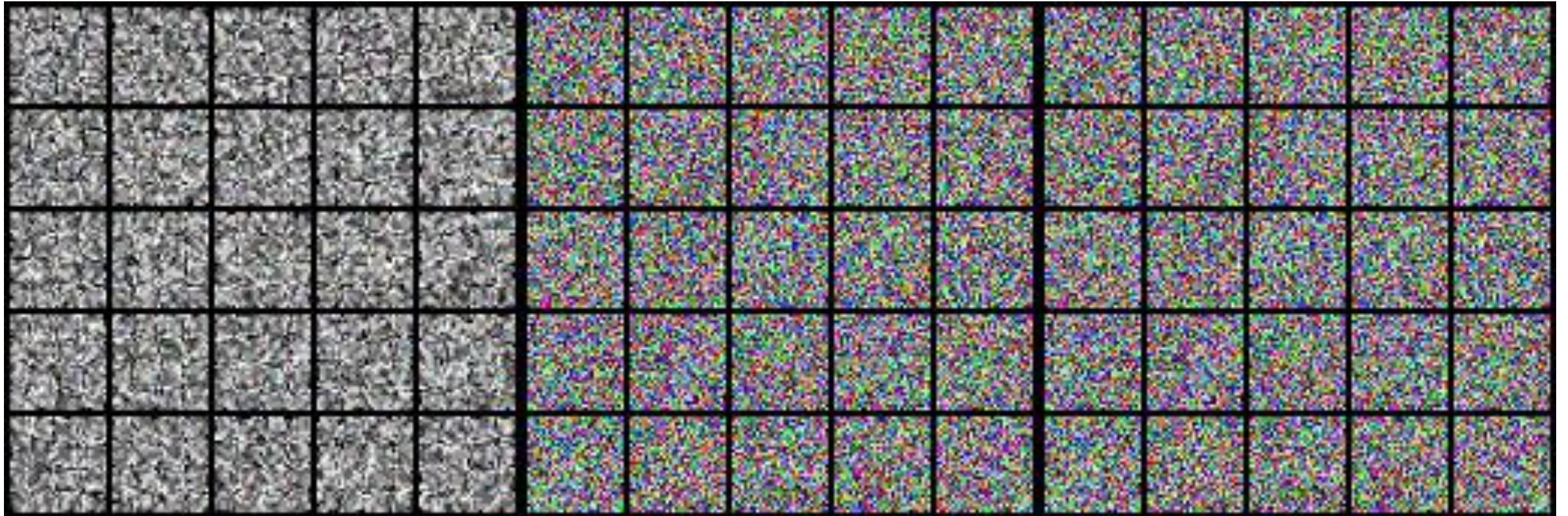


σ_2



σ_3

Experiments: Sampling



Experiments: Sample Quality

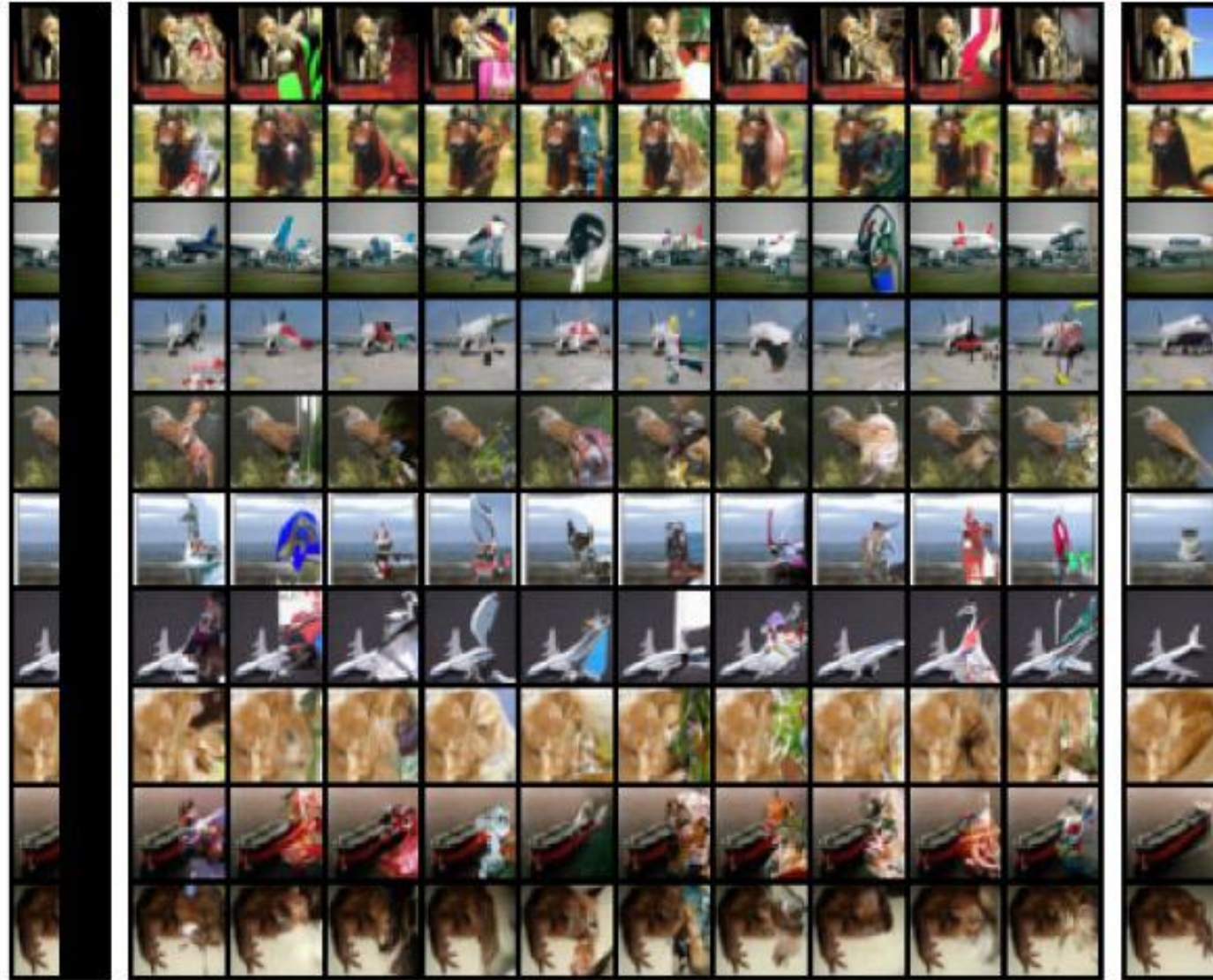
CIFAR-10 Unconditional

Model	Inception Score (higher is better)	FID score (lower is better)
PixelCNN	4.60	65.93
EBM	6.02	40.58
SNGAN	8.22 ± 0.05	21.7
ProgressiveGAN	8.80 ± 0.05	-
NCSN (ours)	8.87 ± 0.12	25.32

Experiments: Inpainting

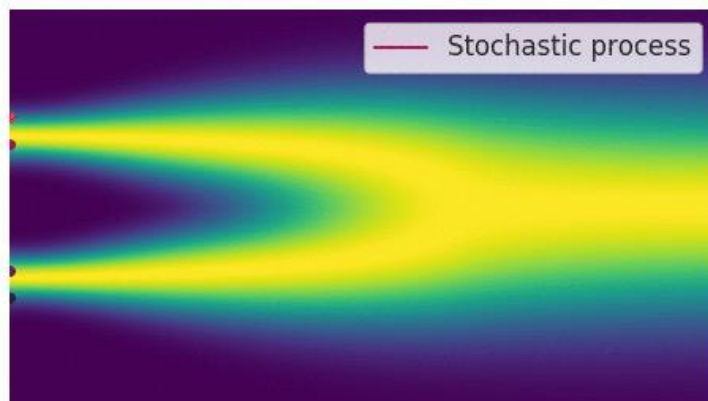
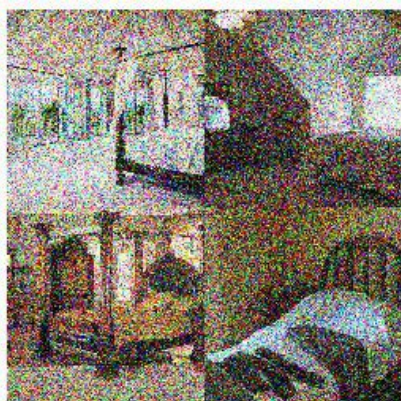
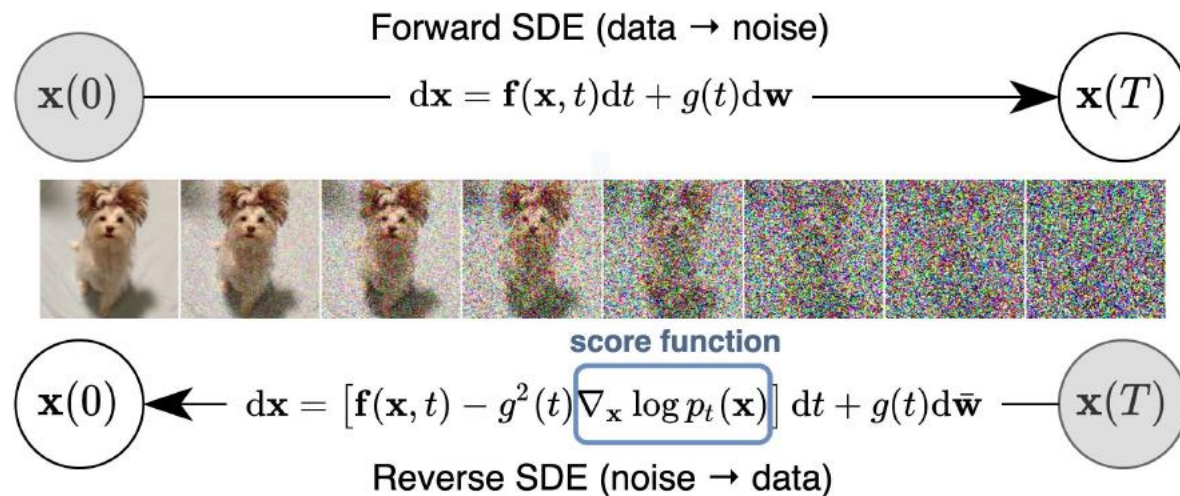


Experiments: Inpainting

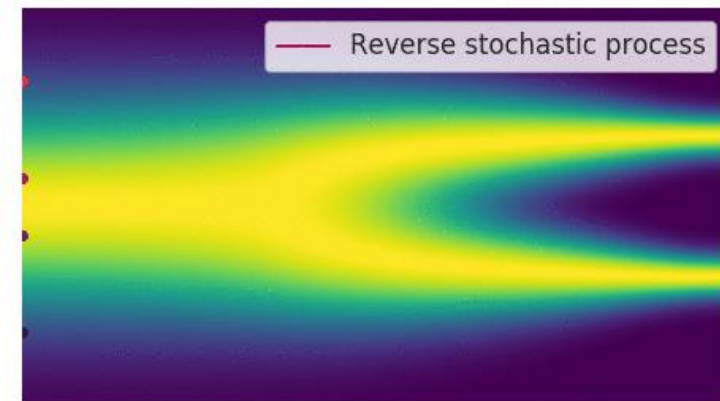
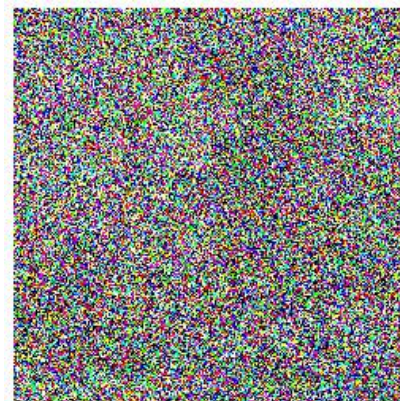


Reversing the SDE for sample generation

- One can generate samples by reversing the perturbation process with annealed Langevin dynamics.
- For infinite noise scales, we can analogously reverse the perturbation process for sample generation by using the reverse SDE.



Perturbing data to noise with a continuous-time stochastic process.



Generating data from noise by reversing the perturbation procedure.

1024 x 1024 samples on FFHQ dataset



256 x 256 samples on LSUN Bedroom



Sample Quality

Table 2: NLLs and FIDs (ODE) on CIFAR-10.

Model	NLL Test ↓	FID ↓
RealNVP (Dinh et al., 2016)	3.49	-
iResNet (Behrmann et al., 2019)	3.45	-
Glow (Kingma & Dhariwal, 2018)	3.35	-
MintNet (Song et al., 2019b)	3.32	-
Residual Flow (Chen et al., 2019)	3.28	46.37
FFJORD (Grathwohl et al., 2018)	3.40	-
Flow++ (Ho et al., 2019)	3.29	-
DDPM (L) (Ho et al., 2020)	$\leq 3.70^*$	13.51
DDPM (L_{simple}) (Ho et al., 2020)	$\leq 3.75^*$	3.17
DDPM	3.28	3.37
DDPM cont. (VP)	3.21	3.69
DDPM cont. (sub-VP)	3.05	3.56
DDPM++ cont. (VP)	3.16	3.93
DDPM++ cont. (sub-VP)	3.02	3.16
DDPM++ cont. (deep, VP)	3.13	3.08
DDPM++ cont. (deep, sub-VP)	2.99	2.92

Table 3: CIFAR-10 sample quality.

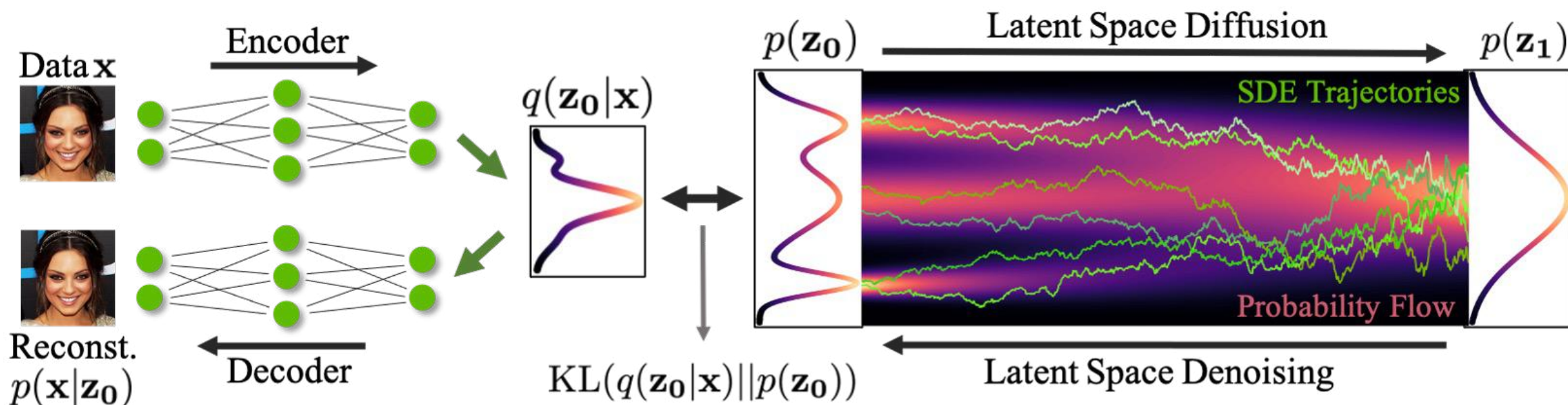
Model	FID ↓	IS ↑
Conditional		
BigGAN (Brock et al., 2018)	14.73	9.22
StyleGAN2-ADA (Karras et al., 2020a)	2.42	10.14
Unconditional		
StyleGAN2-ADA (Karras et al., 2020a)	2.92	9.83
NCSN (Song & Ermon, 2019)	25.32	$8.87 \pm .12$
NCSNv2 (Song & Ermon, 2020)	10.87	$8.40 \pm .07$
DDPM (Ho et al., 2020)	3.17	$9.46 \pm .11$
DDPM++	2.78	9.64
DDPM++ cont. (VP)	2.55	9.58
DDPM++ cont. (sub-VP)	2.61	9.56
DDPM++ cont. (deep, VP)	2.41	9.68
DDPM++ cont. (deep, sub-VP)	2.41	9.57
NCSN++	2.45	9.73
NCSN++ cont. (VE)	2.38	9.83
NCSN++ cont. (deep, VE)	2.20	9.89

Lecture overview

- Energy-based models
- **Score-based Models**
 - Probability Distribution and Score functions
 - Score-based Generative Models
 - Sampling from a Score-based Model
 - **Latent Score-based Generative Models**

Latent Score-based Generative Models

- Score-based generative models (SGMs) are applied directly in data space and often require 1000s of network evaluations for sampling.
- **Idea:** Can we train SGMs in a latent space?



Latent Score-based Generative Models

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \phi, \theta, \psi) &= \mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})}[-\log p_\psi(\mathbf{x} | \mathbf{z}_0)] + \text{KL}(q_\phi(\mathbf{z}_0 | \mathbf{x}) \| p_\theta(\mathbf{z}_0)) \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})}[-\log p_\psi(\mathbf{x}|\mathbf{z}_0)]}_{\text{reconstruction term}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})}[\log q_\phi(\mathbf{z}_0|\mathbf{x})]}_{\text{negative encoder entropy}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})}[-\log p_\theta(\mathbf{z}_0)]}_{\text{cross entropy}}\end{aligned}$$

- A simple expression for the cross-entropy term in the variational loss
- A parameterization of the latent space score function, which mixes a Normal distribution with a learnable SGM.

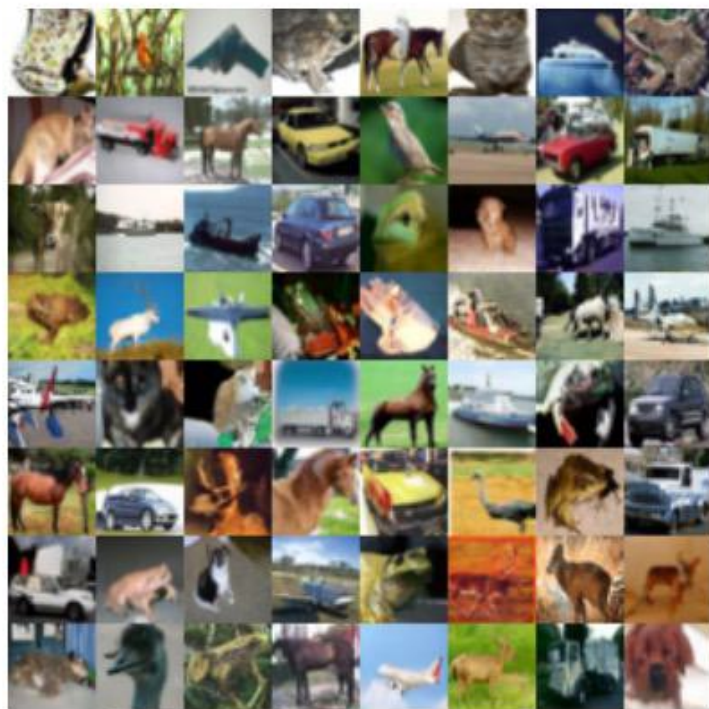
Also employs a SDE-based variance reduction importance sampling schemes to stably train deep LSGMs.

Latent Score-based Generative Models



The evolution of the latent variables under the reverse-time generative process by feeding latent variables from different stages along the process to the decoder to map them back to image space

Latent SGM (LSGM) Samples



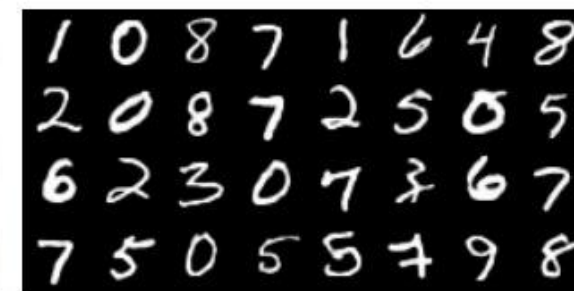
(a) CIFAR-10



(b) CelebA-HQ-256



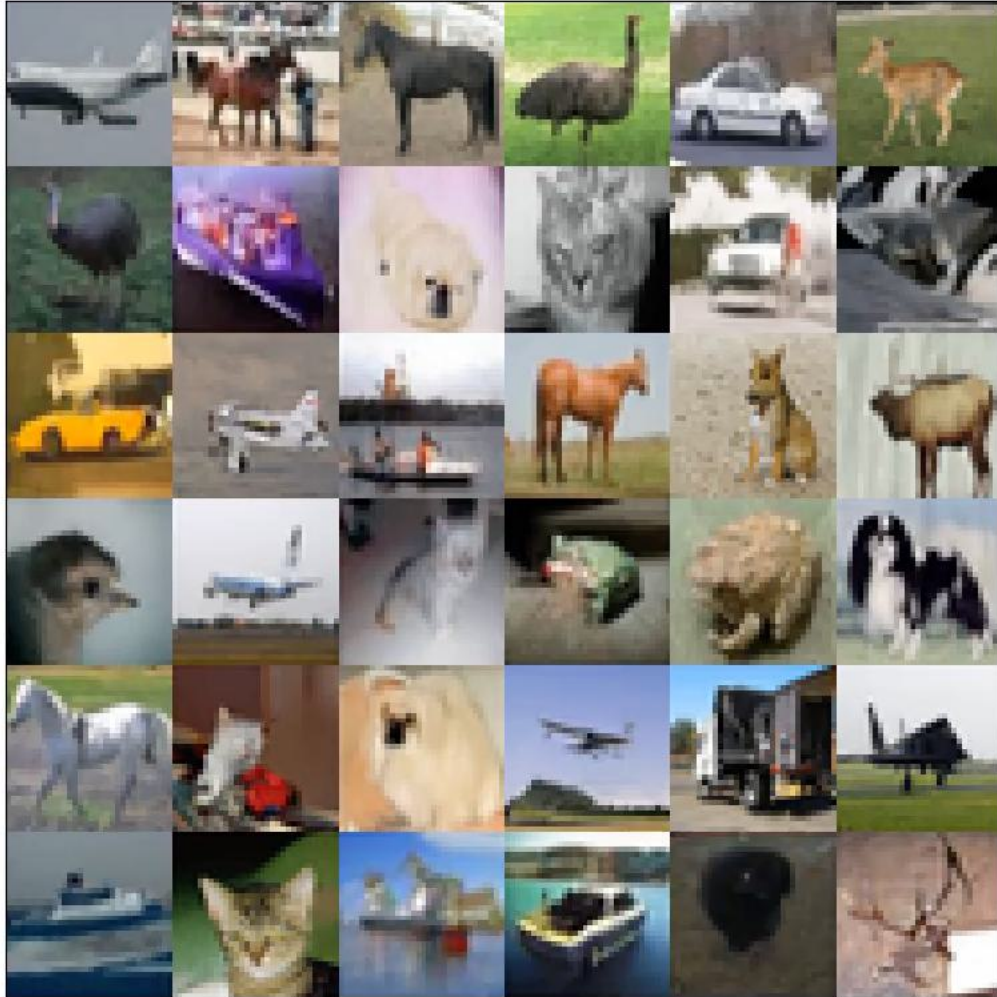
(c) OMNIGLOT



(d) MNIST

Traversing in the latent space of LSGM.

CIFAR-10



CelebA-HQ-256



LSGM Sample Quality

Table 2: Generative performance on CIFAR-10.

	Method	NLL↓	FID↓
Ours	LSGM (FID)	≤3.43	2.10
	LSGM (NLL)	≤ 2.87	6.89
	LSGM (balanced)	≤2.95	2.17
	VAE Backbone	2.96	43.18
VAEs	VDVAE [21]	2.87	-
	NVAE [20]	2.91	23.49
	VAEBM [76]	-	12.19
	NCP-VAE [56]	-	24.08
	BIVA [48]	3.08	-
	DC-VAE [77]	-	17.90
	Score	NCSN [3]	-
Rec. Likelihood [40]		3.18	9.36
DSM-ALS [39]		3.65	-
DDPM [1]		3.75	3.17
Improved DDPM [26]		2.94	11.47
SDE (DDPM++) [2]		2.99	2.92
SDE (NCSN++) [2]		-	2.20
Flows	VFlow [19]	2.98	-
	ANF [18]	3.05	-
Aut. Reg.	DistAug aug [78]	2.53	42.90
	Sp. Transformers [79]	2.80	-
	δ-VAE [80]	2.83	-
	PixelSNAIL [81]	2.85	-
	PixelCNN++ [82]	2.92	-
GANs	AutoGAN [83]	-	12.42
	StyleGAN2-ADA [84]	-	2.92

Table 3: Generative results on CelebA-HQ-256.

	Method	NLL↓	FID↓
Ours	LSGM	≤ 0.70	7.22
	VAE Backbone	0.70	30.87
VAEs	NVAE [20]	0.70	29.76
	VAEBM [76]	-	20.38
	NCP-VAE [56]	-	24.79
	DC-VAE [77]	-	15.80
Score	SDE [2]	-	7.23
Flows	GLOW [85]	1.03	68.93
Aut. Reg.	SPN [86]	0.61	-
GANs	Adv. LAE [87]	-	19.21
	VQ-GAN [64]	-	10.70
	PGGAN [88]	-	8.03

- LSGM obtains the SOTA FID score of 2.10 on CIFAR-10, outperforming previous GANs and SGMs.
- On CelebA-HQ-256, it is on a par with previous SGMs while being 50x to 600x faster in sampling.

Conclusion

- Score-based generative modeling
 - **No need to be normalized / invertible**
 - Flexible architecture choices
 - **No minimax optimization**
 - stable training
 - a natural measurement of training progress / model comparison
- **Adding noise and annealing the noise levels** are critical
- Better or comparable sample quality to GANs.

Related Work

- Generative Stochastic Networks (Bengio et al. (2013), Alain et al. (2016))
 - Sampling starts **close to data points**.
 - Need **MCMC during training** with walkback.
- Nonequilibrium Thermodynamics (Sohl-Dickstein et al. (2015)), Infusion Training (Bordes et al. (2017)), Variational Walkback (Goyal et al. (2017))
 - **Likelihood-based training**.
 - Need **MCMC during training**.

Experiments: Nearest Neighbors



Future Directions

- How to apply score-based generative modeling to discrete data?
- Theoretical guidance on how to choose noise levels?
- Better architecture for higher resolution image generation?
- Improved score estimation?

Next lecture: Diffusion Models