

Deep Generative Models

12. Score-Based Models



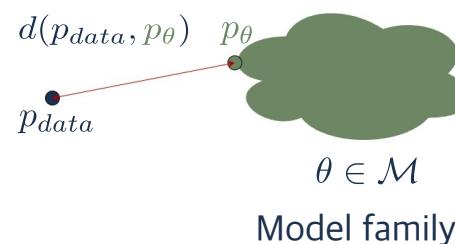
- 국가수리과학연구소 산업수학혁신센터 김민중

Summary

- **Representation:** how do we model the joint distribution of many random variables?
 - Need compact representation
- **Learning:** what is the right way to compare probability distributions?



$$\mathbf{x}_i \sim p_{data} \\ i = 1, 2, \dots, N$$



- **Inference:** how do we invert the generation process

Representation

- Probability density function (p.d.f.) or probability mass function (p.m.f.)

$$p(x)$$

- Autoregressive Models

$$p_{\theta}(x) = \prod_{i=1}^d p_{\theta}(x_i | x_{<i})$$

- Variational Autoencoders

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p(z) dz$$

Representation

- Probability density function (p.d.f.) or probability mass function (p.m.f.)

$$p(x)$$

- Flow Models

$$p_X(x; \theta) = p_Z(f_\theta^{-1}(x)) \left| \det \left(\frac{\partial f_\theta^{-1}(x)}{\partial x} \right) \right|$$

- Energy-based models

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

Representation

- Probability density function (p.d.f.) or probability mass function (p.m.f.)
 $p(x)$
- Pros
 - Maximum likelihood training
 - Principled model comparison via likelihoods
- Cons
 - Special architectures or surrogate losses to deal with intractable partition functions

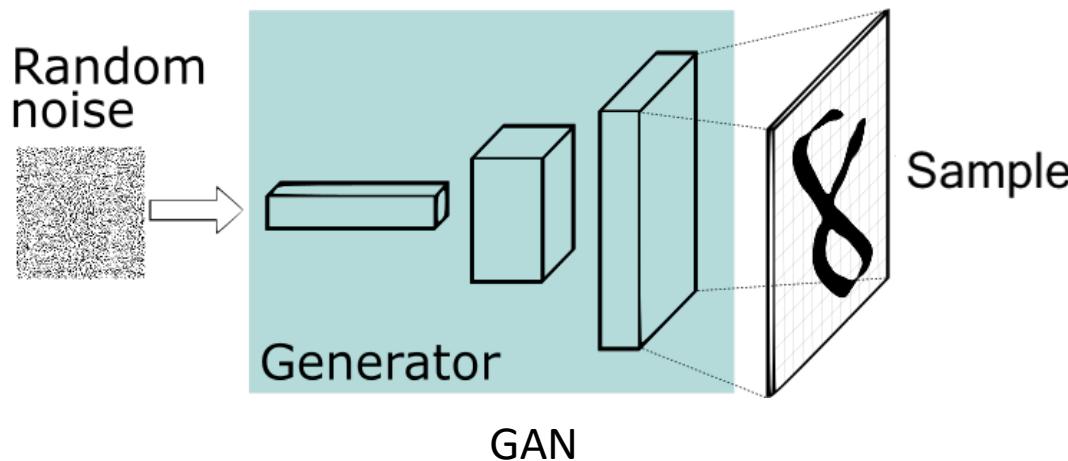
Implicit generative model

- Generative adversarial networks (GANs)

$$z \sim p(z)$$

$$x = G_{\theta}(z)$$

- A two-player minimax game between a generator and a discriminator(two sample test)



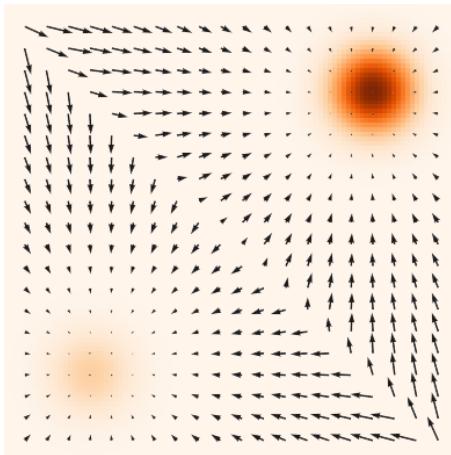
Generative adversarial networks

- Pros
 - Samples typically have better quality
- Cons
 - Require adversarial training
 - Training instability and mode collapse
 - No principled way to compare different models
 - No principled termination criteria for training

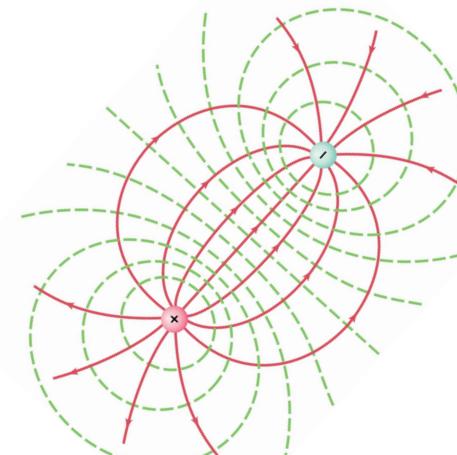
Score functions

- When the pdf is differentiable, we can compute the gradient of a probability density
- Score function

$$\nabla_x \log p(x)$$



(pdf and score)

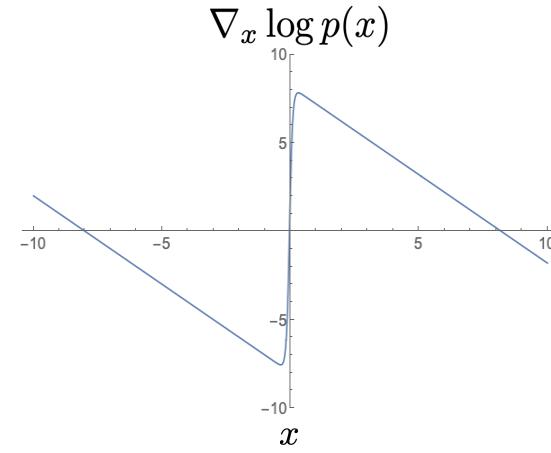
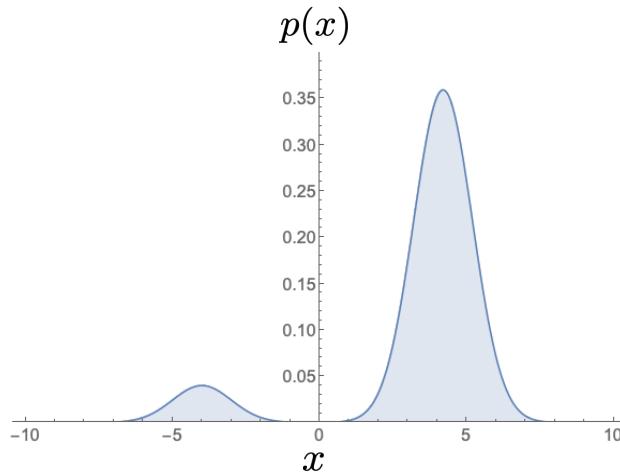


(Electrical potentials and fields)

Score functions

- When the pdf is differentiable, we can compute the gradient of a probability density
- Score function

$$\nabla_x \log p(x)$$



Recap. of Energy-based model

- Energy-based models: $\frac{1}{Z(\theta)} \exp(f_\theta(x))$
 - $Z(\theta)$ is intractable, so no access to likelihood
 - Comparing the probability of two points is easy

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

- Maximum likelihood training:

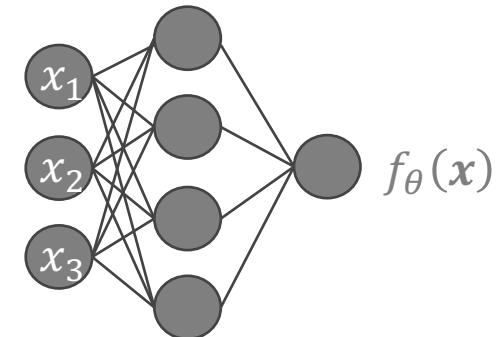
$$\max_{\theta} [f_\theta(x_{train}) - \log Z(\theta)]$$

- Contrastive divergence:

$$\nabla_{\theta} f_{\theta}(x_{train}) - \nabla_{\theta} \log Z(\theta) \approx \nabla_{\theta} f_{\theta}(x_{train}) - \nabla_{\theta} f_{\theta}(x_{sample})$$

- where $x_{sample} \sim p_\theta(x) = \frac{\exp(f_\theta(x))}{Z(\theta)}$

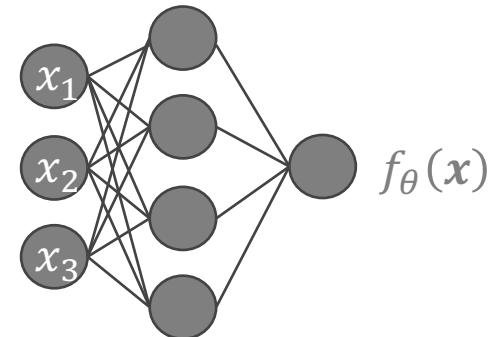
- Requires iterative sampling during training



Recap. of Energy-based model

- Energy-based models

$$p_{\theta}(x) = \frac{1}{Z(\theta)} \exp(f_{\theta}(x))$$
$$f_{\theta}(x) \in \mathbb{R}$$



- Minimizing Fisher divergence

$$\min_{\theta} \frac{1}{2} E_{x \sim p_{data}} [\| \nabla_x \log p_{data}(x) - \nabla_x \log p_{\theta}(x) \|_2^2]$$

- Score matching

$$\frac{1}{2} E_{x \sim p_{data}} [\| \nabla_x \log p_{data}(x) - \nabla_x \log p_{\theta}(x) \|_2^2]$$

$$= E_{x \sim p_{data}} \left[\frac{1}{2} \| \nabla_x \log p_{\theta}(x) \|_2^2 + \text{tr}(\nabla_x^2 \log p_{\theta}(x)) \right] + \text{const.}$$

Score matching for training EBMs

- Score function of EBMs

$$\nabla_x \log p_\theta(x) = \nabla_x f_\theta(x) - \nabla_x \log Z(\theta) = \nabla_x f_\theta(x)$$

- Score matching for EBMs

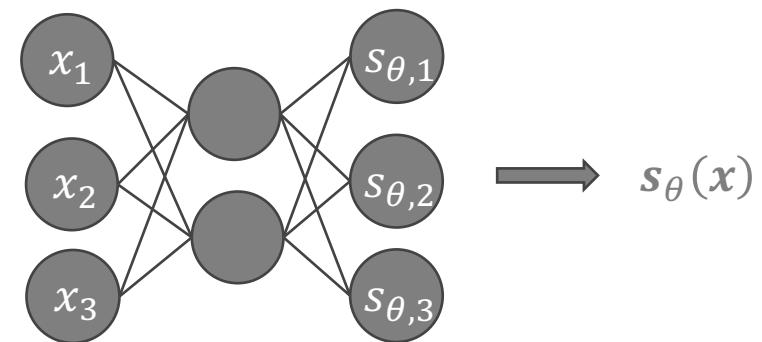
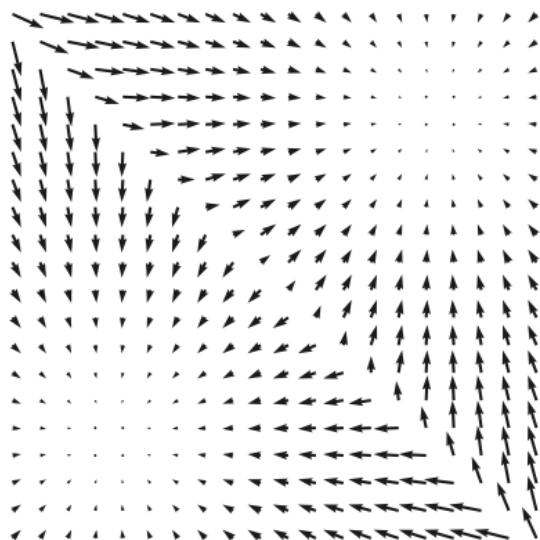
$$\begin{aligned} & E_{x \sim p_{data}} \left[\frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2 + \text{tr}(\nabla_x^2 \log p_\theta(x)) \right] \\ &= E_{x \sim p_{data}} \left[\frac{1}{2} \|\nabla_x f_\theta(x)\|_2^2 + \text{tr}(\nabla_x^2 f_\theta(x)) \right] \end{aligned}$$

Score-based models

- Directly model the vector field of gradients

$$s_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$s_\theta(x) \approx \nabla_x \log p_{data}(x)$$

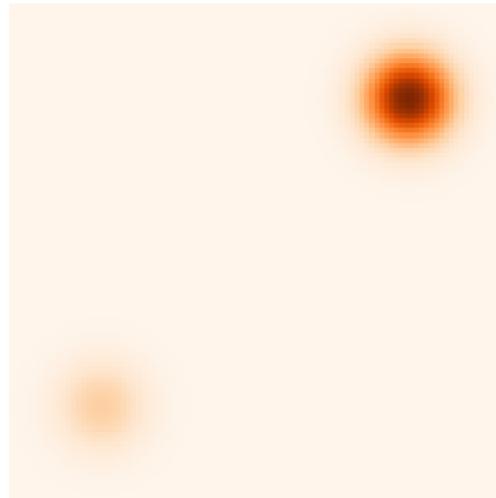


Framework for score estimation

- Score estimation by training score-based models

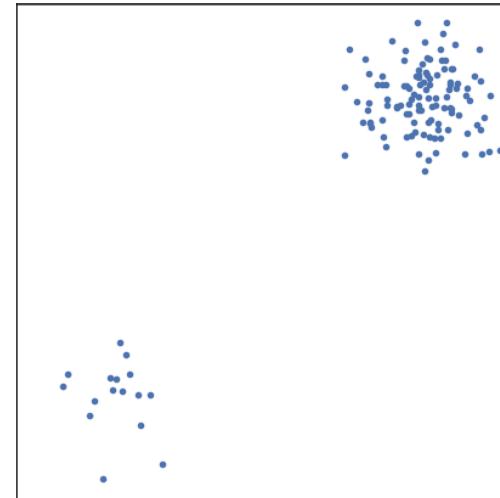
Probability density

$$p_{data}(\mathbf{x})$$



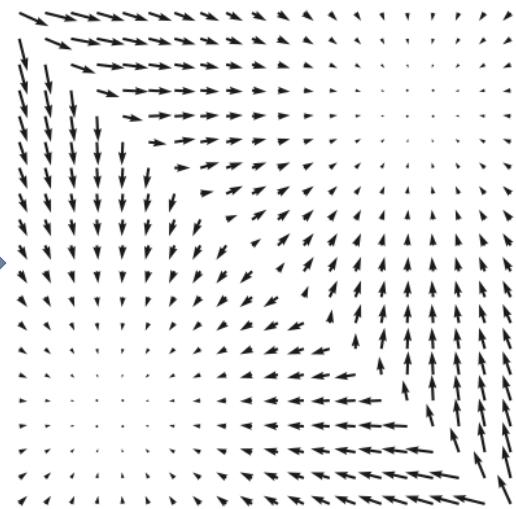
i.i.d. samples

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$$



Score function

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



Score estimation by training score-based models

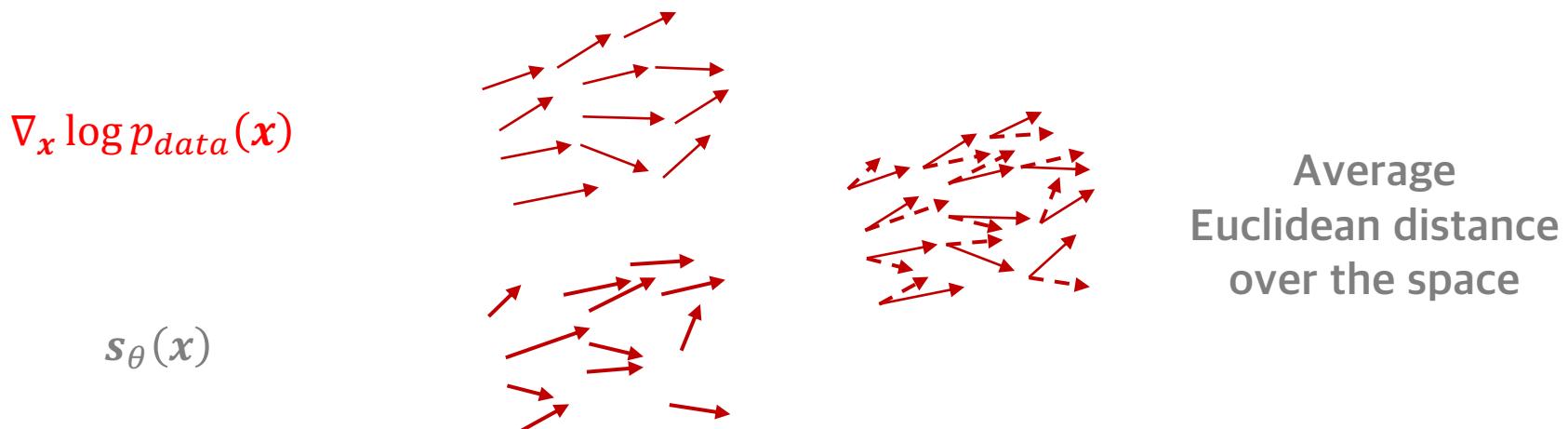
- Given: i.i.d. samples $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\} \sim p_{data}(\mathbf{x})$
- Task: Estimating the score $\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$
- Score model: A learnable vector valued function

$$s_{\theta}(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}^d$$

- Goal

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

- How to compare two vector fields of scores?



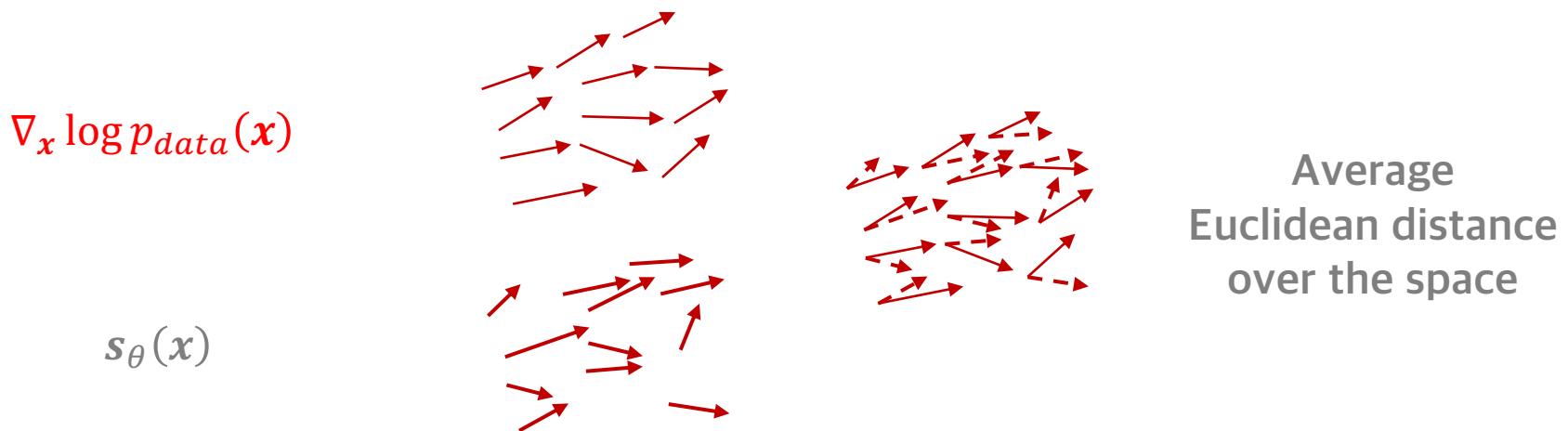
Score estimation by training score-based models

- Objective: Average Euclidean distance over the whole space

$$\frac{1}{2} E_{x \sim p_{data}} [\|\nabla_x \log p_{data}(x) - s_\theta(x)\|_2^2]$$

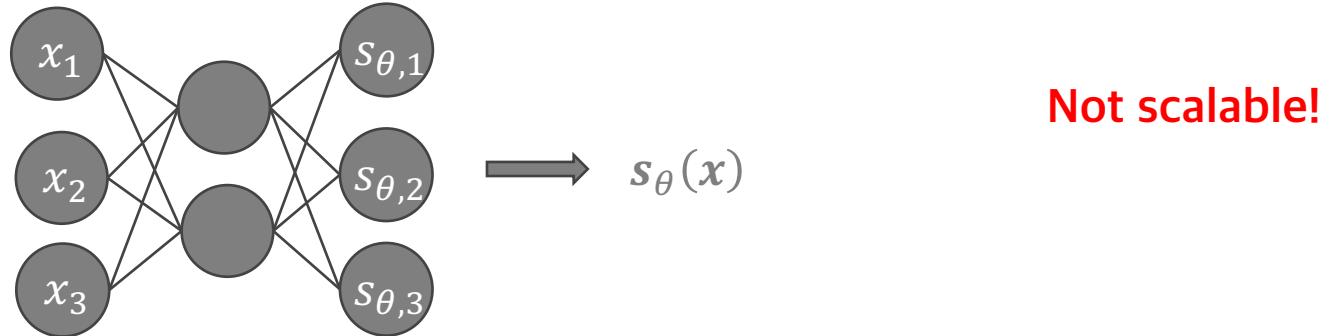
- Score matching

$$E_{x \sim p_{data}} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right]$$



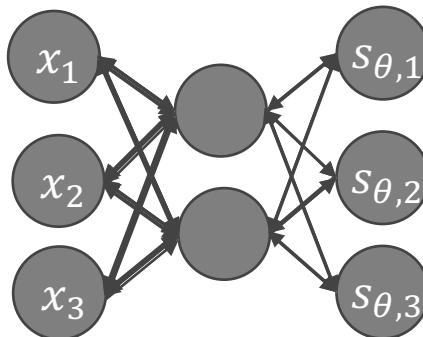
Score matching is not scalable

- Deep neural networks as more expressive score models



- Need to compute $\|s_\theta(x)\|_2^2$ and $\text{tr}(\nabla_x s_\theta(x))$

$$\begin{aligned}\frac{\partial s_{\theta,1}(x)}{\partial x_1} \\ \frac{\partial s_{\theta,2}(x)}{\partial x_2} \\ \frac{\partial s_{\theta,3}(x)}{\partial x_3}\end{aligned}$$



$$\begin{aligned}s_{\theta,1}(x) \\ s_{\theta,2}(x) \\ s_{\theta,3}(x)\end{aligned}$$

$$\nabla_x s_\theta(x) = \begin{pmatrix} \frac{\partial s_{\theta,1}(x)}{\partial x_1} & \frac{\partial s_{\theta,1}(x)}{\partial x_2} & \frac{\partial s_{\theta,1}(x)}{\partial x_3} \\ \frac{\partial s_{\theta,2}(x)}{\partial x_1} & \frac{\partial s_{\theta,2}(x)}{\partial x_2} & \frac{\partial s_{\theta,2}(x)}{\partial x_3} \\ \frac{\partial s_{\theta,3}(x)}{\partial x_1} & \frac{\partial s_{\theta,3}(x)}{\partial x_2} & \frac{\partial s_{\theta,3}(x)}{\partial x_3} \end{pmatrix}$$

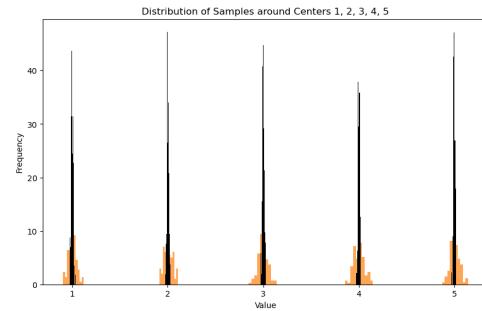
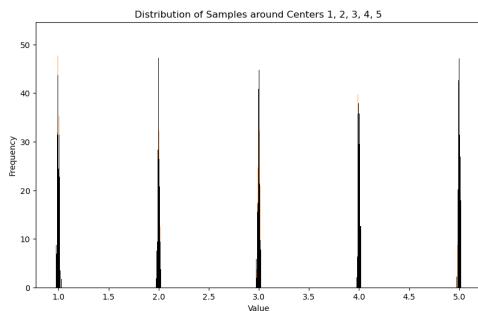
Scalable approaches to score estimation

- Denoising score matching (Vincent, 2011)
- Sliced score matching (Song & Garg, 2019)
- Need to eliminate the trace of Jacobian term in the training objective

Denoising score matching (Vincent, 2011)

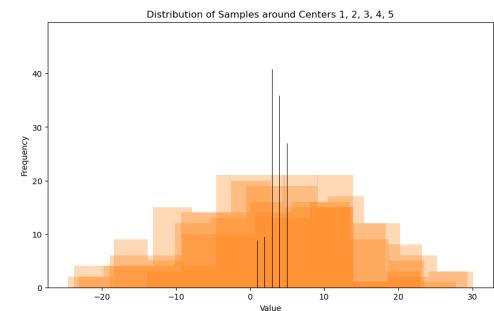
- Consider the perturbed distribution

$$q_\sigma(\tilde{x}|x) := N(\tilde{x}|x, \sigma^2 I), \quad q_\sigma(\tilde{x}) = \int p_{data}(x)q_\sigma(\tilde{x}|x)dx$$



$p_{data}(x)$

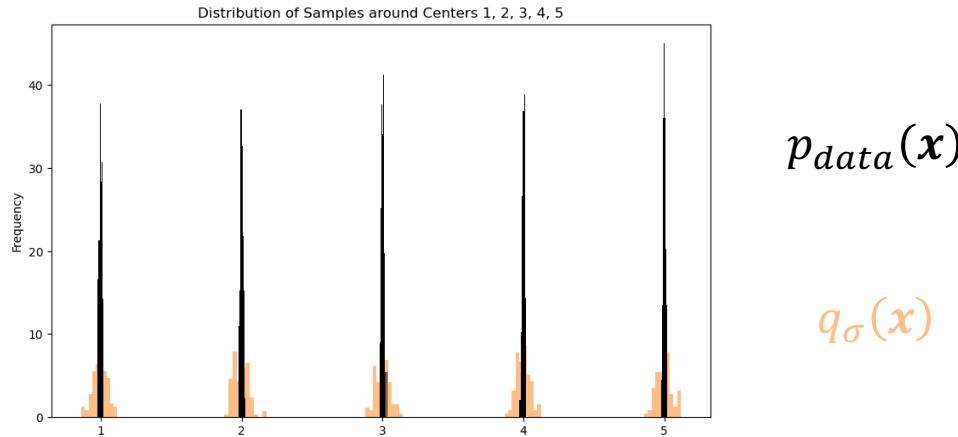
$q_\sigma(x)$



Denoising score matching (Vincent, 2011)

- Consider the perturbed distribution

$$q_\sigma(\tilde{x}|x) := N(\tilde{x}|x, \sigma^2 I), \quad q_\sigma(\tilde{x}) = \int p_{data}(x) q_\sigma(\tilde{x}|x) dx$$



- Score estimation for $\nabla_{\tilde{x}} \log q_\sigma(\tilde{x})$ is easier
- If the noise level is small, this is a good approximation $q_\sigma \approx p_{data}$

Denoising score matching

- Matching the score of a noise-perturbed distribution
- Denoising score matching: $s_\theta(\tilde{x}) \approx \nabla_{\tilde{x}} \log q_\sigma(\tilde{x})$

$$\frac{1}{2} E_{\tilde{x} \sim q_\sigma} [\| \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x}) \|_2^2]$$

- However, $q_\sigma(\tilde{x}) = \int p_{data}(x) q_\sigma(\tilde{x}|x) dx$ is unknown

Denoising score matching

- Fisher divergence

$$\frac{1}{2} E_{\tilde{x} \sim q_\sigma} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x})\|_2^2]$$

$$= \frac{1}{2} \int q_\sigma(\tilde{x}) \|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x})\|_2^2 d\tilde{x}$$

$$= \frac{1}{2} \int q_\sigma(\tilde{x}) \|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x})\|_2^2 d\tilde{x} + \frac{1}{2} \int q_\sigma(\tilde{x}) \|s_\theta(\tilde{x})\|_2^2 d\tilde{x}$$

$$- \int q_\sigma(\tilde{x}) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x})^T s_\theta(\tilde{x}) d\tilde{x}$$

$$= \text{const.} + \frac{1}{2} \int q_\sigma(\tilde{x}) \|s_\theta(\tilde{x})\|_2^2 d\tilde{x} - \int q_\sigma(\tilde{x}) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x})^T s_\theta(\tilde{x}) d\tilde{x}$$

Denoising score matching

$$\begin{aligned} \frac{1}{2} \int q_\sigma(\tilde{x}) \|s_\theta(\tilde{x})\|_2^2 d\tilde{x} &= \frac{1}{2} \int \int p_{data}(x) q_\sigma(\tilde{x}|x) \|s_\theta(\tilde{x})\|_2^2 d\tilde{x} dx \\ &= \frac{1}{2} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\|s_\theta(\tilde{x})\|_2^2] \end{aligned}$$

Denoising score matching

$$\begin{aligned} & - \int q_\sigma(\tilde{x}) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x})^T s_\theta(\tilde{x}) d\tilde{x} = - \int \nabla_{\tilde{x}} q_\sigma(\tilde{x})^T s_\theta(\tilde{x}) d\tilde{x} \\ &= - \int \nabla_{\tilde{x}} \left(\int p_{data}(x) q_\sigma(\tilde{x}|x) dx \right)^T s_\theta(\tilde{x}) d\tilde{x} \\ &= - \int \left(\int p_{data}(x) \nabla_{\tilde{x}} q_\sigma(\tilde{x}|x) dx \right)^T s_\theta(\tilde{x}) d\tilde{x} \\ &= - \int \left(\int p_{data}(x) q_\sigma(\tilde{x}|x) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) dx \right)^T s_\theta(\tilde{x}) d\tilde{x} \\ &= - \int \int p_{data}(x) q_\sigma(\tilde{x}|x) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)^T s_\theta(\tilde{x}) d\tilde{x} dx \end{aligned}$$

Denoising score matching

$$\begin{aligned} & - \int q_\sigma(\tilde{x}) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x})^T s_\theta(\tilde{x}) d\tilde{x} = \\ &= - \int \int p_{data}(x) q_\sigma(\tilde{x}|x) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)^T s_\theta(\tilde{x}) d\tilde{x} dx \\ &= -E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)^T s_\theta(\tilde{x})] \end{aligned}$$

Denoising score matching

$$\begin{aligned} & \frac{1}{2} E_{\tilde{x} \sim q_\sigma} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x})\|_2^2] \\ &= \text{const.} + \frac{1}{2} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\|s_\theta(\tilde{x})\|_2^2] \\ &\quad - E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)^T s_\theta(\tilde{x})] \\ &= \text{const.} + \frac{1}{2} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) - s_\theta(\tilde{x})\|_2^2] \\ &\quad - \frac{1}{2} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)\|_2^2] \\ &= \frac{1}{2} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) - s_\theta(\tilde{x})\|_2^2] + \text{const.} \end{aligned}$$

Denoising score matching

- Estimate the score of a noise-perturbed distribution

$$\min_{\theta} E_{\tilde{x} \sim q_{\sigma}} [\|\nabla_{\tilde{x}} \log q_{\sigma}(\tilde{x}) - s_{\theta}(\tilde{x})\|_2^2]$$

$$\Leftrightarrow \min_{\theta} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_{\sigma}(\tilde{x}|x)} [\|\nabla_{\tilde{x}} \log q_{\sigma}(\tilde{x}|x) - s_{\theta}(\tilde{x})\|_2^2]$$

- $\nabla_{\tilde{x}} \log q_{\sigma}(\tilde{x}|x)$ is easy to compute

$$q_{\sigma}(\tilde{x}|x) = N(\tilde{x}|x, \sigma^2 I), \quad \nabla_{\tilde{x}} \log q_{\sigma}(\tilde{x}|x) = -\frac{1}{\sigma^2}(\tilde{x} - x)$$

- **Pros:** efficient to optimize even for very high dimensional data, and useful for optimal denoising.
- **Con:** cannot estimate the score of clean data (noise-free)

Training for denoising score matching

- Sample a minibatch of datapoints $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ from p_{data}
- Sample a minibatch of perturbed datapoints $\tilde{\mathbf{x}}^{(1)}, \tilde{\mathbf{x}}^{(2)}, \dots, \tilde{\mathbf{x}}^{(n)}$ from $q_\sigma(\tilde{\mathbf{x}})$
$$\tilde{\mathbf{x}}^{(i)} \sim q_\sigma(\tilde{\mathbf{x}}^{(i)} | \mathbf{x}^{(i)}) = N(\tilde{\mathbf{x}}^{(i)} | \mathbf{x}^{(i)}, \sigma^2 \mathbf{I})$$
- Estimate the denoising score matching loss with empirical means

$$\frac{1}{2n} \sum_{i=1}^n \|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}^{(i)} | \mathbf{x}^{(i)}) - s_\theta(\tilde{\mathbf{x}}^{(i)})\|$$

- I.e.,

$$\frac{1}{2n} \sum_{i=1}^n \left\| \frac{1}{\sigma^2} (\tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)}) + s_\theta(\tilde{\mathbf{x}}^{(i)}) \right\|$$

- Stochastic gradient descent
- Need to choose a small σ

Pitfall of denoising score matching

- The loss variance will increase drastically as $\sigma \rightarrow 0$
- Denoising score matching loss for Gaussian perturbations

$$\begin{aligned} & E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} \left[\left\| \frac{1}{\sigma^2} (\tilde{x} - x) + s_\theta(\tilde{x}) \right\|_2^2 \right] \\ &= E_{x \sim p_{data}(x)} E_{z \sim N(0,I)} \left[\left\| \frac{1}{\sigma} z + s_\theta(x + \sigma z) \right\|_2^2 \right] \\ &= E_{x \sim p_{data}(x)} E_{z \sim N(0,I)} \left[\|s_\theta(x + \sigma z)\|_2^2 + 2s_\theta(x + \sigma z)^T \frac{z}{\sigma} + \frac{\|z\|_2^2}{\sigma^2} \right] \end{aligned}$$

- If we choose very small $\sigma \rightarrow 0$

$$\text{Var}\left(\frac{z}{\sigma}\right) \rightarrow \infty,$$

$$E_{z \sim N(0,I)} \left[\frac{\|z\|_2^2}{\sigma^2} \right] \rightarrow \infty$$

Denoising score matching

- Consider the perturbed distribution

$$q_\sigma(\tilde{x}|x) := N(\tilde{x}|x, \sigma^2 I), \quad q_\sigma(\tilde{x}) = \int p_{data}(x)q_\sigma(\tilde{x}|x)dx$$

- Score estimation for $\nabla_{\tilde{x}} \log q_\sigma(\tilde{x})$
- Score matching

$$\min_{\theta} E_{\tilde{x} \sim q_\sigma} [\|\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}) - s_\theta(\tilde{x})\|_2^2]$$

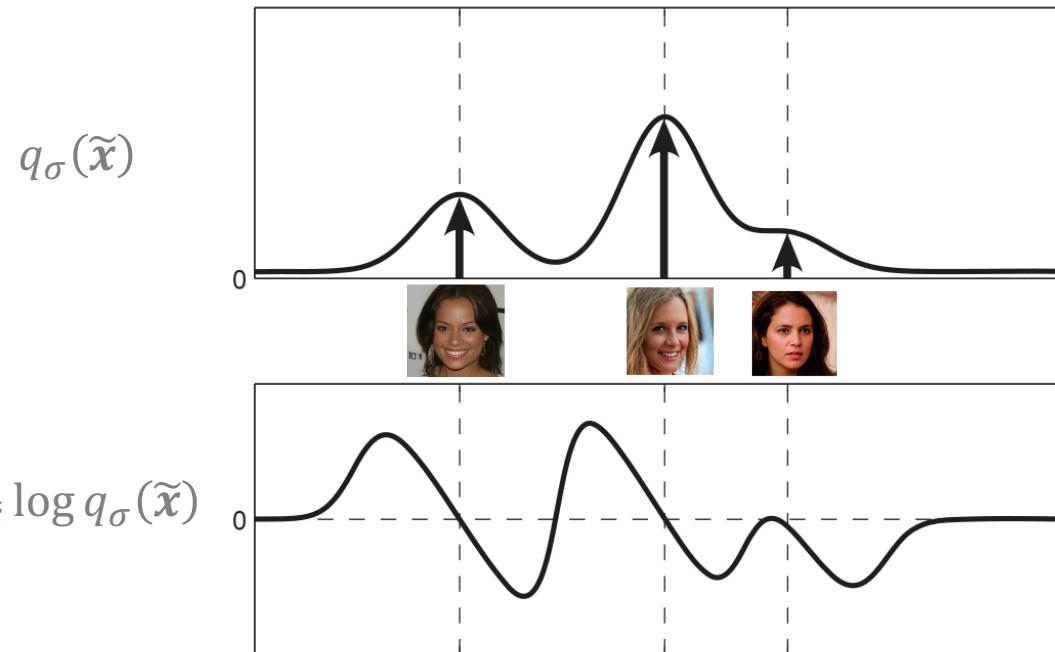
$$\Leftrightarrow \min_{\theta} E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} \left[\left\| \frac{1}{\sigma^2} (\tilde{x} - x) + s_\theta(\tilde{x}) \right\|_2^2 \right]$$



- $s_\theta(\tilde{x})$ tries to estimate the noise that was added to produce \tilde{x}

Tweedie's formula

- Score estimation for $\nabla_{\tilde{x}} \log q_\sigma(\tilde{x})$ is equivalent to denoising
$$E_{x \sim p_{data}(x)} E_{\tilde{x} \sim q_\sigma(\tilde{x}|x)} \left[\left\| \frac{1}{\sigma^2} (\tilde{x} - x) + s_\theta(\tilde{x}) \right\|_2^2 \right]$$
- Tweedie's formula:** Optimal denoising strategy is to follow the gradient (score):



Tweedie's formula and denoising score matching

- Denoising score matching is suitable for optimal denoising
- Given $p_{data}(\mathbf{x})$ and $q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})$, we can define the posterior $p(\mathbf{x}|\tilde{\mathbf{x}})$ with Bayes' rule

$$p(\mathbf{x}|\tilde{\mathbf{x}}) \propto p_{data}(\mathbf{x})q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})$$

- Recall that

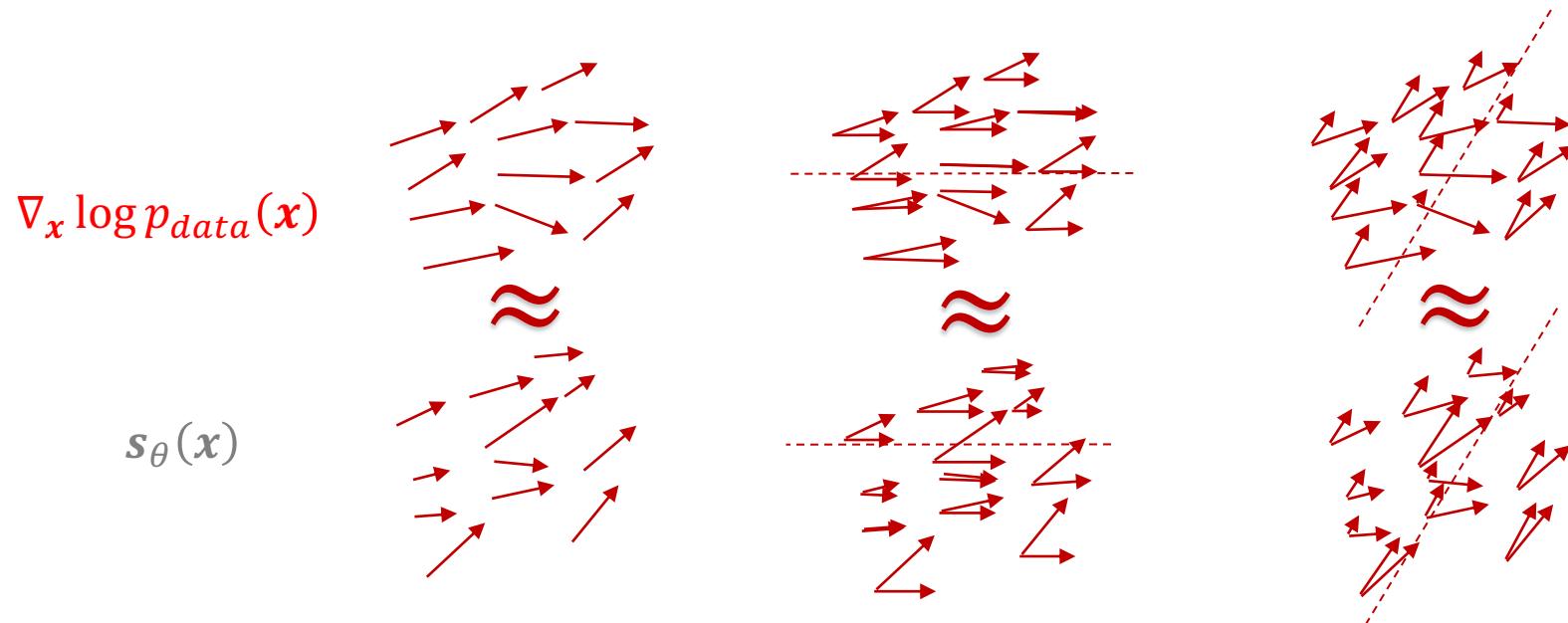
$$q_\sigma(\tilde{\mathbf{x}}) = \int p_{data}(\mathbf{x})q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})d\mathbf{x}$$

- Tweedie's formula:

$$E_{\mathbf{x} \sim p(\mathbf{x}|\tilde{\mathbf{x}})}[\mathbf{x}] = \tilde{\mathbf{x}} + \sigma^2 \nabla_{\mathbf{x}} \log q_\sigma(\tilde{\mathbf{x}}) \approx \tilde{\mathbf{x}} + \sigma^2 s_\theta(\tilde{\mathbf{x}})$$

Sliced score matching

- One dimensional problems is easier
- Consider projections onto random directions



Song*, Garg*, Shi, Ermon. "Sliced Score Matching: A Scalable Approach to Density and Score Estimation." UAI 2019.

Sliced score matching

- **Objective:** Sliced Fisher Divergence

$$\frac{1}{2} E_{\mathbf{v} \sim p_{\mathbf{v}}} E_{\mathbf{x} \sim p_{\text{data}}} \left[(\mathbf{v}^T \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x}))^2 \right]$$

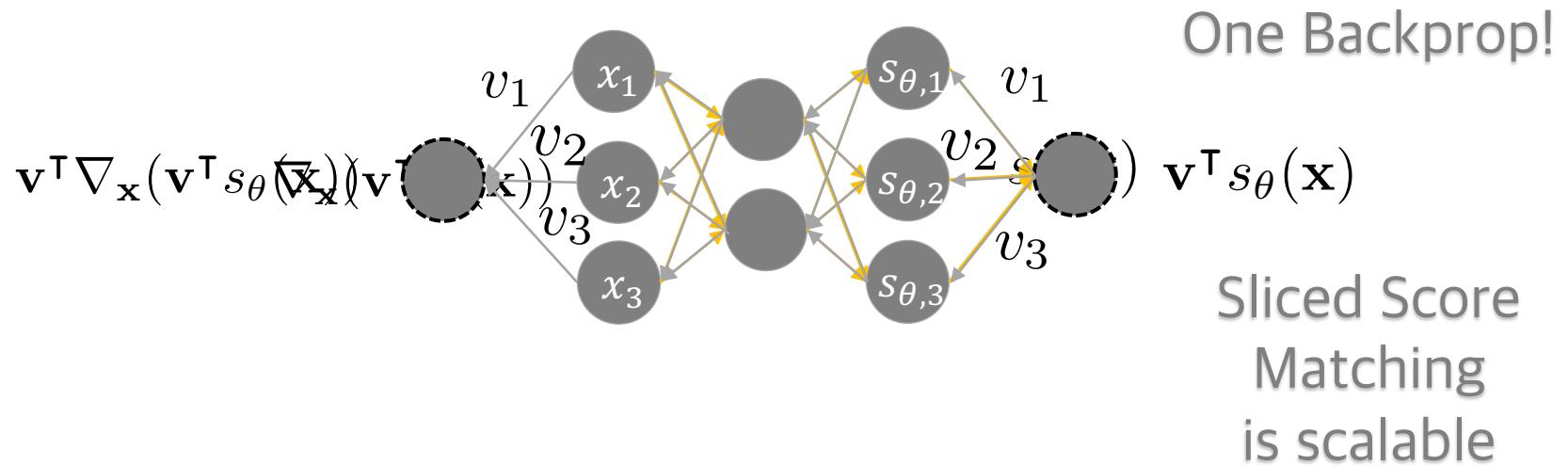
- Integration by parts

$$E_{\mathbf{v} \sim p_{\mathbf{v}}} E_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{2} (\mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x}))^2 + \mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v} \right]$$

$$(\nu_1 \quad \nu_2 \quad \nu_3) \begin{pmatrix} \frac{\partial s_{\theta,1}(\mathbf{x})}{\partial x_1} & \frac{\partial s_{\theta,1}(\mathbf{x})}{\partial x_2} & \frac{\partial s_{\theta,1}(\mathbf{x})}{\partial x_3} \\ \frac{\partial s_{\theta,2}(\mathbf{x})}{\partial x_1} & \frac{\partial s_{\theta,2}(\mathbf{x})}{\partial x_2} & \frac{\partial s_{\theta,2}(\mathbf{x})}{\partial x_3} \\ \frac{\partial s_{\theta,3}(\mathbf{x})}{\partial x_1} & \frac{\partial s_{\theta,3}(\mathbf{x})}{\partial x_2} & \frac{\partial s_{\theta,3}(\mathbf{x})}{\partial x_3} \end{pmatrix} \begin{pmatrix} \nu_1 \\ \nu_2 \\ \nu_3 \end{pmatrix}$$

Computing Jacobian-vector products is scalable

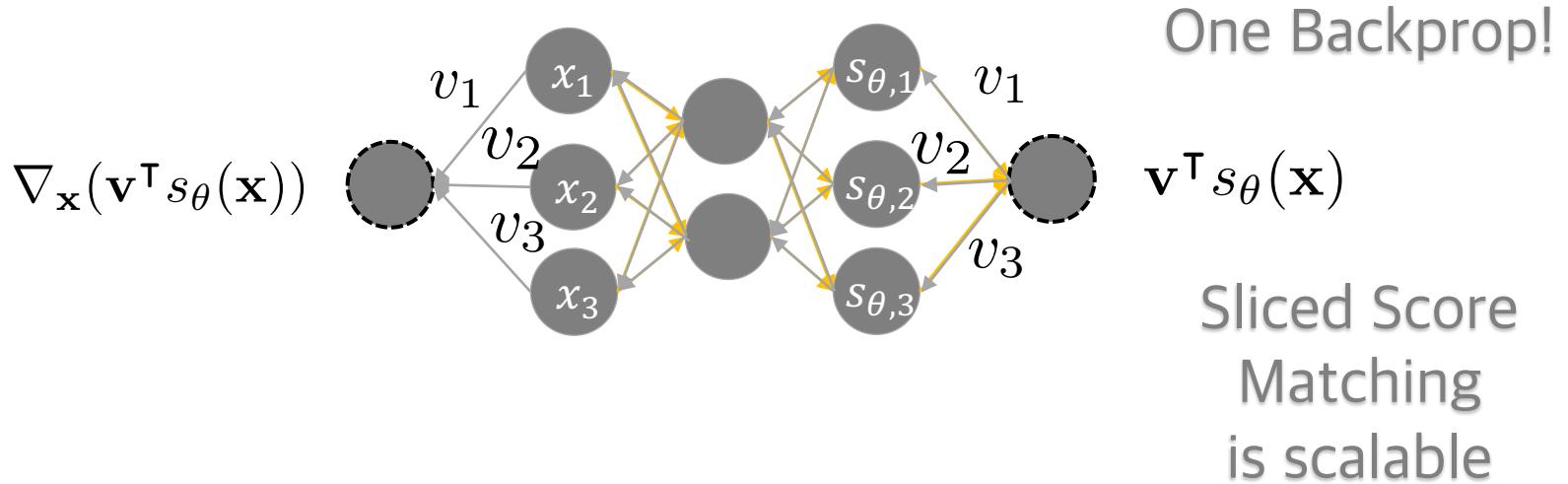
$$\mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v} = \mathbf{v}^T \nabla_{\mathbf{x}} (\mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x}))$$



- Slightly slower than denoising score matching

Computing Jacobian-vector products is scalable

$$\mathbf{v}^T \nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x}) \mathbf{v} = \mathbf{v}^T \nabla_{\mathbf{x}} (\mathbf{v}^T \mathbf{s}_{\theta}(\mathbf{x}))$$



- Slightly slower than denoising score matching

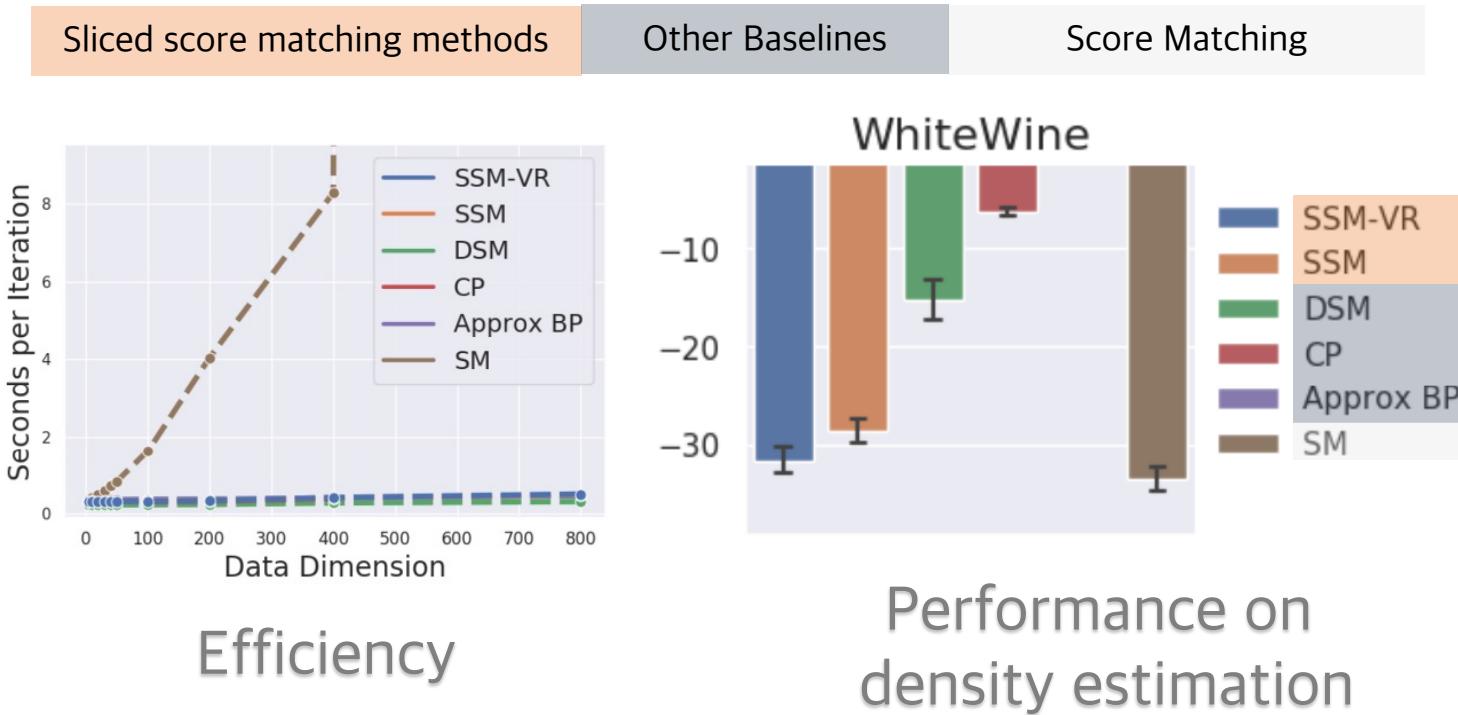
Training of sliced score matching

- Sample a minibatch of datapoints $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ from p_{data}
- Sample a minibatch of projection directions $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(n)}$ from p_v
- Estimate the sliced score matching loss with empirical means

$$\frac{1}{n} \sum_{i=1}^n \mathbf{v}^{(i)T} \nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}^{(i)}) \mathbf{v}^{(i)} + \frac{1}{2} \left(\mathbf{v}^{(i)T} s_{\theta}(\mathbf{x}^{(i)}) \right)^2$$

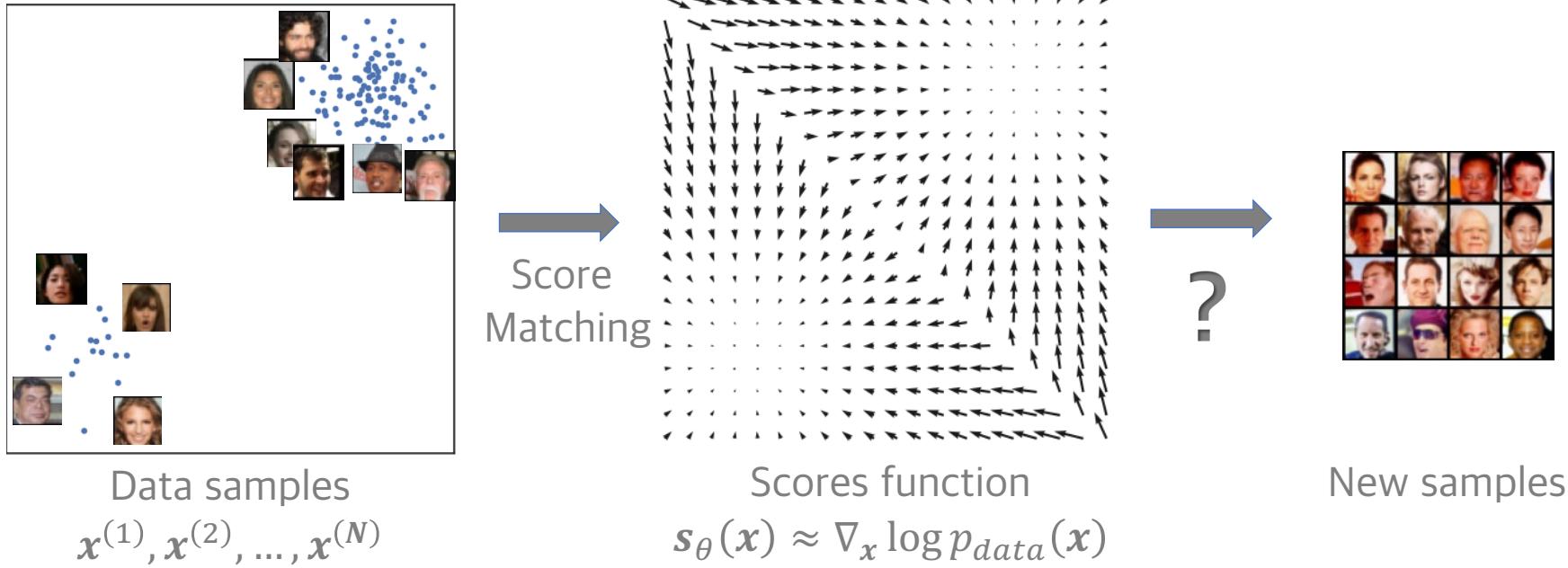
- The projection distribution is typically Gaussian or Rademacher
- Stochastic gradient descent
- Can use more projections per datapoint to boost performance

Experimental results for sliced score matching

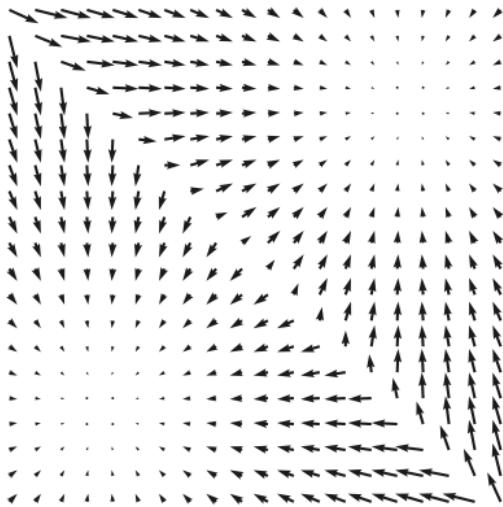


Song*, Garg*, Shi, Ermon. "Sliced Score Matching: A Scalable Approach to Density and Score Estimation." UAI 2019.

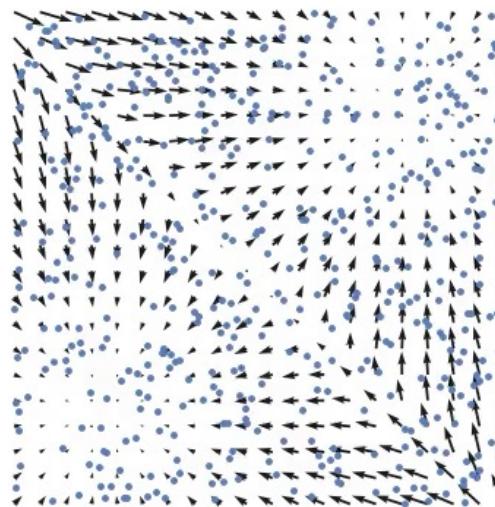
Score-based generative modeling



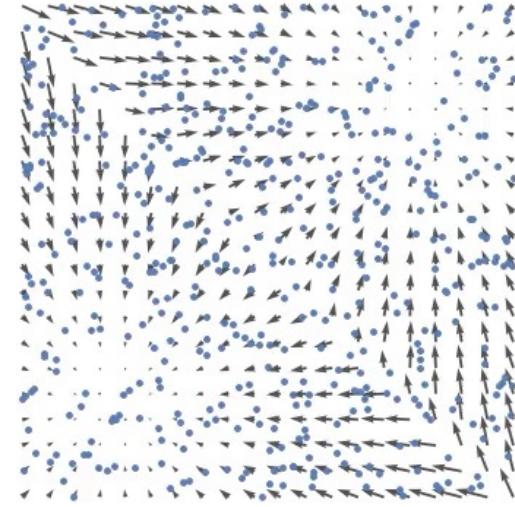
From scores to samples: Langevin MCMC



Scores
 $s_\theta(x)$



Follow the scores
$$x^{t+1} = x^t + \frac{\epsilon}{2} s_\theta(x)$$



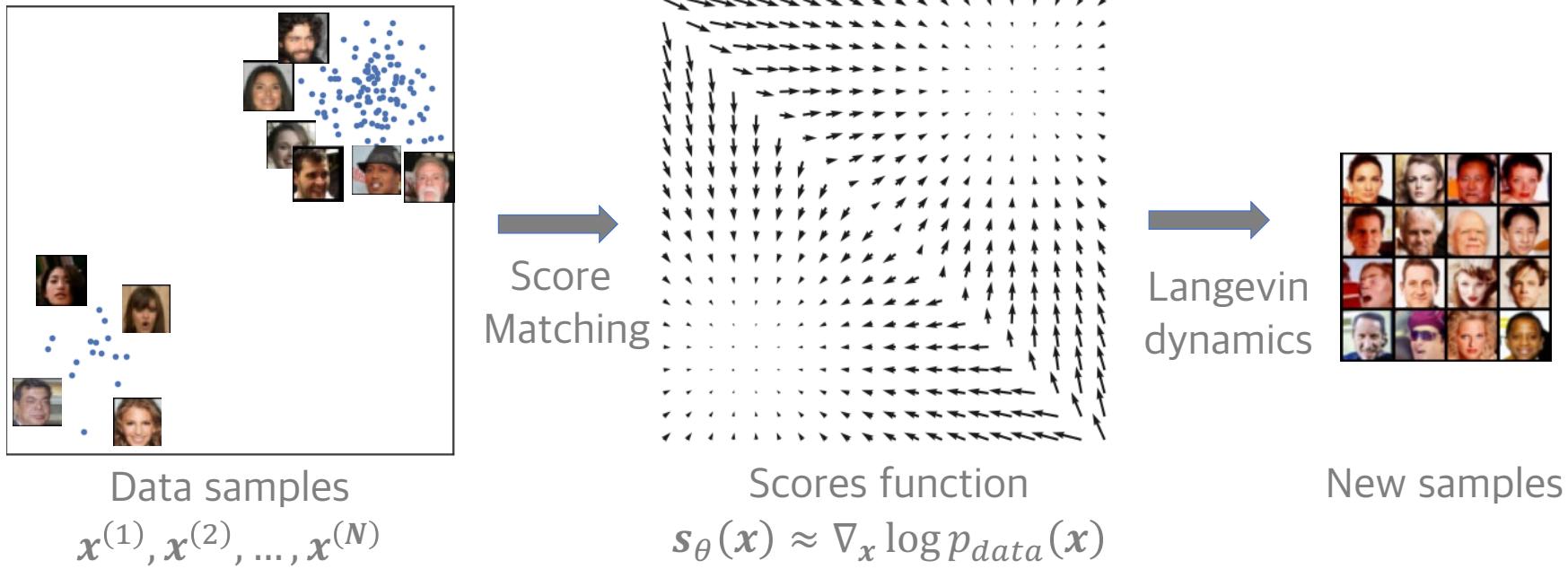
Follow noisy scores:
Langevin MCMC
 $z \sim N(\mathbf{0}, I)$
$$x^{t+1} = x^t + \frac{\epsilon}{2} s_\theta(x) + \sqrt{\epsilon} z$$

Langevin dynamics sampling

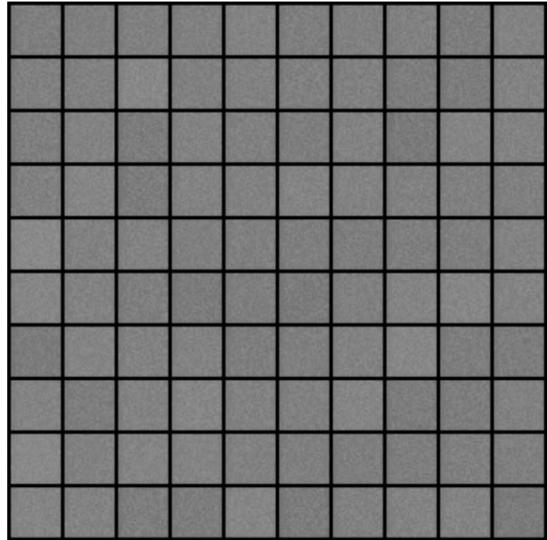
- Sample from $p_{data}(\mathbf{x})$ using only the scores $\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$
 - Initialize $\mathbf{x}^0 \sim \pi(\mathbf{x})$
 - Repeat for $t = 0, \dots, T - 1$
 - $\mathbf{z} \sim N(0, I)$
 - $\mathbf{x}^{t+1} = \mathbf{x}^t + \frac{\epsilon}{2} \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})|_{\mathbf{x}=\mathbf{x}^t} + \sqrt{\epsilon} \mathbf{z}$
 - If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, then we have \mathbf{x}^T converges to a sample from p_{data}
- Langevin dynamics + score estimation

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

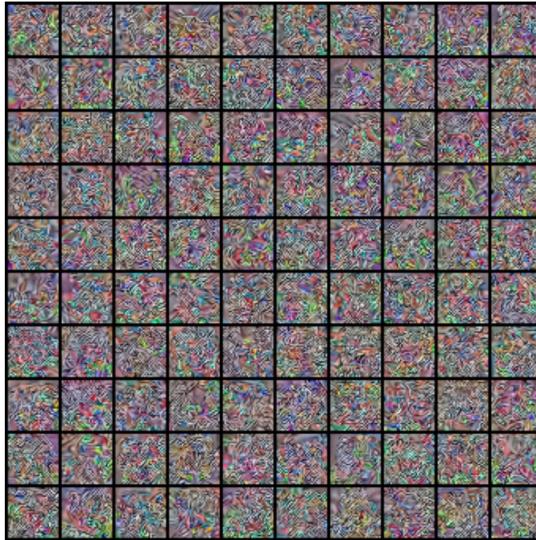
Score-based generative modeling



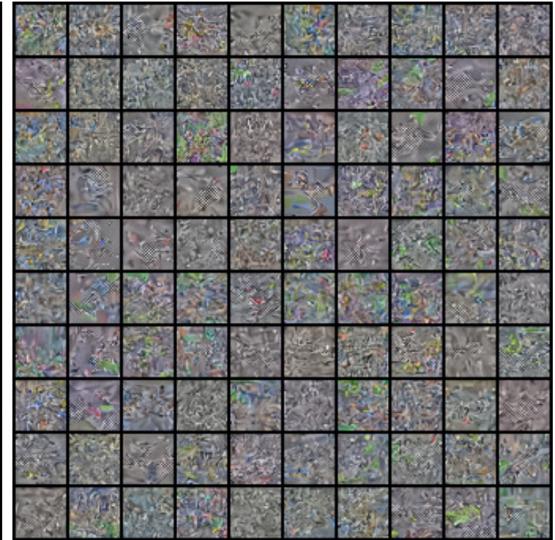
Score-based generative modeling: results



(a) MNIST



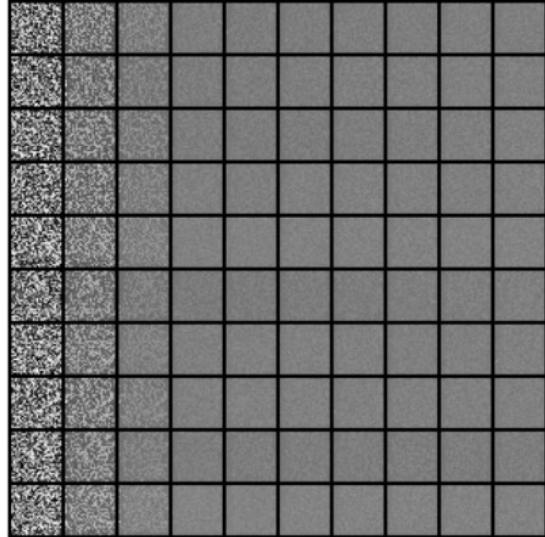
(b) CelebA



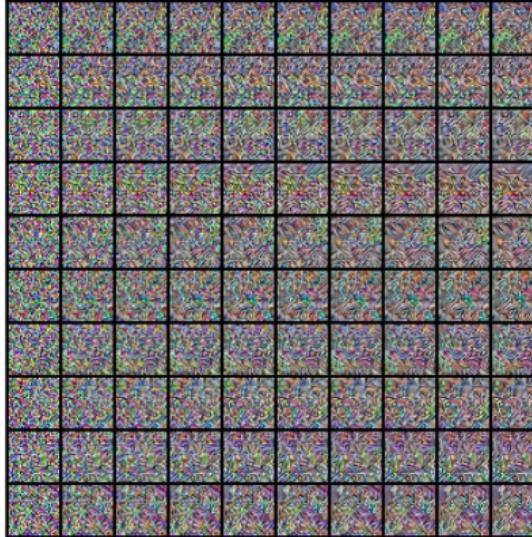
(c) CIFAR-10

Final samples

Score-based generative modeling: results



(a) MNIST



(b) CelebA

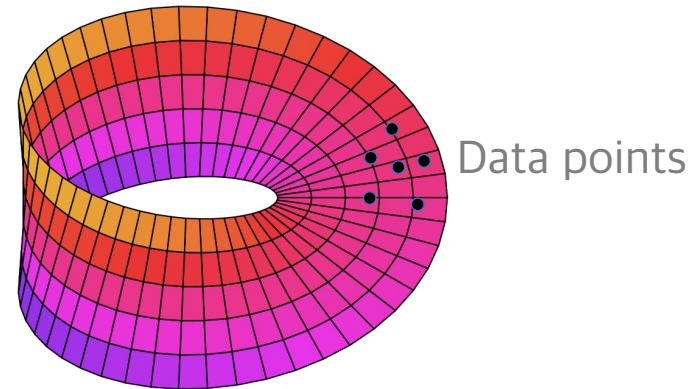


(c) CIFAR-10

Langevin sampling process

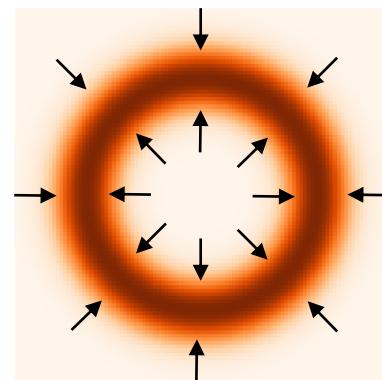
Pitfall 1: manifold hypothesis

- Manifold hypothesis



- Data score is undefined

$$\nabla_x \log p_{data}(x)$$



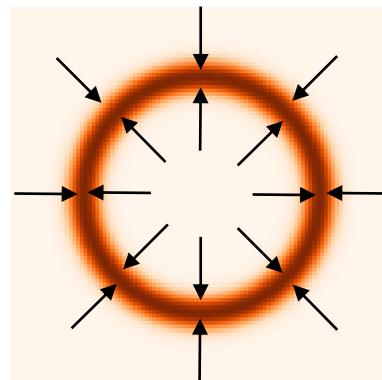
Pitfall 1: manifold hypothesis

- Manifold hypothesis



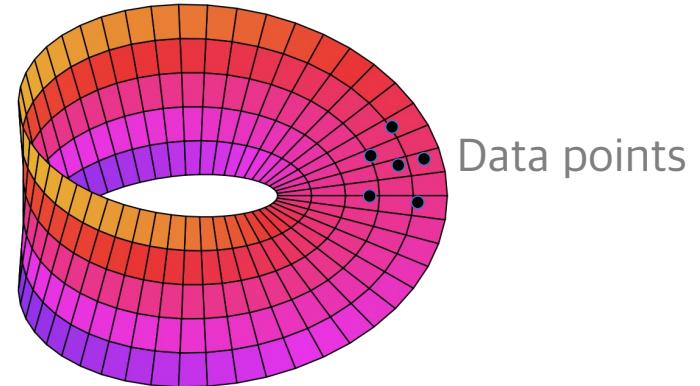
- Data score is undefined

$$\nabla_x \log p_{data}(x)$$



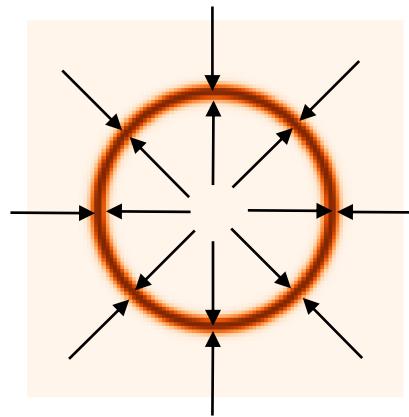
Pitfall 1: manifold hypothesis

- Manifold hypothesis



- Data score is undefined

$$\nabla_x \log p_{data}(x)$$

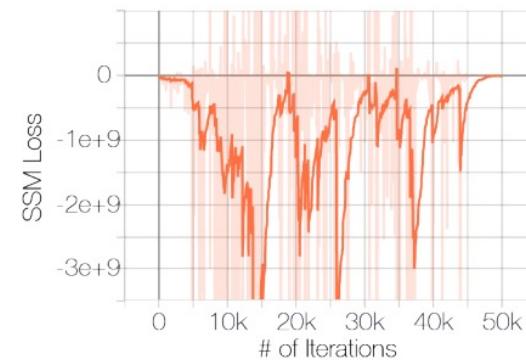


Pitfall 1: manifold hypothesis

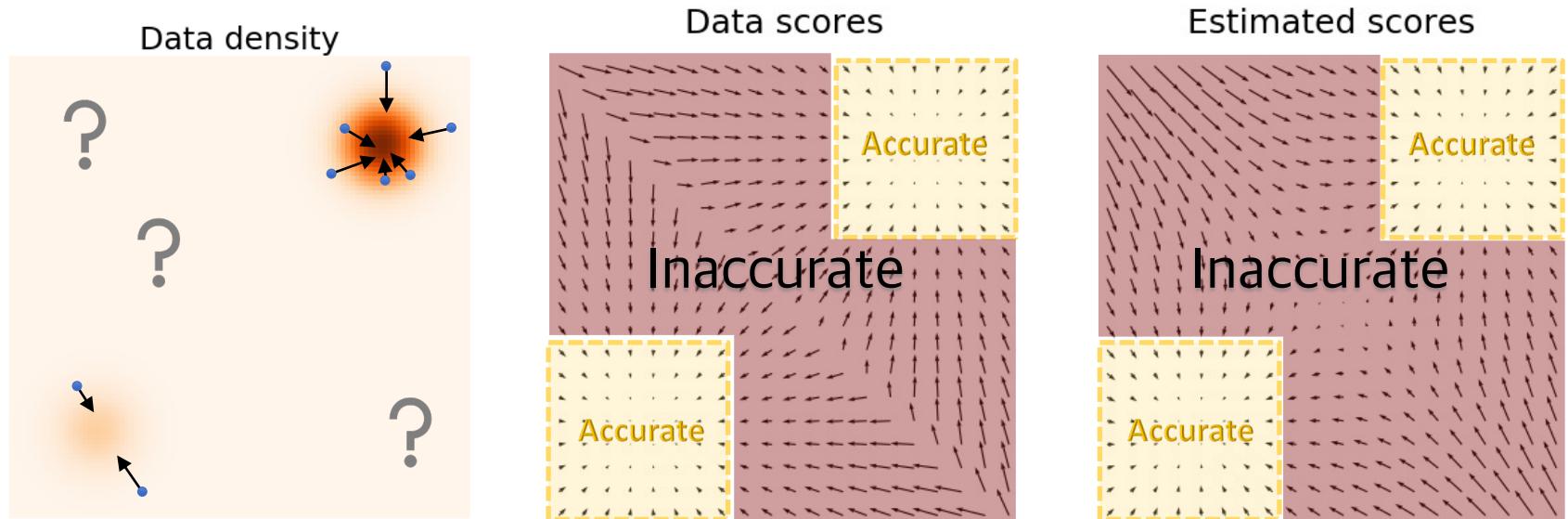
- Fitting the data with a low-dimensional linear manifold (PCA)



- Score estimation on CIFAR-10



Pitfall 2: challenge in low data density regions



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

Langevin MCMC will have trouble
exploring low density regions

Song and Ermon. “Generative Modeling by Estimating
Gradients of the Data Distribution.” NeurIPS 2019.

Pitfall 3: slow mixing of Langevin dynamics between data modes

- Suppose the data distribution has two modes with disjoint supports:

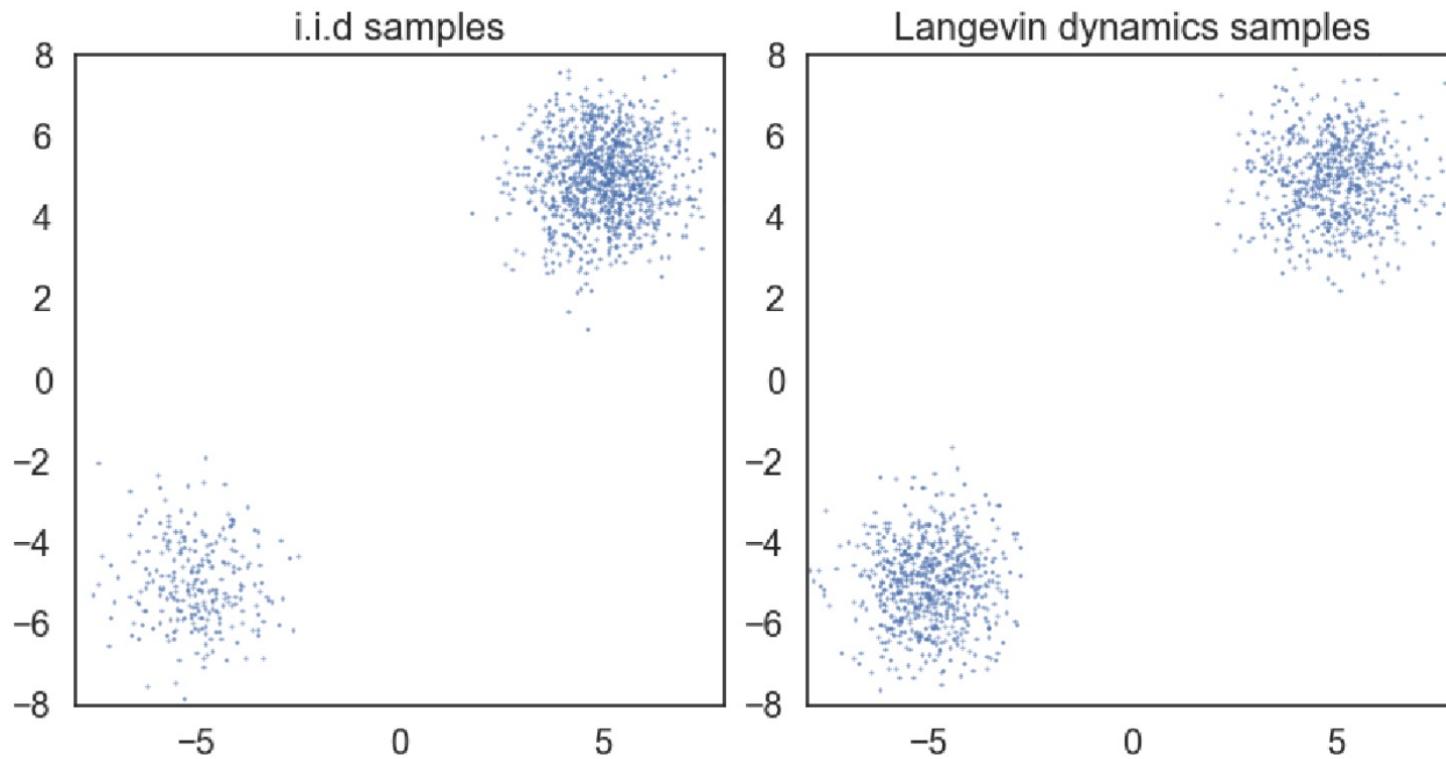
$$p_{data}(x) = \pi p_1(x) + (1 - \pi)p_2(x)$$
$$A \cap B = \emptyset, \quad p_{data}(x) = \begin{cases} \pi p_1(x) & x \in A \\ (1 - \pi)p_2(x) & x \in B \end{cases}$$

- Data score function:

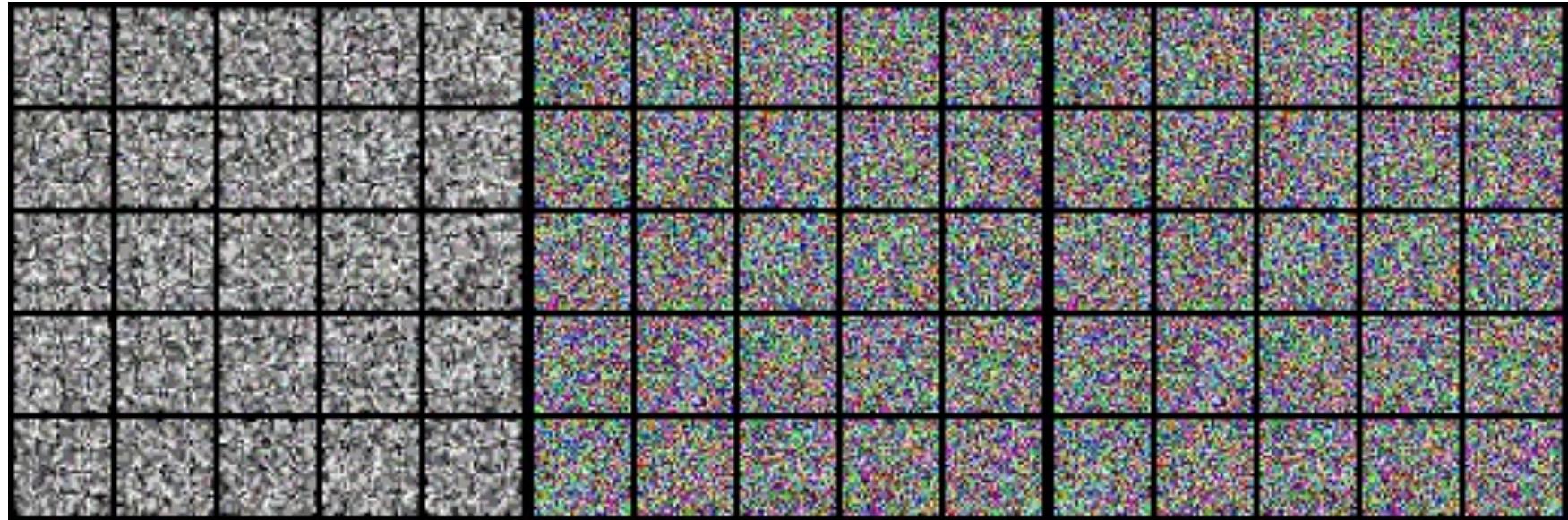
$$\nabla_x \log p_{data}(x) = \begin{cases} \nabla_x \log p_1(x) & x \in A \\ \nabla_x \log p_2(x) & x \in B \end{cases}$$

- The score function has no dependence on π
- Langevin sampling will not reflect π

Pitfall 3: slow mixing of Langevin dynamics between data modes



After fixing pitfalls



Song, Yang, and Stefano Ermon. "Generative Modeling by Estimating Gradients of the Data Distribution." NeurIPS 2019

Thanks
