# **Deep Generative Models**

## 13. 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?



• Inference: how do we invert the generation process

## **Representation: score functions**

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

 $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$ 



## **Score-based models**

• Directly model the vector field of gradients  $s_{\theta}(x) \colon \mathbb{R}^{d} \to \mathbb{R}^{d}$  $s_{\theta}(x) \approx \nabla_{x} \log p_{data}(x)$ 





## **Score estimation**

- Training the score-based model from data points
- Score matching

$$\frac{1}{2} E_{\boldsymbol{x} \sim p_{data}} \left[ \| \nabla_{\boldsymbol{x}} \log p_{data}(\boldsymbol{x}) - \boldsymbol{s}_{\theta}(\boldsymbol{x}) \|_{2}^{2} \right]$$
$$= E_{\boldsymbol{x} \sim p_{data}} \left[ \frac{1}{2} \| \boldsymbol{s}_{\theta}(\boldsymbol{x}) \|_{2}^{2} + tr(\nabla_{\boldsymbol{x}} \boldsymbol{s}_{\theta}(\boldsymbol{x})) \right] + \text{conts.}$$

Not scalable for deep score-based models and high dimensional data

## **Denoising score matching**

 $x \sim p_{data}(x)$ Data distribution



 $\widetilde{\mathbf{x}} \sim q_{\sigma}(\widetilde{\mathbf{x}})$ Noise-perturbed data distribution

 $E_{\widetilde{\mathbf{x}} \sim q_{\sigma}} \left[ \| \nabla_{\widetilde{\mathbf{x}}} \log q_{\sigma}(\widetilde{\mathbf{x}}) - s_{\theta}(\widetilde{\mathbf{x}}) \|_{2}^{2} \right]$ =  $E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\widetilde{\mathbf{x}} \sim q_{\sigma}(\widetilde{\mathbf{x}}|\mathbf{x})} \left[ \| \nabla_{\widetilde{\mathbf{x}}} \log q_{\sigma}(\widetilde{\mathbf{x}}|\mathbf{x}) - s_{\theta}(\widetilde{\mathbf{x}}) \|_{2}^{2} \right] + \text{const.}$ =  $E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})} \left[ \left\| \frac{1}{\sigma} \mathbf{z} + s_{\theta}(\mathbf{x} + \sigma \mathbf{z}) \right\|_{2}^{2} \right] + \text{const.}$ 

- Pros
  - more scalable than score matching
  - reduces score estimation to a denoising task
- Con: cannot estimate the score of clean data (noise-free)

$$\boldsymbol{s}_{\theta}(\boldsymbol{x}) \approx \nabla_{\boldsymbol{x}} \log q_{\sigma}(\boldsymbol{x}) \neq \nabla_{\boldsymbol{x}} \log p_{data}(\boldsymbol{x})$$

## Sliced score matching

• **Objective**: Sliced Fisher Divergence

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

- Projection distribution  $\bar{p}_{v}$  can be Gaussian or Rademacher
- Pros
  - Much more scalable than score matching
  - Estimates the true data score
- Con
  - Slower than denoising score matching

## Score-based generative modeling



## Langevin dynamics sampling

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

 $\boldsymbol{s}_{\theta}(\boldsymbol{x}) \approx \nabla_{\boldsymbol{x}} \log p_{data}(\boldsymbol{x})$ 

## Langevin dynamics sampling

• Using

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

- Initialize  $x^0 \sim \pi(x)$
- Repeat for  $t = 0, \dots, T 1$ 
  - $z \sim N(0, I)$
  - $x^{t+1} = x^t + \frac{\epsilon}{2} s_{\theta}(x)|_{x=x^t} + \sqrt{\epsilon} z$

## Pitfalls

- Manifold hypothesis: data score is undefined  $\nabla_x \log p_{data}(x)$
- Score matching fails in low data density regions



• Langevin dynamics converges very slowly





Deep Generative Models | mjgim@nims.re.kr |

NIMS & AJOU University

## **Gaussian perturbation**

• The solution to all pitfalls: Gaussian perturbation

• Manifold + noise



• Score matching on noisy data



 $N(0, 10^{-4}I)$ 

### **Gaussian perturbation**

 $q_{\sigma}(\widetilde{\boldsymbol{x}}|\boldsymbol{x}) \coloneqq N(\widetilde{\boldsymbol{x}}|\boldsymbol{x}, \sigma^2 \boldsymbol{I}),$ 

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



large  $\sigma$ 

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

 $q_{\sigma}(\mathbf{x})$ 

#### mjgim@nims.re.kr | Deep Generative Models NIMS & AJOU University

## Improving score estimation by adding noise

Perturbed density



#### Perturbed scores



#### Estimated scores



 $x \sim p_{data}(x)$ Data distribution



 $\widetilde{x} \sim q_{\sigma}(\widetilde{x})$ Noise-perturbed data distribution

## Improving score estimation by adding noise





 $\widetilde{x} \sim q_{\sigma}(\widetilde{x})$ Noise-perturbed data distribution

## Multi-scale noise perturbation

• How much noise to add?



• Multi-scale noise perturbations

 $\sigma_L > \sigma_{L-1} > \cdots > \sigma_2 > \sigma_1$ 



## Trading off data quality and estimation accuracy



Worse data quality!

Better score estimation!

## Using multiple noise scales



## Annealed Langevin dynamics: Joint scores to samples

- Sample using  $\sigma_L > \sigma_{L-1} > \cdots > \sigma_2 > \sigma_1$  sequentially with Langevin dynamics
- Anneal down the noise level

**Deep Generative Models** 

• Samples used as initialization for the next level



mjgim@nims.re.kr

NIMS & AJOU University

# Joint score estimation via noise conditional score networks



## **Annealed Langevin dynamics**

**Algorithm 1** Annealed Langevin dynamics. **Require:**  $\{\sigma_i\}_{i=1}^L, \epsilon, T.$ 1: Initialize  $\tilde{\mathbf{x}}_0$ 2: for  $i \leftarrow 1$  to L do 3:  $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2 \qquad \triangleright \alpha_i$  is the step size. 4: for  $t \leftarrow 1$  to T do 5: Draw  $\mathbf{z}_t \sim \mathcal{N}(0, I)$  $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_{\theta}(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$ 6: 7: end for  $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$ 8: 9: end for return  $\tilde{\mathbf{x}}_T$ 

## **Comparison to the vanilla Langevin dynamics**



Langevin dynamics



Annealed Langevin dynamics

## **Comparison to the vanilla Langevin dynamics**



Langevin dynamics

Annealed Langevin dynamics

## Training noise conditional score networks

- Denoising score matching is naturally suitable, since the goal is to estimate the score of perturbed data distributions
- Weighted combination of denoising score matching losses

$$\frac{1}{L} \sum_{i=1}^{L} \lambda(\sigma_{i}) E_{\widetilde{\mathbf{x}} \sim q_{\sigma_{i}}(\widetilde{\mathbf{x}})} \left[ \left\| \nabla_{\widetilde{\mathbf{x}}} \log q_{\sigma_{i}}(\widetilde{\mathbf{x}}) - s_{\theta}(\widetilde{\mathbf{x}}, \sigma_{i}) \right\|_{2}^{2} \right] \\
= \frac{1}{L} \sum_{i=1}^{L} \lambda(\sigma_{i}) E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\widetilde{\mathbf{x}} \sim q_{\sigma_{i}}(\widetilde{\mathbf{x}}|\mathbf{x})} \left[ \left\| \nabla_{\widetilde{\mathbf{x}}} \log q_{\sigma_{i}}(\widetilde{\mathbf{x}}|\mathbf{x}) - s_{\theta}(\widetilde{\mathbf{x}}, \sigma_{i}) \right\|_{2}^{2} \right] \\
+ \text{ const.}$$

## **Choosing noise scales**

- Maximum noise scale
  - $\sigma_L \approx$  maximum pairwise distance between datapoints



• Minimum noise scale:  $\sigma_1$  should be sufficiently small so that noise in final samples is negligible. I.e.,  $q_{\sigma_1}(\mathbf{x}) \approx p_{data}(\mathbf{x})$ 

## **Choosing noise scales**

- Key intuition: adjacent noise scales should have sufficient overlap to facilitate transitioning across noise scales in annealed Langevin dynamics
- A geometric progression with sufficient length

$$\sigma_L > \sigma_{L-1} > \dots > \sigma_2 > \sigma_1$$
$$\frac{\sigma_1}{\sigma_2} = \frac{\sigma_2}{\sigma_3} = \dots = \frac{\sigma_{L-1}}{\sigma_L}$$

## **Choosing the weighting function**

- Weighted combination of denoising score matching losses
  - $\frac{1}{L}\sum_{i=1}^{L}\lambda(\sigma_{i})E_{\boldsymbol{x}\sim p_{data}(\boldsymbol{x})}E_{\widetilde{\boldsymbol{x}}\sim q_{\sigma_{i}}(\widetilde{\boldsymbol{x}}|\boldsymbol{x})}\left[\left\|\nabla_{\widetilde{\boldsymbol{x}}}\log q_{\sigma_{i}}(\widetilde{\boldsymbol{x}}|\boldsymbol{x}) \boldsymbol{s}_{\theta}(\widetilde{\boldsymbol{x}},\sigma_{i})\right\|_{2}^{2}\right]$  $=\frac{1}{L}\sum_{i=1}^{L}\lambda(\sigma_{i})E_{\boldsymbol{x}\sim p_{data}(\boldsymbol{x})}E_{\boldsymbol{z}\sim N(\boldsymbol{0},\boldsymbol{I})}\left[\left\|\frac{\boldsymbol{z}}{\sigma_{i}} + \boldsymbol{s}_{\theta}(\boldsymbol{x}+\sigma_{i}\boldsymbol{z},\sigma_{i})\right\|_{2}^{2}\right]$
- How to choose the weighting function  $\lambda: \mathbb{R}_+ \to \mathbb{R}_+$ ?
- Goal: balancing different score matching losses in the sum

## **Choosing the weighting function**

- How to choose the weighting function  $\lambda: \mathbb{R}_+ \to \mathbb{R}_+$ ?
- Goal: balancing different score matching losses in the sum
- Since  $\sigma_i^2 \propto 1/E \left[ \left\| \nabla_{\widetilde{\mathbf{x}}} \log p_{\sigma_i}(\widetilde{\mathbf{x}} | \mathbf{x}) \right\|_2^2 \right], \, \lambda(\sigma_i) \coloneqq \sigma_i^2$  $\frac{1}{L}\sum_{i}^{L}\sigma_{i}^{2}E_{\boldsymbol{x}\sim p_{data}(\boldsymbol{x})}E_{\boldsymbol{z}\sim N(\boldsymbol{0},\boldsymbol{I})}\left[\left\|\frac{\boldsymbol{z}}{\sigma_{i}}+\boldsymbol{s}_{\theta}(\boldsymbol{x}+\sigma_{i}\boldsymbol{z},\sigma_{i})\right\|_{2}^{2}\right]$  $= \frac{1}{L} \sum_{i=1}^{L} E_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})} E_{\boldsymbol{z} \sim N(\boldsymbol{0}, \boldsymbol{I})} [\|\boldsymbol{z} + \sigma_{i} \boldsymbol{s}_{\theta}(\boldsymbol{x} + \sigma_{i} \boldsymbol{z}, \sigma_{i})\|_{2}^{2}]$  $= \frac{1}{L} \sum_{i=1}^{L} E_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})} E_{\boldsymbol{z} \sim N(\boldsymbol{0}, \boldsymbol{I})} [\|\boldsymbol{z} + \boldsymbol{\epsilon}_{\theta}(\boldsymbol{x} + \sigma_{i} \boldsymbol{z}, \sigma_{i})\|_{2}^{2}]$ • where  $\epsilon_{\theta}(\cdot, \sigma_i) \coloneqq \sigma_i s_{\theta}(\cdot, \sigma_i)$

## Training noise conditional score networks

- Sample a mini-batch of datapoints  $x^{(1)}, x^{(2)}, \cdots, x^{(n)}$  from  $p_{data}$
- Sample a mini-batch of noise scale indices

 $i_1, i_2, ..., i_n \sim U(1, 2, \cdots, L)$ 

- Sample a mini-batch of Gaussian noise z<sup>(1)</sup>, z<sup>(2)</sup>, ..., z<sup>(n)</sup> from N(0, I)
- Estimate the weighted mixture of score matching losses

$$\frac{1}{n}\sum_{k=1}^{n} \left[ \left\| \boldsymbol{z}^{(k)} + \boldsymbol{\epsilon}_{\theta} \left( \boldsymbol{x}^{(k)} + \sigma_{i_{k}} \boldsymbol{z}^{(k)}, \sigma_{i_{k}} \right) \right\|_{2}^{2} \right]$$

- Stochastic gradient descent
- As efficient as training one single non-conditional score-based model

## **Generation with annealed Langevin dynamics**

• For each  $q_{\sigma_i}(\mathbf{x})$  with  $\sigma_1 < \sigma_2 < \cdots < \sigma_L$ , Song & Ermond run T steps of Langevin MCMC to get a sample sequentially

$$\boldsymbol{x}_{i}^{t} \coloneqq \boldsymbol{x}_{i}^{t-1} + \frac{\alpha_{i}}{2} \boldsymbol{s}_{\theta^{*}} (\boldsymbol{x}_{i}^{t-1}, \sigma_{i}) + \sqrt{\alpha_{i}} \boldsymbol{z}, \qquad t = 1, 2, \dots, T$$

• where  $\alpha_i > 0$  is the step size and  $z \sim N(0, I)$ 

$$\alpha_i \coloneqq \epsilon \frac{\sigma_i^2}{\sigma_1^2}$$

•  $\epsilon > 0$ 

## Using multiple noise levels



## **Experiments:** sampling



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

# Thanks