

Generative model with gradient...

① Definition of score: $\nabla_x \log p(x)$.

Two ingredients:

- Langevin dynamics.
- Score Matching.

② Langevin dynamics: $\tilde{x}_0 \sim \pi(x)$. π , a prior distribution.

$$\tilde{x}_t = \tilde{x}_{t-1} + \frac{\varepsilon}{2} \nabla_x \log p(\tilde{x}_{t-1}) + \sqrt{\varepsilon} \tilde{z}_t.$$

$\tilde{z}_t \sim N(0, 1)$.

recursively compute. when $\varepsilon \rightarrow 0$ and $T \rightarrow \infty$

when ε is small and T is large, $\tilde{x}_t \sim p(x)$.

Instead of compute $p(x)$, only requires score function.

First train the score network. $S_\theta(x) \approx \nabla_x \log p(x)$.

③ Score Matching: $\mathcal{L} \in \mathbb{R}$. Basic score matching.

Objective function: $\mathcal{L} = \frac{1}{2} E[\|\nabla_x \log p_{\text{data}}(x) - S_\theta(x)\|_2^2]$.

which is equivalent to:

$$\mathcal{L} = E_{p(x)}[\text{tr}(\nabla_x S_\theta(x)) + \frac{1}{2} \|S_\theta(x)\|_2^2].$$

To avoid jacobian of $S_\theta(x)$.

use denoising score matching.

$$\text{denoising score matching: } \frac{1}{2} E[\|\nabla_x \log p_{\text{data}}(x) - S_\theta(x)\|_2^2]$$

$$= \frac{1}{2} E[\underbrace{\|\nabla_x \log p_{\text{data}}(x)\|_2^2}_{\text{constant}} - \underbrace{2 \nabla_x \log p_{\text{data}}(x) \cdot S_\theta(x)}_{\mathcal{L}_2} - \underbrace{\|S_\theta(x)\|_2^2}_{\mathcal{L}_3}].$$

$$L_2 = \int p_{\text{data}}(x) \nabla_x \log p_{\text{data}}(x) \cdot S_\theta(x) dx \Rightarrow \underline{x^T y = \frac{1}{n} \sum_{i=1}^n x_i y_i}$$

$$= \int p_{\text{data}}(x) \sum_{x_i} \nabla_{x_i} \log p_{\text{data}}(x_i) \cdot \sum_{x_i} S_\theta(x_i) dx$$

$$= \sum_{x_i} \int p_{\text{data}}(x_i) \nabla_{x_i} \log p_{\text{data}}(x_i) S_\theta(x_i) dx \quad \underline{\text{求导}}$$

$$= \sum_{x_i} \int p_{\text{data}}(x_i) \frac{\log p_{\text{data}}(x_i)}{p_{\text{data}}(x_i)} S_\theta(x_i) dx$$

$$= \sum_{x_i} \int \log p_{\text{data}}(x_i) S_\theta(x_i) dx \quad \rightarrow \text{分部积分:}$$

$$= \sum_{x_i} \left(\int_{\text{data}} p(x) S_\theta(x) - \int p_{\text{data}} d S_\theta(x) \right) \quad \begin{matrix} \int u dv \\ = uv - \int v du \end{matrix}$$

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

$$\psi_{\theta_i} \left(\frac{x}{\sigma} \right)$$

④. Denoising score Matching:

First perturb x with a pre-specified noise distribution.

$p_{\tilde{x}}(\tilde{x})$, then use score matching.

可以理解成
diffusion model?

$$p_{\tilde{x}}(\tilde{x}) \triangleq \int p_{\text{data}}(x) p(\tilde{x}|x) dx \quad \leftarrow \text{可以理解为 } p(x_t|x_{t+1})$$

Objective function:

$$\frac{1}{2} E_{p_{\tilde{x}}(\tilde{x}|x) p(x)} \left[\|S_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log p_{\tilde{x}}(\tilde{x}|x)\|_2^2 \right]$$

$$S_\theta^*(x) \approx \nabla_x \log p_{\tilde{x}}(x)$$

$$S_\theta^*(x) \approx \nabla_x \log p(x)$$

Only when noise very small.

在数据密度较低的地方, 采样不准确.

④ Noise conditional Score Network (NCSN).

给密度较低的地方加较大扰动噪声, 但加太大会失去数据原本特征.

所以用一个 conditional 的 Network., with multiple scale noise of perturbations simultaneously.

When using Langevin dynamics to generate samples, first uses score corresponding to large noise, then gradually anneal down the noise level.

Add isotropic Gaussian noise: $\{\sigma_i\}_{i=1}^L, \sigma_1 < \sigma_2 < \dots < \sigma_L$.

$$\mathbf{p}_{\sigma_i} \sim N(0, \sigma_i^2 \mathbf{I}) \quad i=1, 2, \dots, L$$

$$p_{\sigma}(x) = \int p_{\text{data}}(x) N(x|t, \sigma^2 \mathbf{I}) dt$$

Train a network to jointly estimate the score of all perturbed data distributions S_{θ} . Conditional Score Network (CSN).