

<Target> : 확산모델 (다양성 확보!)

DDPM

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

ICML 2020

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.

<Reference>

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

<https://cvpr2022-tutorial-diffusion-models.github.io/>

Denoising Diffusion Models

Learning to generate by denoising

DDPM (2020)



DDIM (2020)



Improved DDPM (2021)
(Nichol et al.)



Latent Diffusion model
(2021)



stable Diffusion (2022)

<history ⇒ past ~ to SOTA>

1. Denoising Diffusion Probabilistic Models (DDPM, 2020):

- 제안자: Ho et al.
- 개요: 확산 과정과 역확산 과정을 통해 이미지 데이터를 점진적으로 노이즈화하고, 이를 역으로 복원하여 이미지를 생성하는 모델입니다.
- 주요 특징: Markovian 확산 과정과 비슷한 역확산 과정을 사용합니다.

2. Denoising Diffusion Implicit Models (DDIM, 2020):

- 제안자: Song et al.
- 개요: DDPM의 효율성을 개선한 모델로, 직접적인 역확산을 통해 빠른 샘플링을 가능하게 합니다.
- 주요 특징: 비-Markovian 샘플링 방법을 사용하여 더 빠르고 효율적인 이미지 생성을 구현합니다.

3. Improved Denoising Diffusion Probabilistic Models (Improved DDPM, 2021):

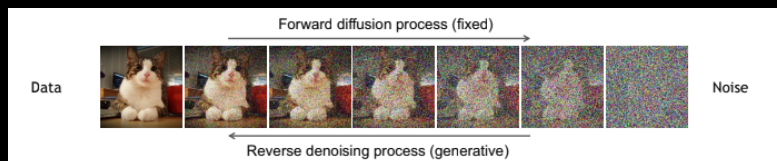
- 제안자: Nichol and Dhariwal
- 개요: DDPM의 성능을 개선하기 위해 하이퍼파라미터 튜닝과 아키텍처 수정 등을 도입한 모델입니다.
- 주요 특징: 성능 향상과 더 높은 품질의 이미지 생성을 목표로 합니다.

4. Latent Diffusion Models (LDM, 2021):

- 제안자: Rombach et al.
- 개요: 고해상도 이미지를 생성하기 위해 잠재 공간(latent space)에서 확산 과정을 수행하는 모델입니다.
- 주요 특징: 효율성을 높이기 위해 잠재 공간에서의 연산을 활용합니다.

5. Stable Diffusion (2022):

- 제안자: Stability AI
- 개요: LDM의 발전된 형태로, 높은 품질의 이미지를 빠르게 생성할 수 있는 모델입니다.
- 주요 특징: 일반 사용자도 쉽게 사용할 수 있도록 접근성을 높인 모델입니다.



Learning to Generate by Denoising

이는 정제된 기본 SD(2.2.2)의
3가지 변형인 DDPM(2.2.2.1)

< What Diffusion Model Does? >

· 확산 모델은 노이즈를 생생하게 정적으로 학습에 추가하고
이를 역으로 복원하면서 생생한 데이터를 생생 (출처 GPT-4)

학습 과정 (DDPM)

① Forward Diffusion Process (정적 노이즈화 과정)

입력 데이터를 정적으로 노이즈화 시켜 x_T 에 도달함

($t=1, \dots, T$)

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t) I)$$

여기서 α_t 는 노이즈의 크기 조절 파라미터

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}) \dots (1)$$

→ 전체 노이즈화 과정

② Reverse Diffusion process (역 확산 과정)

노이즈화된 데이터 x_T 에서 원본 데이터 x_0 로 복원 하는 과정

여기서 역확산 과정의 확률 분포를 모델링 함

$$(2) \dots \rightarrow p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

($p_\theta, \mu_\theta, \Sigma_\theta$) 가 생명을 통해 학습되는 파라미터들.

③ 파라미터 ($p_\theta, \mu_\theta, \Sigma_\theta$)를 Gradient Descent로 최적화.

역확산 과정의 파라미터를 변분 방법 (VB Variational Bound)
를 최소화 하는 방향

Very similarly to VAE and thus we can use the variational lower bound to optimize the negative log-likelihood.

$$\begin{aligned}
 -\log p_{\theta}(\mathbf{x}_0) &\leq -\log p_{\theta}(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{1:T}|\mathbf{x}_0)) \geq 0 \quad \leftarrow \text{Point !!!} \\
 &= -\log p_{\theta}(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})/p_{\theta}(\mathbf{x}_0)} \right] \\
 &= -\log p_{\theta}(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} + \log p_{\theta}(\mathbf{x}_0) \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \quad \text{여기에서 } q \text{ 는 } \mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0) \\
 &\quad = q(x_1, x_2, \dots, x_T)
 \end{aligned}$$

05_03 Lecture note

이 16 짜리

<loss function>

$$\text{Let } L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \geq -\mathbb{E}_{q(\mathbf{x}_0)} \log p_{\theta}(\mathbf{x}_0)$$

$$\mathcal{L}_{\text{no}} = \mathbb{E}_{\mathbf{x}_{0:T}} \left[-\log p_{\theta}(\mathbf{x}_{0:T}) + \log q(\mathbf{x}_{1:T}|\mathbf{x}_0) \right] \dots (3)$$

→ (1) A forward diffusion process

→ (2) ⇐ (Reverse diffusion process)

why? \mathcal{L}_{no} ?

→ \mathbf{x}_0 가 주어졌을 때 $D_{\text{KL}}(q \| p_{\theta})$ 를 최소화 하는게 목표임

(p_{θ} 를 q 분포와 최대한 가깝게 하는것임)

$$\rightarrow \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \text{ 최소화 } \dots (4)$$

$$= D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{1:T}|\mathbf{x}_0)) \text{를 최소화}$$

→ $q(\mathbf{x}_{1:T}|\mathbf{x}_0) \approx p_{\theta}(\mathbf{x}_{1:T}|\mathbf{x}_0)$ 를 만드는것

\mathbf{x}_0 가 Given 되었을 때!

< \mathcal{L}_{VLB} 식짜기. 주(4)를 풀어서 쓰면 >

$$\begin{aligned}
 \mathcal{L}_{VLB} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\
 &= \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left(\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[\underbrace{\log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)}}_{\text{term 1}} + \sum_{t=2}^T \underbrace{\log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}}_{\text{term 2}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{\text{term 3}} \right]
 \end{aligned}$$

\mathbf{x}_T 은 이제 거의 Noise이다.
 $q(\mathbf{x}_T|\mathbf{x}_0)$ 나 $p_\theta(\mathbf{x}_T)$ 나 거의 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ 로
 의미가 없다.

실제로 최소화하는 항이다.

μ 은 매우 작게 주면
 \mathbf{x}_T 은 거의 \mathbf{x}_0 만큼
 관계가 가능한 항

$$\mathcal{L}_{VLB} = \mathbb{E}_q \left[\text{상수} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \text{상수} \right]$$

우리가 $\mathbf{x}_0, \mathbf{x}_t$ 가 있는 항을 구할 수 있음

$$\begin{aligned}
 q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N} \left(\frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_t \right), \tilde{\beta}_t \mathbf{I} \right) \\
 p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) &= \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)) \\
 &= \mathcal{N} \left(\frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(\mathbf{x}_t, t) \right), \Sigma_\theta(\mathbf{x}_t, t) \right)
 \end{aligned}$$

대입

$$\mathbb{E}_q \left[\sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

$$\Rightarrow \mathcal{L}_{t+1} = \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[\frac{(1-\alpha_t)^2}{2\alpha_t(1-\alpha_t) \|\Sigma_\theta\|_F^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\epsilon_t, t)\|^2 \right]$$

우리는 \mathcal{L}_{VLB} 를 계산할 예정.

< Simplified L_{no} : Ignoring the weighting term >

training. However, we found it beneficial to sample quality (and simpler to implement) to train on the following variant of the variational bound:

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\left\| \epsilon - \epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t) \right\|^2 \right] \quad (14)$$

↳ 그냥 이렇게 식도 증명은 Empirically 증명.

Ho et al. [NeurIPS 2020](#) observe that simply setting $\lambda_t = 1$ improves sample quality. So, they propose to use:

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(1, T)} \left[\left\| \epsilon - \underbrace{\epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t)}_{\mathbf{x}_t} \right\|^2 \right]$$

For more advanced weighting see [Choi et al., Perception Prioritized Training of Diffusion Models, CVPR 2022.](#)

weighting term은 아래 참조

< Summary > → ① Training
② Sample Generation (Inference)

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t) \right\|^2$ 
6: until converged
  
```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
  
```