



Denoising Diffusion Probabilistic Models (2020-NeurlPS)





Jonathan Ho

Unknown affiliation Verified email at berkeley.edu - <u>Homepage</u> Artificial Intelligence Machine Learning

TITLE	CITED BY	YEAR
Denoising diffusion probabilistic models J Ho, A Jein, P Abbeel Advances in neural information processing systems 33, 6848-6851	22120	2020
Photorealistic text-to-image diffusion models with deep language understanding C Saharia, W Chan, S Saxona, L U, J Whang, EL Denton, K Ghasemipour, Advances in neural information processing systems 35, 36479-36494	6574	2022
Classifier-free diffusion guidance J Ho, T Salimans arXiv preprint arXiv:2207.12598	4344	2022
Generative adversarial imitation learning J Ho., S Ermon Advances In Neural Information Processing Systems, 4565-4573	4100	2016
Image super-resolution via iterative refinement C. Saharis, J. Ho, W. Chan, T. Salimans, DJ. Fleet, M. Norouzi IEEE transactions on pattern analysis and machine intelligence 45 (4), 4713-4726	2126	2022
Evolution strategies as a scalable alternative to reinforcement learning T Salmans, J Ho, X Chan, S Stider, I Salakavar arXiv preprint arXiv:1703.03864	1932	2017
Video diffusion models J Ho, T Salimans, A Gritsonko, W Chan, M Noreuzi, DJ Floet Advances in Neural Information Processing Systems 35, 8633-8646	1813	2022
Palette: Image-to-image diffusion models C.Saharia, W.Chan, H.Chang, C.Lee, J.Ho, T.Salimans, D.Fleet, ACM SISGRAPH 2822 conference proceedings, 1-18	1667	2022
Progressive distillation for fast sampling of diffusion models T Safimans, J Ho arXiv preprint arXiv:2202.00512	1365	2022
Cascaded Diffusion Models for High Fidelity Image Generation J Ho, C Saharia, W Chan, TJ Float, M Norouzi, T Salimans arXiv preprint arXiv:2106.15282	1315	2021
Variational diffusion models 0 Kingma, T Salimann, 8 Poole, J Ho Advances in neural information processing systems 34, 21696-21787	1206	2021
Motion planning with sequential convex optimization and convex collision checking J Schulman, Y Duan, J Ho, A Lee, 1 Awwal, H Bradlow, J Part, S Patil, The fotemational Journal of Robotics Research 33 (9), 1251-1270	1047	2014

F FOLLOW

Contents.



- 1 Background
- 2 Diffusion
- **3** Diffusion Process
- 4 Result & Summary
- 5 Code



1 Background

Background

Markov Chain

- 시간에 따라서 상태가 변하는 시스템을 모델링하는 방법
- Markov Property : 특정 시점 t+1 에서의 상태는 오직 바로 직전시점 t 상태에만 의존하고 이전 시점과는 독립적이다.

$$P(X_{t+1} = x_{t+1} | X_t = x_t, X_{t-1} = x_{t-1}) = P(X_{t+1} = x_{t+1} | X_t = x_t)$$

KL-Divergence

• 두 확률분포 P, Q가 얼마나 다른지를 측정하는 척도

$$KL(P||Q) = E_p[lop\frac{P(x)}{Q(x)}] = -\sum P(x)log\frac{P(x)}{Q(x)} = H(P,Q) - H(P)$$

$$P, Q \cap Cross-Entropy$$

MLE(Maximum Likelihood Estimation)

- Likelihood : 주어진 데이터 X가 특정 확률분포로부터 나왔을 가능성을 나타내는 값
- MLE : 관측된 데이터 X를 가장 잘 설명하는(Likelihood를 최대로 만드는) θ 를 찾는 방법

Background

Bayes Rule

• Bayes Rule : 주어진 정보를 바탕으로 사전 확률을 업데이트하여 사후 확률을 얻는 방법

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

• P(A): 사전확률 (사건 B에 대한 정보 없이 사건 A가 발생할 확률)

• *P(B)*: 사건 B가 발생할 확률

• P(B|A): likelihood

• P(A|B) : 사후 확률

→ 순방향 과정을 바탕으로 역방향을 추론한다.

확률의 연쇄법칙

• 확률의 변수가 여러 개일때 확률을 조건부 확률의 곱으로 분해하는 방법

$$P(X_1, X_2, \dots, X_n) = P(X1) * P(X2|X1) * P(X_3|X_1, X_2) * \dots * P(X_n|X_{1,\dots}X_{n-1})$$



2 Diffusion

Diffusion

Overview

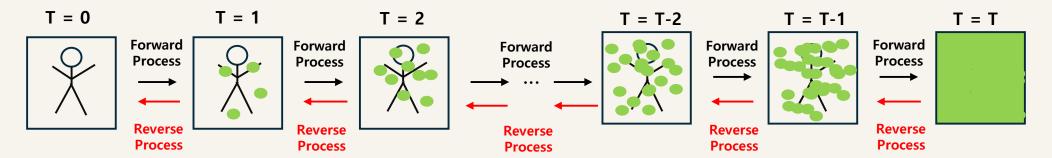


- 잉크 분자들은 시간이 지남에 따라서 물 전체에 고르게 퍼져(diffusion) 나감 (Uniform 분포)
- 잉크 분자들이 움직이는 단계를 확대해서 보면 움직임의 확률 분포는 Gaussian 분포를 따르는것으로 근사할 수 있음
- 잉크 분자들의 움직임의 확률분포를 알 수 있다면 다시 되돌리는것도 가능함
- 따라서 움직임의 확률분포를 Neural Net으로 근사하여 해당 확률분포를 학습하는것이 Diffusion의 목표



Overview

- 패턴을 학습하기 위해 고의적으로 패턴을 붕괴 (Nosing), 이를 다시 복원하도록(reverse) Neurl Network를 학습
- Nosing 과정: Forward(Diffusion) Process (가우시안 노이즈를 추가)
- Denosing 과정 : Reverse Process



- Reverse process : $P_{ heta}(X_{0:T}) \coloneqq P(X_T) \prod_{t=1}^T P_{ heta}(X_{t-1}|X_t)$ 노이즈를 제거하여 원본을 복원하는 과정

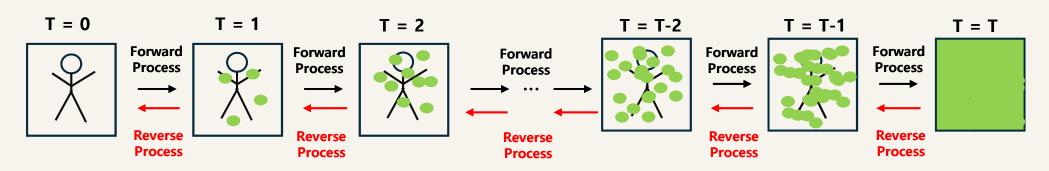
Goal : 확률 분포 q에서 관측한 값으로 $P_{ heta}$ 의 likelihood를 구했을때 그 likelihood가 최대가 되도록 하는 heta를 찾는것



Diffusion Process



Diffusion Process



- Forward process : $q(X_{1:T}|X_o)\coloneqq egin{pmatrix} \prod_{t=1}^T q(X_t|X_{t-1}) \\$ 추가하는 과정
 - → 노이즈를 추가하는 과정(Forward Process)은 모델 X 따라서 <mark>학습을 진행하지 않는다.</mark>
- Reverse process : $P_{\theta}(X_{0:T}) \coloneqq P(X_T) \prod_{t=1}^T P_{\theta}(X_{t-1}|X_t)$ 노이즈를 제거하여 원본을 복원하는 과정
 - → 노이즈를 제거하는 과정(Reverse Process)을 Neurl Net으로 모델링하여 학습을 진행

Forward Process

Forward Process

- 원본 이미지 X_o 에 아주 조금씩 여러단계(T) 에 걸쳐 가우시안 노이즈를 추가하여 <mark>완전한 가우시안 노이즈 (X_T) 로 만드는 과정</mark>
- 해당 과정은 Markov Chain에 정의되며, 각 단계의 상태 X_T 는 오직 X_{T-1} 단계에만 의존한다
- $q(X_{1:T}|X_0) = \prod_{t=1}^T q(X_t|X_{t-1})$
- 네트워크 정의 X -> 훈련 파라미터가 존재하지 않는다.

수식 표현

$$q(X_t|X_{t-1}) \coloneqq N(x_t; \sqrt{1 - B_t}X_{t-1}, B_t I)$$

- $N(x_t; \sqrt{1-B_t}X_{t-1}, B_tI)$: 평균이 $\sqrt{1-B_t}X_{t-1}$, 분산이 B_tI 인 정규분포
- B_t : 분산 스케줄 (t 단계에서 얼마나 많은 노이즈를 추가할지 결정하는 파라미터, 보통 t가 클수록 B_t 도 커지게 설정)
- $B_t I$: 추가되는 노이즈의 분산

$$X_t = \sqrt{1 - B_t} X_{t-1} + \sqrt{B_t} z_{t-1}$$

- 이전 시점 (t-1) 에서 현재 시점(t) 이미지(노이즈를 섞은) 를 만드는 과정
- $z_{t-1}: N(0,I)$ 에서 샘플링 된 가우시안 노이즈

Forward Process

Forward Process

$$X_t = \sqrt{1 - B_t} X_{t-1} + \sqrt{B_t} z_{t-1}$$

$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} z_{t-1} \quad (\alpha_t := 1 - B_t)$$

- 이전 시점 (t-1) 에서 현재 시점(t) 이미지(노이즈를 섞은) 를 만드는 과정
- z_{t-1} : N(0,I)에서 샘플링 된 노이즈
- 이렇게 진행할 시 0부터 T 까지 모든 스텝을 진행해야하므로 샘플링 시간이 매우 오래걸림
- 따라서 DDPM은 중간 단계를 거치지 않고 원본이미지 X_0 로부터 임의의 시점 t의 X_t 를 한 번에 샘플링 할 수 있게 함(학습의 효율)

Forward Process

증명

1.
$$\alpha_t \coloneqq 1 - B_t \ (0 < B_t < 1)$$
, $\overline{\alpha_t} \coloneqq \prod_{s=1}^t \alpha_s$

2.
$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} Z_{t-1}$$

3.
$$X_1 = \sqrt{\alpha_1} X_0 + \sqrt{1 - \alpha_1} z_0$$

4.
$$X_2 = \sqrt{\alpha_2} X_1 + \sqrt{1 - \alpha_2} z_1$$

$$5. \quad X_1$$
를 X_2 에 대입 : $X_2 = \sqrt{\alpha_2 \alpha_1} X_0 + \sqrt{\alpha_2 (1 - \alpha_0)} z_0 + \sqrt{1 - \alpha_2} z_1$

각 가우시안 분포 z_0, z_1 에서 샘플링된 노이즈

$$6. \quad X_2 = \sqrt{\alpha_2 \alpha_1} X_0 + \sqrt{1 - \bar{\alpha}_2} \bar{z}_1$$

7.
$$X_t = \sqrt{\overline{a_t}} X_0 + \sqrt{1 - \overline{a_t}} \overline{z_t}$$

- $ightharpoonup \cdot X_0$ 에서 바로 X_t 를 샘플링 가능 (샘플링 시간 단축)
 - t가 커질수록 점점 원본 이미지 X_0 는 0에 가까워져 노이즈만 남는다

두 가우시안 분포 합치기

$$z_0 \sim N(0, I)$$
 , $z_1 \sim N(0, I)$

 $-> a_0z_0 + b_0z_1$ 은 평균이 0 이고 분산이 $(a^2 + b^2)I$ 인 가우시안 분포를 따른다.

$$a^{2} + b^{2} = \alpha_{2}(1 - \alpha_{0}) + 1 - \alpha_{2}$$

$$= \alpha_{2} - \alpha_{2}\alpha_{1} + 1 - \alpha_{2}$$

$$= 1 - \alpha_{2}\alpha_{1}$$

$$= 1 - \overline{\alpha_{2}}$$





실제 모델내부 동작과정

$$X_t = \sqrt{\overline{a_t}} X_0 + \sqrt{1 - \overline{a_t}} \overline{z_t}$$

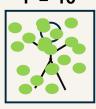




$$T = 3$$



$$T = 10$$



$$T = 20$$



$$T = t-1$$



$$T = t-2$$



$$T = T$$

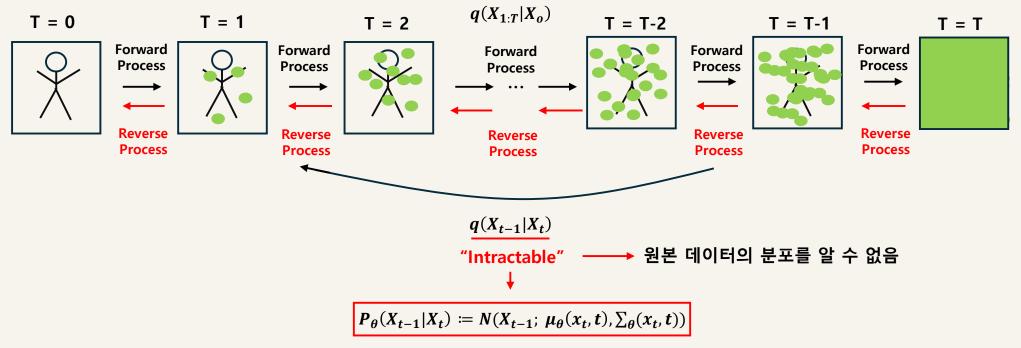


- ightharpoonup 1. 1~T 까지의 다양한 t 값을 랜덤하게 뽑아 X_t 를 만든다
 - 2. 만들어진 X_t 로 신경망은 노이즈를 제거하는 패턴을 학습한다 (Reverse process)



Reverse Process

• 노이즈 (X_t) 를 이미지 (X_0) 로 만드는 과정 (Denosing)



Neural Net으로 모델링하여 디노이징 과정을 학습한다

Goal : 확률 분포 q에서 관측한 값으로 $P_{ heta}$ 의 likelihood를 구했을때 그 likelihood가 최대가 되도록 하는 heta를 찾는것



DDPM Loss

$$\begin{split} E_{q(X_{t}|X_{0})}[-\log(p_{\theta}(X_{0})] &= \int (-\log p_{\theta}(X_{0})) * q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{p_{\theta}(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{p_{\theta}(X_{t}|X_{0})} * \frac{q(X_{t}|X_{0})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{p_{\theta}(X_{t}|X_{0})} * \frac{q(X_{t}|X_{0})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{q(X_{t}|X_{0})} * q(X_{t}|X_{0}) + \log \frac{q(X_{t}|X_{0})}{p_{\theta}(X_{t}|X_{0})}) * q(X_{t}|X_{0})]dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{q(X_{t}|X_{0})} * q(X_{t}|X_{0})dX_{T} - D_{KL}[q||p] \\ &\leq \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0}))dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0},X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0}))dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0}))dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0})dX_{T} \\ &= \int (-\log \frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{$$

Reverse Process

- 노이즈 (X_t) 를 이미지 (X_0) 로 만드는 과정 (Denosing)
- Gaussian Markov chain의 모수인 $\mu_{ heta}$ 를 학습하고자 함

$$E = E_{q(X_t|X_0)}[-log p_{\theta}(X_t)] + E_{q(X_t|X_0)}[-log \frac{p_{\theta}(X_0|X_t)}{q(X_t|X_0)}]$$

DDPM Loss

$$\begin{split} E_{q(X_{t}|X_{0})}[-\log(p_{\theta}(X_{0})] & \leq \int (-\log\frac{p_{\theta}(X_{0},X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0}))dX_{T} \\ & = \int (-\log\frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})} * p_{\theta}(X_{t})) * q(X_{t}|X_{0}))dX_{T} \\ & = \int (-\log\frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}) * q(X_{t}|X_{0})dX_{T} + \int (-\log p_{\theta}(X_{t})) q(X_{t}|X_{0}))dX_{T} \\ & = E_{q(X_{t}|X_{0})}[-\log p_{\theta}(X_{t})] + E_{q(X_{t}|X_{0})}[-\log\frac{p_{\theta}(X_{0}|X_{t})}{q(X_{t}|X_{0})}] \end{split}$$

이대로 학습을 진행하면 P_{θ} 는 q에 따라서 노이즈를 추가하는 방향으로 학습이 진행된다.

따라서 $q(X_t|X_0)$ 를 $q(X_0|X_t)$ 로 바꾸어 학습을 진행해야함

DDPM Loss

$$\begin{split} E_{q(X_{1:T}|X_0)}[-\log(p_{\theta}(X_0)] &= E_{q(X_{1:T}|X_0)}[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{p_{\theta}(X_1,X_2,\dots,X_T|X_0)}] \\ &= E_{q(X_{1:T}|X_0)}[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{p_{\theta}(X_1,X_2,\dots,X_T|X_0)} * \frac{q(X_{1:T}|X_0)}{q(X_{1:T}|X_0)}] \\ &= E_{q(X_{1:T}|X_0)}[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{q(X_{1:T}|X_0)} * \frac{q(X_{1:T}|X_0)}{p_{\theta}(X_1,X_2,\dots,X_T|X_0)}] \\ &= E_{q(X_{1:T}|X_0)}\left[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{q(X_{1:T}|X_0)}\right] - D_{KL}((q_{X_1:T}|X_0)||p_{\theta}(X_1,X_2,\dots,X_t|X_0)) \\ &\leq E_{q(X_{1:T}|X_0)}\left[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{q(X_{1:T}|X_0)}\right] \longrightarrow \text{Markov chain & Baye's} \\ &= E_{q(X_{1:T}|X_0)}\left[-\logp_{\theta}(X_T) - \sum_{t=1}^T \log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_t|X_{t-1})}\right] \\ &= E_{q(X_{1:T}|X_0)}\left[-\logp_{\theta}(X_T) - \sum_{t=1}^T \log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_t|X_{t-1})}\right] \longrightarrow \mathsf{t} = \mathsf{1} \text{ $\ensuremath{\mbox{$\mb$$

DDPM Loss

$$\begin{split} E_{q(X_T|X_0)}[-\log(p_{\theta}(X_0)] & \leq & E_{q(X_T|X_0)}\left[-\log\frac{p_{\theta}(X_0,X_1,\dots,X_T)}{q(X_{1:T}|X_0)}\right] \\ & = & E_{q(X_{1:T}|X_0)}\left[-\log\frac{p_{\theta}(\frac{|X_T|\prod_{t=1}^Tp_{\theta}(X_{t-1}|X_t)}{q(X_{1:T}|X_0)}}\right] \longrightarrow \text{Markov chain} \\ & = & E_{q(X_{1:T}|X_0)}\left[-\log p_{\theta}(X_T) - \sum_{t=1}^T\log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_t|X_{t-1})}\right] \\ & = & E_{q(X_{1:T}|X_0)}\left[-\log p_{\theta}(X_T) - \sum_{t=2}^T\log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_t|X_{t-1})} - \log\frac{p_{\theta}(X_0|X_1)}{q(X_1|X_0)}\right] \longrightarrow \text{t} = 1 \text{ 항을 분리하기} \\ & \vdots \\ & = & E_{q(X_{1:T}|X_0)}[-\log\frac{p_{\theta}(X_T)}{q(X_T|X_0)} - \sum_{t=2}^T\log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_{t-1}|X_t,X_0)} - \log p_{\theta}(X_0|X_1)] \end{split}$$



DDPM Loss

$$E_{q(X_{1:T}|X_{0})}[-log\frac{p_{\theta}(X_{T})}{q(X_{T}|X_{0})} - \sum_{t=2}^{T}log\frac{p_{\theta}(X_{t-1}|X_{t})}{q(X_{t-1}|X_{t},X_{0})} - logp_{\theta}(X_{0}|X_{1})]$$

$$\boxed{1}$$

$$\boxed{2}$$

- \bigcirc $-lograc{p_{ heta}(X_T)}{q(X_T|X_0)}$: 학습을 진행하는 파라미터 heta에 대해서 무관하기때문에 무시
- ② $\sum_{t=2}^{T} log \frac{p_{\theta}(X_{t-1}|X_t)}{q(X_{t-1}|X_t,X_0)}$: 주된 학습 목표 (reverse process)
- (3) $logp_{\theta}(X_0|X_1)$: 수많은 T(Time-step)에서 극히 일부분이므로 값이 매우 작음(학습과정에서 무시)

결국 ② 를 minimize 하는 문제

DDPM Loss

$$(2)E_{q(X_{1:T}|X_0)}\sum_{t=2}^{T}-log\frac{p_{\theta}(X_{t-1}|X_t)}{q(X_{t-1}|X_t,X_0)}$$
 : 주된 학습 목표 (reverse process)

$$\int p(x)log p(x)dx = -\frac{1}{2}(1 + \log(2\pi\sigma_1^2))$$

$$\int p(x)log q(x)dx = -\frac{1}{2}\log(2\pi\sigma_2^2) - (\frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2})$$

$$\longrightarrow$$
 분포가 가우시안(p,q)

$$D_{KL}(p||q) = \int p(x)logp(x)dx - \int p(x)logq(x)dx$$

$$= -\frac{1}{2} + \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} \longrightarrow \text{DDPM에서는 분산이 고정값이므로 } \frac{\sigma_2}{\sigma_1} = \mathbf{1}$$

$$= E_q\left[\frac{1}{2\sigma_t^2}\big||\mu_1 - \mu_2|\big|^2\right] + C$$

DDPM Loss

$$2E_{q(X_{1:T}|X_0)} \sum_{t=2}^{T} -log \frac{p_{\theta}(X_{t-1}|X_t)}{q(X_{t-1}|X_t,X_0)} : 주된 학습 목표 (reverse process) D_{KL}(p||q) = E_q \left[\frac{1}{2\sigma_t^2} \left| |\mu_1 - \mu_2| \right|^2 \right] + C$$

$$= E_q \left[\frac{1}{2\sigma_t^2} \left| \frac{|\widetilde{\mu_t}(X_t, X_0) - \underline{\mu_\theta(X_t, t)}|}{\downarrow}^2 \right] + C$$

q,p 분포의 평균값 $(\widetilde{\mu_t},\mu_{ heta})$ 만 구하면 로스 계산이 가능

DDPM Loss

$$E_{q} \left[\frac{1}{2\sigma_{t}^{2}} \left| \left| \frac{\widetilde{\mu}_{t}(X_{t}, X_{0})}{\downarrow} - \underline{\mu}_{\underline{\theta}}(X_{t}, t) \right| \right|^{2} \right] + C$$

$$q \qquad \qquad p$$

- → Goal : q, p 의 평균값을 구하기
- 1. $q(X_t|X_0)$

$$q(X_t|X_{t-1})=N(X_t;A=\sqrt{1-eta_t}X_{t-1},eta_t I)$$
 — Reparameterization trick 사용하여 표현 $lpha_t\coloneqq 1-eta_t \ and \ ar{a}_t\coloneqq \prod_{s=1}^t lpha_s$

 X_{t-1} 을 같은 방식으로 구하여 대입 후 정리

$$\longrightarrow \quad X_t = \sqrt{\overline{\alpha}_t} X_0 + \sqrt{1 - \overline{\alpha}_t} \epsilon$$

ightarrow 즉 $\operatorname{q}(X_t|X_0)$ 는 평균이 $\sqrt{\overline{\alpha}_t}$, 표준편차가 $\sqrt{1-\overline{\alpha}_t}$ 인 분포

$$q(X_t|X_{t-1}) = N(X_t; A = \sqrt{1 - \beta_t} X_{t-1}, \beta_t I)$$

$$X \sim N(\mu, \sigma^2) \to x = \mu + \sigma \epsilon,$$

$$\epsilon \sim N(0, I)$$

1.
$$\alpha_t := 1 - B_t \ (0 < B_t < 1) \ , \ \overline{\alpha_t} := \prod_{s=1}^t \alpha_s$$

2.
$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} Z_{t-1}$$

3.
$$X_1 = \sqrt{\alpha_1} X_0 + \sqrt{1 - \alpha_1} z_0$$

$$4. X_2 = \sqrt{\alpha_2} X_1 + \sqrt{1 - \alpha_2} z_1$$

$$S$$
. X_1 를 X_2 에 대입 : $X_2 = \sqrt{\alpha_2 \alpha_1} X_0 + \sqrt{\alpha_2 (1 - \alpha_0)} z_0 + \sqrt{1 - \alpha_2} z_1$

$$6. X_2 = \sqrt{\alpha_2 \alpha_1} X_0 + \sqrt{1 - \overline{\alpha}_2} \overline{z_1}$$

7.
$$X_t = \sqrt{\overline{a_t}} X_0 + \sqrt{1 - \overline{a_t}} \overline{z_t}$$

DDPM Loss

$$(2) E_{q(X_{1:T}|X_{0})} \sum_{t=2}^{T} -log \frac{p_{\theta}(X_{t-1}|X_{t})}{q(X_{t-1}|X_{t},X_{0})} = E_{q} \left[\frac{1}{2\sigma_{t}^{2}} \left| \left| \widetilde{\mu_{t}}(X_{t},X_{0}) - \mu_{\theta}(X_{t},t) \right| \right|^{2} \right] + C$$

1.
$$q(X_{t-1}|X_t,X_0)$$

$$q(X_{t-1}|X_t, X_0) = q(X_t|X_{t-1}, X_0) \frac{q(X_{t-1}|X_0)}{q(X_t|X_0)}$$

$$\rightarrow$$
 $q(X_t|X_0) = N(X_t; \sqrt{\overline{\alpha}_t}X_0, (1-\overline{\alpha}_t)I)$ 대입

$$\vdots f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
에 대입 후 계산

$$\longrightarrow \quad \tilde{u}_t(x_t, x_0) \coloneqq \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} X_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} X_t , \tilde{\beta}_t \coloneqq \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

 $: X_{\tau}$ 에 대해서 정리

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(X_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \overline{\alpha_t}}} \epsilon_t \right)$$

- 2. $p_{\theta}(X_{t-1}|X_t)$
- → 근본적으로 q에 근사를 해야한다. 다만, θ (학습 가능한 파라미터) 를 도입하여 학습 가능하도록 설 정한다

$$\mu_{\theta}(X_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(X_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \overline{\alpha_t}}} \epsilon_{\theta}(X_t, t) \right)$$

DDPM Loss

$$\widetilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(X_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \overline{\alpha_t}}} \epsilon_t \right) \qquad \mu_{\theta}(X_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(X_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \overline{\alpha_t}}} \epsilon_{\theta}(X_t, t) \right)$$

: 위 두 식을 대입하여 정리

최종
$$Diffusion\ Loss = E_{t,X_0,\epsilon}[\left|\left|\epsilon - \epsilon_{\theta}\left(\sqrt{\overline{\alpha}_t}X_0 + \sqrt{1-\overline{\alpha}_t}\epsilon,t\right)\right|\right|^2$$

Neural Network는 원본 노이즈 ϵ 과 시점 t에서 모델이 예측한 노이즈의 L2 Norm을 최소화 하는 방향으로 학습하게 된다(노이즈 패턴의 학습)





평가지표 : FID score, IS score

Table 1: CIFAR10 results. NLL measured in bits/dim.

Model	IS	FID	NLL Test (Train)
Conditional			
EBM [11]	8.30	37.9	
JEM [17]	8.76	38.4	
BigGAN [3]	9.22	14.73	
StyleGAN2 + ADA (v1) [29]	10.06	2.67	
Unconditional			
Diffusion (original) [53]			≤ 5.40
Gated PixelCNN [59]	4.60	65.93	$3.\overline{03} (2.90)$
Sparse Transformer [7]			2.80
PixelIQN [43]	5.29	49.46	
EBM [11]	6.78	38.2	
NCSNv2 [56]		31.75	
NCSN [55]	8.87 ± 0.12	25.32	
SNGAN [39]	8.22 ± 0.05	21.7	
SNGAN-DDLS [4]	9.09 ± 0.10	15.42	
StyleGAN2 + ADA (v1) [29]	9.74 ± 0.05	3.26	
Ours $(L, \text{ fixed isotropic } \Sigma)$	7.67 ± 0.13	13.51	$\leq 3.70 (3.69)$
Ours (L_{simple})	9.46 ± 0.11	3.17	$\leq 3.75 (3.72)$

FID Score : 실제 이미지와 생성된 이미지를 Pre-trained 된 InceptionNet 을 이용하여 피처맵을 추출한 뒤 추출된 두개의 피처맵의 분포 차이로 계산

IS Score: 실제 이미지와 생성된 이미지를 Pre-trained 된 InceptionNet 을 이용하여 분류한 뒤 분류 확률값을 이용하여 KL-divergence 를 계산

NLL Test: 모델이 데이터의 분포를 얼마나 <mark>효율적으로 학습했는지</mark>를 보여주는 지표, 한 픽셀을 표현할때 얼마나 많은 비트가 사용되는지



생성 샘플





LSUN(Church), FID: 7.89

LSUN(Bedroom), FID: 4.90

서로 다른 두 이미지 보간





Summary

- Diffusion은 Forward Process와 Reverse Process로 나뉜다.
- Forward Process는 원본 이미지에 Gaussian noise를 주입하여 완전한 Gaussian noise로 만든다
- Reverse Process는 Denoise 방향으로 학습이 진행된다.



5 Code

Expreiment

Simple diffusion process code

