

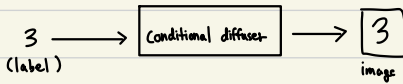
## Day 6, 7: Conditional Diffusion Model

- Since now, we have modeled generative model that models population of  $x$ ,  $p(x)$ . What if we model  $P(x|y)$  where  $y$  is a condition.  
⇒ 조건  $y$  (텍스트, 라벨, 이미지 등) 를 비워두어 생성하고자 하는  $x$ 를 지정한다!

### \* Set condition $y$ as label. (Conditional Diffusion Model)

- Train with MNIST images and give condition (label) when generating!

fig)



- Conditional Diffusion Model where the neural network predicts  $u_\theta(x_t, t)$  directly.

## 1. Modeling $P_\theta(x)$

$$p_\theta(x_0) = \int p_\theta(x_0, x_1, x_2 \dots x_T) dx_{1:T} \text{ (marginalization)}$$

latent variables

이것도 latent variables로  $P_\theta(x_T)$ 이나,  $P(x_T) = N(x_T; 0, I)$ 로 가정. (prior를 가정한다! 현재  $x_T$ 가 어떤 분포를 따를지 가정)

$$= \int p_\theta(x_0 | x_{1:T}) p_\theta(x_1 | x_{2:T}) \dots p_\theta(x_{T-1} | x_T) P(x_T) dx_{1:T} \text{ (chain rule)}$$

$$= \int p_\theta(x_0 | x_1) p_\theta(x_1 | x_2) \dots p_\theta(x_{T-1} | x_T) P(x_T) dx_{1:T} \text{ (Markov property) where } P(x_T) = N(x_T; 0, I)$$

$$p_\theta(x_{t-1} | x_t) = N(x_{t-1}; u_\theta(x_t, t), S_f^2(t)I)$$

## 2. Modeling $P_\theta(x_0 | y)$

$$p_\theta(x_0 | y) = \int p_\theta(x_0, x_1, \dots x_T | y) dx_{1:T} \text{ (marginalize)}$$

added condition

latent variables

$$= \int p_\theta(x_0 | x_{1:T}, y) p_\theta(x_1 | x_{2:T}, y) \dots p_\theta(x_{T-1} | x_T, y) p_\theta(x_T | y) dx_{1:T} \text{ (chain rule)}$$



Markov (t-1 point t-2, 전 정보 + 현재 point만) :  $p_\theta(x_{t-1} | x_{t:T}, y) \approx p_\theta(x_{t-1} | x_t, y)$

$$= \int p_\theta(x_0 | x_1, y) p_\theta(x_1 | x_2, y) \dots p_\theta(x_{T-1} | x_T, y) p_\theta(x_T | y) dx_{1:T} \text{ (Markov property)}$$

일반적으로 latent variable 이용.. But, diffusion model assumes  $p_\theta(x_T) \rightarrow P(x_T) = N(x_T; 0, I)$

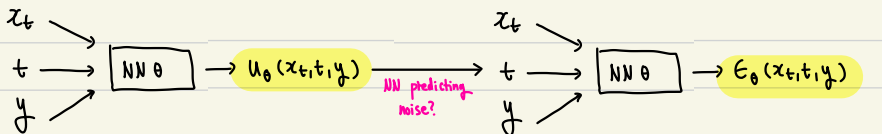
And also in conditional diffusion model  $p_\theta(x_T | y) \rightarrow P(x_T | y) = P(x_T) = N(x_T; 0, I)$



또한 현재 이미 동일한 Gaussian noise를 가정. (prior를 가정한다)

$$p_\theta(x_{t-1} | x_t, y) = N(x_{t-1}; u_\theta(x_t, t, y), S_f^2(t)I) \text{ (added condition } y \text{ for Neural Network's input)}$$

fig)



=> 순전파는 늘 cca 이렇게 condition은 주는 행위로 노이즈 예측 혹은 내분출을 직접 예측하는 방식으로 상당히 간편함.

=> forward diffusion process에 label 제공 X

## \* Score function

- At previous step, just provided condition  $y$ ...

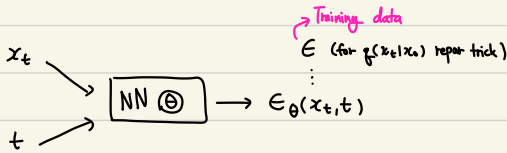
→ issue: The condition 'y' can be ignored.

Is there any way we can strengthen the condition?

"Guidance Method"

In diffusion model, the neural network predicted noise at each step  $\Rightarrow \epsilon_\theta(x_t, t)$

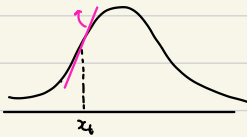
fig)



Where  $\epsilon \approx -\sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log P(x_t) \dots$  ① (중요한  $\epsilon$ 가 score function  $\nabla_{x_t} \log P(x_t)$ 와 관련 됨.)

①: A gradient of  $x_t$  related to  $\log P(x_t) \Rightarrow$  score function / score

비슷하면 약간 다른

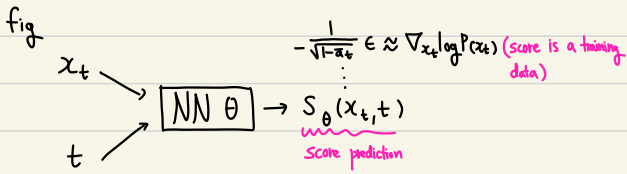


①에서  $\epsilon$ 는  $\nabla_{x_t} \log P(x_t)$  (score)로 근사할 수 있음.

$$\epsilon \approx -\sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log P(x_t)$$

$\epsilon$ 가 score와 관련 있음 때문에 붙여 제 이름.

$\Rightarrow$  그럼  $\epsilon$  말고, score'!  $-\frac{\epsilon}{\sqrt{1-\bar{\alpha}_t}} = \nabla_{x_t} \log P(x_t)$ 를 학습 데이터 쓰자! (신경망이 score를 예측하도록 한다)



$\Rightarrow$  "Diffusion model where neural network predicts score  $S_\theta(x_t, t) \rightarrow \nabla_{x_t} \log P(x_t) \approx -\frac{\epsilon}{\sqrt{1-\bar{\alpha}_t}}$ "

\* Why ① in page 3?

$$\epsilon \approx -\sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log P(x_t) \dots \textcircled{1}$$

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I)$$

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon \quad (\epsilon \sim \mathcal{N}(0, I)) \dots \textcircled{2}$$

• Tweedie's Formula

For  $x \sim \mathcal{N}(x; u, \Sigma)$ ,  $\nabla \log P(x)$  is

$$E[u|x] = x + \Sigma \nabla_x \log P(x) \quad (\nabla \log P(x) \text{ is the score function})$$

• Adopt Tweedie's formula to Diffusion

$$x_t \sim \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I) \quad \textcircled{a}$$

$$E[\sqrt{\bar{\alpha}_t} x_0 | x_t] = x_t + (1-\bar{\alpha}_t) I \nabla_{x_t} \log P(x_t) \quad \textcircled{b}$$

Since  $\textcircled{a}$ ,  $x_t$  is constant

$$E[\sqrt{\bar{\alpha}_t} x_0 | x_t] = x_t + (1-\bar{\alpha}_t) I \nabla_{x_t} \log P(x_t)$$

$$\Leftrightarrow E[x_t - \sqrt{1-\bar{\alpha}_t} \epsilon | x_t] = x_t + (1-\bar{\alpha}_t) I \nabla_{x_t} \log P(x_t)$$

$$\Leftrightarrow E[x_t | x_t] - E[\sqrt{1-\bar{\alpha}_t} \epsilon | x_t] = x_t + (1-\bar{\alpha}_t) I \nabla_{x_t} \log P(x_t)$$

$$\Leftrightarrow x_t - \sqrt{1-\bar{\alpha}_t} E[\epsilon | x_t] = x_t + (1-\bar{\alpha}_t) I \nabla_{x_t} \log P(x_t)$$

$$\Leftrightarrow E[\epsilon | x_t] = -\sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log P(x_t)$$

$$\Leftrightarrow \text{monte carlo for } E[\epsilon | x_t] : \epsilon \approx -\sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log P(x_t)$$

(one sample  $x_t$ )

## \* Neural Network predicting score function 부설명

- Tweedie 공식으로 현재 얻어있는  $\epsilon$ 를  $-\sqrt{1-a_t} \nabla_{x_t} (\log P(x_t))$ 에 근사하여  $\frac{\epsilon}{-\sqrt{1-a_t}}$  (Score function)를 주어진 데이터로 사용  $\rightarrow$  신경망이 score를 예측할 수 있게끔.
- 이 모든 과정은  $\nabla_{x_t} \log P(x_t)$ 를 해석적으로 풀 수 없게끔 얻어오기 때문에 주어진 score를 근사할 때에 대한 방법임.  $\rightarrow$  Tweedie를 통해 score를 근사할 방법을 찾음.

구체)  $x_t$ 는 한 데이터  $x_0$ 에 샘플  $t$ 에 해당하는 noise가 섞인 데이터 ( $x_t$  ->  $x_0$ 를 sampling함) ①

Tweedie 공식으로 주어진  $x_t$ 와  $x_0$ 에 대한 score function은 기본적으로 ②  $a_t$  hyperparam

스케일된 원본  $\sqrt{a_t} x_0$ 에 대한 기댓값을 풀 수 있다 (사실 평균)



트weedie 공식이 가능하게끔  $\sqrt{a_t} x_0$ 를 노이즈  $\epsilon$ 에 대한 식으로 바꿀 수 있고,  $\epsilon$ 의  $(E[\sqrt{a_t} x_0 | x_t] \rightarrow E[\epsilon | x_t])$ 에 대한 기댓값으로 바꿀 수 있다:  $E[\epsilon | x_t]$



$$E[\epsilon | x_t] = -\sqrt{1-a_t} \nabla_{x_t} \log P(x_t) \text{ 라는 식에 의존한다}$$

③

monte carlo를 통해  $E[\epsilon | x_t]$ 를 근사하는 것  
③을 근사하는 것과 같음.

$\Rightarrow$  몬테카를로를 통해 주어진  $x_t$ 에 대한  $\epsilon$ 의 랜덤 샘플에서  $\epsilon$ 를 여러개 얻은 1개 평균이기 ③라는 Tweedie formula로부터 이런식으로 도출된 식을 근사.

$P(\epsilon | x_t)$ 에서  $\epsilon$ 를 랜덤한 sampling하면,  $\epsilon \approx -\sqrt{1-a_t} \nabla_{x_t} \log P(x_t)$ 가 나온다.

이들 통해  $\frac{\epsilon}{-\sqrt{1-a_t}}$ 를 보낼 수 있게끔 할 수 있게끔 Score function의 예측을 할 수 있게됨. (이론적 제1차 forward diffusion process에 해당하는  $\epsilon$ 를 score function으로 활용함!)

## \* 조건부 기댓값

\*\*조건부 기댓값(conditional expectation)\*\*의 정의는 다음과 같습니다.

- 조건부 기댓값은 어떤 확률 변수  $X$ 와  $Y$ 에 대해,  $X = x$ 라는 정보가 주어졌을 때  $Y$ 의 평균(기댓값)을 의미합니다. 즉, " $X = x$ 라는 상황 아래에서  $Y$ 가 가질 수 있는 값들의 평균"입니다. ① ② ③ ④.
- 조건부 확률 분포를 사용하기 때문에, 조건이 주어진 상태에서의 기댓값이라고도 하고, 조건에 따라 값이 바뀌는 함수입니다.

이산형(Discrete) 확률변수의 경우

$$E[Y|X=x] = \sum_y y \cdot P(Y=y | X=x)$$

연속형(Continuous) 확률변수의 경우

$$E[Y|X=x] = \int y \cdot f_{Y|X}(y|x) dy$$

여기서  $f_{Y|X}(y|x)$ 는 조건부 확률밀도함수입니다. ⑤ ⑥.

즉, 조건부 기댓값은 어떤 사건이나 랜덤변수의 값이 주어진 상황에서, 다른 변수의 평균적인 값을 의미합니다. 이는 데이터 예측, 통계 모델링, 머신러닝 등 거의 모든 확률적 모델에서 중요한 개념입니다.

# \* Guidance (with aspect of neural net predicting score in diffusion)

- More details of Classifier Guidance & Classifier Free guidance will be dealt in Github README.md files!

① Classifier Guidance : Noise를 어떻게 하면 손쉽게 이를 GAN에 생성 가능!

- Score function predicting diffusion model :  $\nabla_{x_t} \log P(x_t)$
- Conditional Score function predicting diffusion model :  $\nabla_{x_t} \log P(x_t|y)$

· Induction of conditional score function

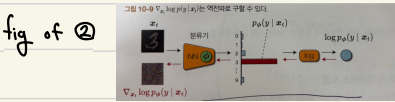
$$P(x_t|y) = \frac{P(x_t)P(y|x_t)}{P(y)}$$

$$\begin{aligned} \nabla_{x_t} \log P(x_t|y) &= \nabla_{x_t} \log \frac{P(x_t)P(y|x_t)}{P(y)} \\ &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(y|x_t) - \cancel{\nabla_{x_t} \log P(y)}^0 \\ &= \underbrace{\nabla_{x_t} \log P(x_t)}_{\substack{\text{score prediction} \\ \text{neural network} \\ \textcircled{1}}} + \underbrace{\nabla_{x_t} \log P(y|x_t)}_{\substack{\text{Classifier} \\ \text{gradient} \\ \textcircled{2}}} \end{aligned}$$

로 각각 관련이 있는 것을 조합할 수 있음!

① : Diffusion neural network를 통해 score 찾기

② : Classifier를 y에 대한 log 값을 계산 후 입력  $x_t$ 에 대한 gradient 찾기



→ 이 둘을 조합하기 위해 분기 기호를 조합 . (생성 n class y 값으로 만들어 줌)

$$\Rightarrow \nabla_{x_t} \log P(x_t|y) = \nabla_{x_t} \log P(x_t) + \rightarrow \nabla_{x_t} \log P(y|x_t)$$

## \* Summary of classifier guidance of diffusion model

- 일반적인 점 추정 혹은 모델에 분류를 강하게 하면 전체 점들 또한 전체 픽셀 모델은 많은 에러를 생!

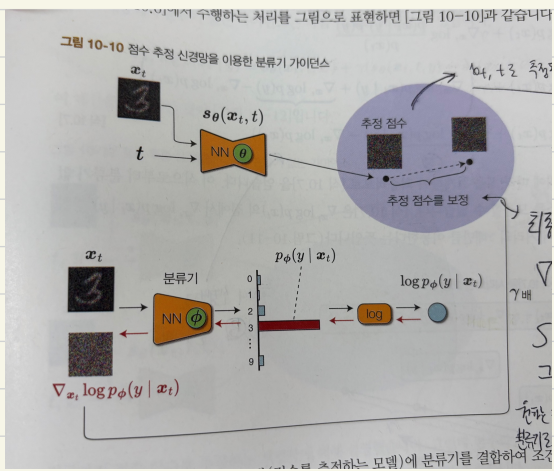
①

②

- 전체 픽셀 score 예측은 학습된 classifier에서 생는 gradient를 통해 생

- 1- 생는 시에 원본의 강도로 y label 방향으로 생는 지에 따른 강도로 나뉨.

- 생는 법 fig



## ② Classifier free Guidance

- ①이 불완전 cuz 조건이 없어야 해서 실패함 ↓

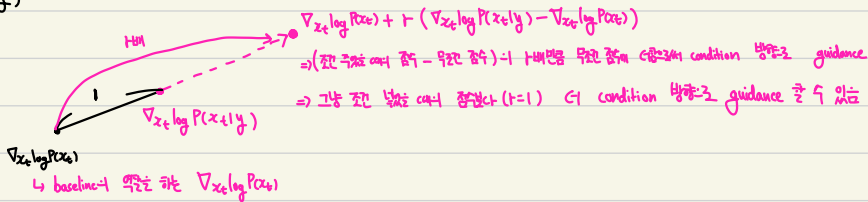
- Give guidance without classifier

$$\begin{aligned}
 \nabla_{x_t} \log P(x_t | y) &= \nabla_{x_t} \log \frac{P(x_t) P(y | x_t)}{P(y)} \\
 &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(y | x_t) - \cancel{\nabla_{x_t} \log P(y)}^0 \\
 &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(y | x_t) \quad \text{C.G. induction} \\
 &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log \frac{P(y) P(x_t | y)}{P(x_t)} \\
 &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(y) + \nabla_{x_t} \log P(x_t | y) - \nabla_{x_t} \log P(x_t) \\
 &= \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(x_t | y) - \nabla_{x_t} \log P(x_t)
 \end{aligned}$$

bypass 할 수 있는 거 있음!

$$\therefore \nabla_{x_t} \log P(x_t | y) = \nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(x_t | y) - \nabla_{x_t} \log P(x_t)$$

fig)



$\Rightarrow$  "HMM" Diffusion: neural network로 조건점수와 무조건 점수를 모두 출력하게끔 훈련시키고, generation시  $\nabla_{x_t} \log P(x_t) + \nabla_{x_t} \log P(x_t | y) - \nabla_{x_t} \log P(x_t)$ 를 통해 condition guiding을 줌 (HMM은 condition 빼면 guidance 할 수 있음)



• How to predict  $\nabla_{x_t} \log P(x_t), \nabla_{x_t} \log P(x_t|y)$

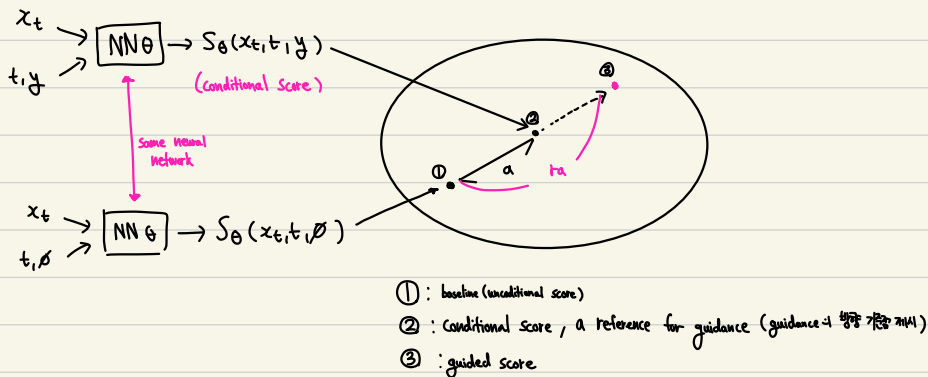
Predict with Single Neural Network! ( $S_\theta$ )

- ①  $\nabla_{x_t} \log P(x_t) = S_\theta(x_t, t, \emptyset)$  ↗ zero embedded  
 ②  $\nabla_{x_t} \log P(x_t|y) = S_\theta(x_t, t, y)$  ↗ condition embedded

⇒ 같은 시에 ①, ②를 다 예측할 수 있도록 한층-로 만들어줌

$$\therefore \nabla_{x_t} \log P(x_t|y) = S_\theta(x_t, t, \emptyset) + \gamma (S_\theta(x_t, t, y) - S_\theta(x_t, t, \emptyset))$$

Fig) Generation with Classifier Free Guidance



• What if Neural Network predicts noise  $\epsilon$ ?

Since  $\epsilon \approx -\sqrt{1-\alpha_t} \nabla_{x_t} \log P(x_t)$ , Score와 상반반면에 채워지지 않음.

⇒  $\epsilon_{\text{uncond}}$  가 있으면,  $\epsilon_\theta(x_t, t, \emptyset) / \epsilon_\theta(x_t, t, y)$ 로 unconditional/conditional 점수를 예측할 수 있음!

★ Summary) 훈련 시에 동일한 network가 unconditional / conditional score를 모두 학습하도록 train, generation 시에 unconditional score와 혼합한 guidance를 제공

baseline

direction  
reference

conditional score를 모두 계산하여  $uncon + t(con - uncon)$  방식으로 baseline uncon score를 보정하여 혼합한 condition 방향으로 보정

• Appendix

$$S_{guided} = S_{uncon} + t(S_{con} - S_{uncon})$$



$t=0$  : unconditional generation, ①

$t=1$  : 전역 score만 반영된 conditional generation

↗ 인공, 부분적

$t>1$  : conditional signal을 강화, 더 "컨트롤" 전역 score 생성  $\Rightarrow$  t이 너무 크면 artifact가 생길 위험

"Please also refer to README.md"