

Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
 - A. 이름: **Maximum Likelihood Training of Score-Based Diffusion Models**
 - B. 저널: **NeurIPS**
 - C. 도메인: **Diffusion**
 - D. 출판연도: **2021**
 - E. 저자: **Yang Song, Stefano Ermon**
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. 논문 **Figure**를 그대로 따라 그리면 안됩니다.
 - A. [선행연구] 2021 Yang Song 논문에서는, DDPM과 SMLD를 Score Based Diffusion Model 형태로 결합한 프레임워크를 제시했음.
 - B. [선행연구] 그러나, Score Based Model의 본질적인 문제는, Score Matching을 Loss로 사용하여 최적화를 한다는 점인데, 이는 Likelihood에 영향을 주지 못한다는 것임.
 - C. [선행연구] 따라서, Likelihood가 필요한 Task를 해결할 수 없었을 뿐 더러, 정확한 Likelihood 값을 계산하지 못하는 문제로 인해 다른 Generative Model과의 비교가 어렵다는 한계가 있었음.
 - D. 이 논문에서는 Score Based Diffusion Model의 Loss를 크게 벗어나지 않고도, Likelihood를 maximum하게 할 수 있는 Likelihood Weight를 제안함.
 - E. 기존의 Score Matching에서는 Weighted Combination을 사용했는데, 이것의 역할은 모든 noise에 대해서 동일한 Weight를 주는 것을 방지하고(만약, 모든 noise에 대해서 동일한 Weight를 설정하는 경우에는 Variance가 커질 수 있음), bias를 방지하는 것임. Weighted Combination: $\lambda(t)$
 - F. KL-Divergence를 Minimizing 하는 것이 Likelihood를 Maximizing하는 것과 동일하며, $\lambda(t) = \sigma^2$ 일 때, KL-Divergence는 Score Matching에 대한 Upper Bound로 작용한다는 것을 증명함.

- G. 결과적으로(실험 + 증명), Likelihood를 증가시킬 수 있는 구조를 입증함
3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
- A. Score Matching을 Loss로 사용하는 경우 생기는 문제점은, Model과 특정 Distribution의 Score 값의 차이를 줄이는 것을 목표로 한다는 것임. 다시 말하자면, 해당 값을 최적화한다고 하더라도, Likelihood에는 큰 영향을 주지 못한다는 한계가 있다는 것이고, Likelihood를 직접적으로 활용하는 Task를 해결하지 못한다는 것을 뜻함.
- B. 그러나 해당 논문에서는, Likelihood Weight를 접목하여, Upper Bound를 기반으로, 해당 값을 최적화하면, Likelihood 또한 상승하도록 하는 구조를 제안함.
- C. 결과적으로, Score Matching이라는 loss의 큰 틀은 벗어나지 않은 채, Likelihood를 상승시킬 수 있는 방안을 제시함.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
- A. 수식 유도과정 확인
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. Sampling을 하기 위해서는 여러 번 Model을 호출할 필요가 있음
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
- A. Diffusion 선행연구
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
- A. 아니요

날짜: 2025-07-07

이름: 신준원

Maximum Likelihood Training of Score-based Diffusion models 01/01.

Introduction

Score-based diffusion model

- Diffusing process (SDE)

$\Rightarrow P_{\text{data}} \sim \mathcal{N}$ (noise distribution)

- reverse diffusion process (ODE)

\Rightarrow SDE + Score function \Rightarrow CNF (continuous normalizing flow)

tractable likelihood

* Score function

\Rightarrow 연속적인 시간변화 흐름 (random X \Rightarrow ODE $\frac{dx}{dt} = \dots$)

$\nabla \log p(x) \Rightarrow$ Gradient field (= Score) 리딩
(= $\frac{\partial \log p(x)}{\partial x}$)

real data p_0

noise data p_T

$\frac{\partial \log p(x)}{\partial x}$
(score)

높은 = high score

낮은 = low score

low

high

ex) $p(x) = \mathcal{N}(0, \sigma^2 I)$, normal distribution.

$p(x) \sim 1$ 위한 Score 리딩

$$\log p(x) = -\frac{1}{2\sigma^2} \|x\|^2 + C$$

$$\nabla_x \log p(x) = -\frac{x}{\sigma^2}$$

\Rightarrow 리딩 리딩 = score 리딩

리딩 리딩 = score 리딩

리딩 \rightarrow 리딩

⊕ Weight Scale 이 필요 이유

NCSN \rightarrow 모든 noise level 에 대한 분포를 학습하기 위해
noise 이나 다른 loss update 가 필요.

weighting

\Rightarrow σ 이 작을수록 작은 sigma 의 경우 \Rightarrow error 가 작을수록 $\left(\frac{\sigma}{\sigma}\right)$

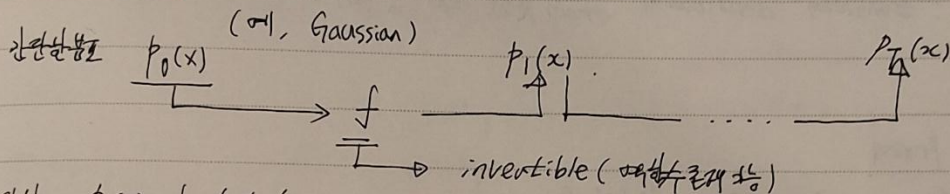
큰 sigma 의 경우 \Rightarrow error 가 작을수록 $\left(\frac{\sigma}{\sigma}\right)$

정확도에 대해 λ 가 클수록 효과가 있음 $\Rightarrow \lambda(\sigma)$

CNF (continuous normalization flow)

: PDF 이 세해서, 연속적인 시간변화로 모델링

Normalizing Flow (이산적인 변수 $\rightarrow t, t+1, \dots$) \rightarrow Continuous ~~st~~



간단한 $p_0(x)$ distribution

\Rightarrow 복잡한 $p_T(x)$ 분포로 변환 (with invertible f)

\rightarrow 복잡한 분포로 만들거나 과정에서 필요한 것은 ODE - solver

즉, 이 step 에 대해서, 어떤 ODE - solver 필요.

계산량 ↑. 그러나, 정확한 계산을 기본으로.

정확한 likelihood 계산에 이용 가능

그렇다면 Score Matching Loss는?

\rightarrow 효과적인 효율적인 학습가능 (그러나 \rightarrow 정확한 likelihood 계산 불가)

* likelihood 를 사용하는 (필요) task 에 적용이 어렵다는 한계 존재

결론 \rightarrow weighted Combination

\Rightarrow likelihood와 무관하게 학습하면 SM-Loss
 \Rightarrow likelihood 값을 추가로 학습가능!

denotation

$p(x)$: dataset distribution ($p = p$)

$p_{\pm}(x)$: marginal distribution of $X(\pm)$

$p_{\pm}(x'|x)$: transition distribution from $X(0)$ to $X(\pm)$

Diffusing process

목적 : $p(x) \rightarrow$ noise distribution.

사용자 : SDE (stochastic differential equation)

$$dx = \underbrace{f(x,t)}_{\text{drift}} dt + \underbrace{g(t)}_{\text{diffusion}} dW \quad \text{--- Wiener}$$

Coefficient

$$p(x) + dx \rightarrow \{X(\pm)\}_{\pm \in [0, T]} \rightarrow \pi(x) \sim \mathcal{N}(0, I).$$

tractable, reverse를 관측하기!

\Rightarrow 따라서, 분포를 완전히 바꾼다.
($p(x)$)

* Diffusion process

$\rightarrow f(x,t), g(t), T$ 를 수정 가능

이 때, 어떻게 수정하느냐에 따라

$\left. \begin{array}{l} VE \\ VP \\ sub-VP \end{array} \right\} 3 \text{가지 유형}$

Sampling with Reverse SDE

reverse time diffusion process

$$dx = [f(x, t) - g(t)^2 \nabla_x \log p_t(x)] dt + g(t) \cdot d\tilde{w}$$

Werner process
(reverse)

diffusion $x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_T$

reverse $x_0 \leftarrow x_1 \leftarrow x_2 \leftarrow \dots \leftarrow x_T$

target : $\nabla_x \log p_t(x)$ (가짜 데이터)

$\hookrightarrow S_\theta(x, t)$: model (Score-based model)

가짜 데이터를 모델의 점수 $\nabla_x \log p_t(x) \approx S_\theta(x, t)$

model \Rightarrow parameterize

likelihood-base : $\nabla \log p_t(x)$

Score-function base : $\nabla \log p_t(x) \rightarrow$ $\nabla \log p_t(x)$ 가 진짜 x \rightarrow $\nabla \log p_t(x)$ 가 진짜 x

실제 Loss

$$J_{gm}(\theta; \lambda(\cdot)) := \frac{1}{T} \int_0^T \mathbb{E}_{p_t(x)} [\lambda(t) \|\nabla_x \log p_t(x) - S_\theta(x, t)\|^2] dt$$

(positive) weighting function

$\arg\min J_{gm}(\theta; \lambda(\cdot)) \rightarrow$ training gm.

$\Rightarrow dx =$ Affine $\Rightarrow p_t(x'|x) = \text{Gaussian}$

\hookrightarrow gm. 학습

\hookrightarrow sampling reverse - gdz $\frac{1}{2}$

(with gdz-golker)

Likelihood of SBDM

Diffusion process
Reverse process \rightarrow likelihood 평가 가능

① $p_{\theta}^{SDE}(x)$

$$\rightarrow d\hat{x} = [f(\hat{x}, t) - g(t) \dot{s}_{\theta}(\hat{x}, t)] dt + g(t) d\bar{w} \quad \hat{x}_{\theta}(T) \sim \pi$$

$$p_{\theta}^{SDE}(x) \Rightarrow \hat{x}_{\theta}(0) \sim p_{\theta}^{SDE} \quad (\text{with GDE-solver})$$

GDE의 경우 정확한 계산이 불가능함
 \Rightarrow Lower-bound 이용

② $p_{\theta}^{ODE}(x) \rightarrow$ (SDE); \rightarrow margin. \rightarrow $p_{\theta}(x)$ 평가 가능.

$$\frac{dx}{dt} = f(x, t) - \frac{1}{\sigma} g(t) \nabla_x \log p_{\theta}(x)$$

특징: deterministic

$$\rightarrow f(x) \rightarrow f^{-1}(x) \Rightarrow \text{정확한 계산 가능} \quad (\text{계산량 } \uparrow, T)$$

approx likelihood 계산 가능

Bounding likelihood

- ① $p_{\theta}^{SDE} \rightarrow$ Likelihood not x
 ② $p_{\theta}^{ODE} \rightarrow$ Likelihood not 0 but expensive.

↓

Score match loss \rightarrow not loss, not likelihood \rightarrow not good.

KL divergence

maximizing Likelihood = minimizing D_{KL} * main idea

정리 1

$$D_{KL}(p \parallel p_{\theta}^{SDE}) \leq J_{SM}(\theta; f(\cdot)^v) + D_{KL}(p_T \parallel \pi)$$

Upper bound

diffusion not, $\lambda(t) = f(t)^2$ \rightarrow $\lambda(t) = f(t)^2$ \rightarrow $\lambda(t) = f(t)^2$

증명 ① denote μ : joint distribution of diffusion process $\{X(t)\}_{t \in [0, T]}$

ν : joint distribution of process $\{\hat{X}_{\theta}(t)\}_{t \in [0, T]}$

② $D_{KL}(\mu \parallel \nu) \geq D_{KL}(p \parallel p_{\theta}^{SDE})$

Left side

Right side of proof

③ $D_{KL}(\mu \parallel \nu) = D_{KL}(p_T \parallel \pi) + E_{p_T(z)} [D_{KL}(\mu(\cdot | X(T)=z) \parallel \nu(\cdot | \hat{X}_{\theta}(T)=z))]$

이제 ① ② ③

→ $\lambda(t) = f(t)^2$

정리 $-E_{p(x)} [\log p_{\theta}^{SDE}(x)] \leq J_{SM}(\theta; f(\cdot)^v) + C_1 = J_{SM}(\theta; f(\cdot)^v) + C_2$

$\Rightarrow \lambda(t) = f(t)^2 =$ likelihood weighting

정리

$$p(x), q(x), p_+, q_+, q_F = \pi, S_\theta(x, z) = \nabla_x \log q_+(x) \cdot \frac{1}{\pi}$$

$$p_\theta^{SDZ} = p_\theta^{ODE} = q$$

$$D_{KL}(p \| p_\theta^{SDZ}) = J_{SM}(\theta | q, \cdot) + D_{KL}(p_T \| \pi)$$

증명

$$S_\theta(x, z) = \nabla_x \log q_+(x) \text{ 일 때, } \nabla p_\theta^{SDZ} = q$$

등식인 $\nabla_x \log q_+(x) \text{ 이므로, } p_\theta^{SDZ} = p_\theta^{ODE} = q$

⊕ Fokker-Planck equation

그러나, $S_\theta(x, z) = \nabla_x \log q_+(x)$ 어려움

⊕ p^{SDZ}, p^{ODE} 어떤 z 의 likelihood 더 좋은지 보장 못함.

⇒ p_θ^{SDZ} 와 p_θ^{ODE} 는 거의 비슷 ($S_\theta(x, z) \approx \nabla_x \log p_+(x)$ 같은 특별한 경우)

⊕ 경험적) training with likelihood weight.

→ $\frac{p^{ODE}}{p^{SDZ}}$ 의 likelihood 상용값 *
p

datapoint 이나 log-likelihood

$$J_{gm} \rightarrow \nabla_x \log p_\theta(x) \text{의 미분값} \Rightarrow \underline{J_{BGM}(\theta, \lambda(\cdot))}$$

사용

따라서, Theorem 1 \Rightarrow training에 대한 이론적 증명 \Rightarrow datapoint 이나 사용하는 것

Theorem 3

대체 가능

$$-\log p_\theta^{SDZ}(x) \leq L_\theta^{gm}(x) \stackrel{\text{대체 가능}}{=} L_\theta^{BGM}(x)$$

증명

$$\text{어떤 분포 } p \rightarrow -\mathbb{E}_{p(x)} [\log p_\theta^{SDZ}(x)] = D_{KL}(p \parallel p_\theta^{SDZ}) + H(p)$$

Numerical stability (수치적 안정성)

instability 한 경우, 작은 값 \Rightarrow 큰 영향력 있다.

$t \rightarrow 0$ (작은 noise) 일때, instability 하려 함을 알 수 있음.

해결방안 $\rightarrow 0$ 이 아닌, 작은 값 사용. 즉 $t \in [t, T]$