

Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
 - A. 이름: Generative Modeling by Estimating Gradients of the Data Distribution
 - B. 저널: NeurIPS
 - C. 도메인: Generative model, NCSN
 - D. 출판연도: 2019.7.12
 - E. 저자: Yang Song, Stefano Ermon
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. **논문 Figure를 그대로 따라 그리면 안됩니다.**
 - A. [선행연구] Generative model을 설계하는 가장 큰 요인은 P_{data} 분포를 따르는 데이터를 Sampling하고자 함에 있음. 같은 목표를 가지고 설계된 다양한 Generative model은 크게 두가지(논문 출판 시기 기준) 나눌 수 있다.
 - i. Log-likelihood based: $p(x; \theta) = \frac{1}{z(x)} q(x; \theta)$ 의 형태를 구하고자 함.
(Transformer-Auto Regression에서는, Chain형식으로 계산되는 구조로 (output -> 다음의 input) $z(x)$ 를 정의할 필요가 없다는 장점이 있음. 다음 예시로는 Energy based model같은 Surrogate loss 기반의 접근방식으로, 가장 큰 특징으로는 $z(x)$ 를 구하는 것이 어렵다는 점 때문에 이를 우회하는 구조를 가지고 있음.)
 - ii. GAN: Likelihood based가 아니며 따라서, PDF를 구하지 않는 구조로, divergences를 최소화하는 형식으로 계산함.
 - B. [선행연구] 앞서 두가지 model들은 특정한 구조를 사용해야 한다는 한계(예: Transformer, EBM 등)와 학습이 불안정하며, Implicit한 구조로 인해 다른 모델과의 성능 비교가 어렵다는 한계를 가짐(예: GAN)
 - C. 본 논문에서는 이러한 한계를 해결할 수 있는, 새로운 Generative model의 구조를 제안함. (NCSNs, Noise Conditional Score Networks)

- D. Model의 설계에서 목표로 하는 것은, 모든 데이터 공간(Low – High dimensional space)에서 일관적인 성능을 뽑아내는 것이라고 할 수 있음. 그러나 해당 과정을 위한 가장 큰 문제는 Low-dimensional에 있는 Data의 양이 매우 적다는 점임 (심지어 Score Function은 데이터의 미분 + log 로 구성됨 -> 즉, Low dimension이 실제로 들어가는 경우 값이 사라지는 경우도 있음) 이러한 문제는 단순히 Model의 low-dimension에 대한 표현력이 낮다는 것을 의미할 뿐만 아니라 Sampling(이 논문에서는 Langevin Dynamics) 성능을 일관적이지 못하게 함.
- E. 따라서, Low-Dimensional한 경우에도 강건한 모델의 성능을 위해 두가지 방안을 제안함.
- i. Perturbed Data
 - ii. Annealed Langevin Dynamics
- F. 앞서 이야기한 바와 동일하게, Low-dimension에 대한 표현력이 떨어지는 문제를 data + noise를 통해 해결하고 있음. (따라서, Denoising Score Matching을 목적함수로 이용)
- G. Langevin Dynamics는 Score Function(p_{data})만 있다면, Sampling을 할 수 있다는 장점이 있으나, 논문에서는 Annealed Langevin Dynamics가 필요하다고 설명하고 있음. $p_{data} = \pi * p_1(x) + (1 - \pi) * p_2(x)$ (e.g., $\pi = 0.7$)로 p_{data} 분포를 구성했을 때, 이상적인 결과는 Score Function (p_{data})의 Vector field에서 p_1 의 분포가 p_2 보다 선명한 경우라고 할 수 있음. 그러나, 한번의 Langevin Dynamics에서는 Score Function를 사용해서 구하다 보니, 결과적으로 $\pi : Scale Factor$ 가 제거되는 문제가 발생함. 즉, p_{data} 를 잘 설명하지 못하고 있다고 할 수 있음.
- H. 따라서, Annealed Langevin Dynamics가 필요한 것인데, 간단하게 이야기하자면 Langevin Dynamics에 대해서 나온 분포를 토대로 다시 Langevin Dynamics를 진행하고 다시 진행한다면 (논문에서, 10번) 모든 Dimension을 잘 반영하고 있는 model sampling이 나올 것이라는 내용임.
- I. 해당 논문에서는 결과적으로, 기존의 Generative model의 한계를 어느정도 해결하면서 우수한 성능(CIFAR-10, MNIST에서 GAN대비)의 model을 제안함. (안정적인 학습- GAN보다, Scalable – 다른 Auto Regressive, EBM 대비)

3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
- A. 문제점을 정확히 파악하고 (Generative model의 한계, Langevin dynamics) 실험적으로 보여줌과 동시에 해결방안을 제시함. 그렇게 제시한 해결방안이 실제 실험에서 우수한 성능임을 보임으로써 우수한 논문임을 보임.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
- A. 기존 모델의 문제점을 파악하고 -> 보완점을 제시한 것인지, 새로운 것을 적용하기 위해 문제점을 찾은 것인가가 궁금함.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. Sampling – 10번(NFE: 10) 즉, 속도가 느릴 수밖에 없음
- B. 뿐만 아니라, Score Matching! = Model. (서로 미분한 값에 대한 비교이므로)
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
- A. Diffusion 선행 연구
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
- A. 아니요

날짜: 2025-06-27

이름: 신준원

2025. 06. 27 Generative modeling by Estimating Gradients of Data Distribution

Langevin Dynamics

- 분자 동역학 + SDE (확률 미분 방정식)
- 임의의 MCMC 샘플링

ex) GCMC matching

① explicit 한 G2 모델

ex) Auto-regressive, surrogate loss

↳ chain rule 기반

Factorization 가능

$P(x)$ 를 표현 가능

($\propto P(x)$)

Introduction

생 (2019)

Generative model

likelihood-based

GAN

↳ implicit 한 G2

↓

likelihood-based → special G2

ex) RNN, transformer or surrogate loss

GAN: implicit → 학습 복잡

실제 GAN A model

GAN B model

학습, 메모리 어려움

추가 → ~~복잡성~~ → noise contrastive estimation

minimum probability flow 등

↳ low dimensional data에 대한 표현

각종 Generative model method: 한계 or 제한 존재

- ⊕ high dimension 만 생성
 - low dimension 만 생성
- 모든 dimension에 대해 동일한 표현

high Dimensional vs low Dimensional

개념: 차원 크기

↳ Generative modeling 방법

→ 이미지의 차원 (ex) mnist → 784 pixel) 은 high dimensional

숫자 생성, 음성은 pixel (차원) 은 low dimensional

즉, low dimensional에 대해서 학습할 때보다 high dimensional에 대해서 학습할 때 Generative model 지킬

⊕ high dimensional

제한 필요

크기, 학습 방법

모든 dimension

에 대해 학습

실용 가능한

Generative model 지킬

고정된 low dimension 에 대한 Gcore matching을 하는 것.

→ low dimension에 대한 노이즈 → 비대칭 (또는 0)

따라서 노이즈를 σ 줄여주는 $\text{noise} + \text{data}$ 을 통해 학습
(sampling) $\rightarrow \text{large noise} \Rightarrow \text{low dimension화}$

이후 Langevin dynamic 에 대해서

→ noise 를 크게 만들어 ~~학습~~ (low dimension)
진행

→ 점차 noise 를 줄여가며 ~~학습~~ (High dimension)
진행

모든 dimension 에 대해서 정보를 가진 sample 생성

+ 시각 x, 모든 Gcore network에 적용 가능, 효과적인 비교 진행 가능

likelihood-based 한계

GAN 한계

Score based generative modeling \rightarrow ~~matrix~~ \star

$$\frac{1}{n} \mathbb{E}_{p_{\text{data}}} [\| \nabla_{\theta} S_{\theta}(x) - \nabla_x \log p_{\text{data}}(x) \|_2^2]$$

$$\approx \mathbb{E}_{p_{\text{data}}} \left[\text{tr}(\nabla_x S_{\theta}(x)) + \frac{1}{n} \| \nabla_{\theta} S_{\theta}(x) \|_2^2 \right]$$

→ 제곱 항의 후
생략 가능

notation
[이제]

Jacobian of $S_{\theta}(x)$

Condition

가산함수

$S_{\theta}^*(x) = \nabla_x \log p_{\text{data}}(x)$

→ Vector valued function

각각에 대한 벡터 값 함수

tr: trace $\rightarrow \sum_{i=1}^d \frac{\partial^2}{\partial \theta_i^2} \log p_{\theta}(x) \rightarrow \text{tr}(\nabla_{\theta}^2 \log p_{\theta}(x))$

background

generative modeling $\rightarrow P_{data}$ 의 분포에서 sampling 하는 것

sampling을 하기 위한? $\rightarrow PDF(P_{data})$ 를 알아야 함.

이때 $P_{data} \rightarrow$ 큰 값 + 작은 값 $p_\theta(x) = \frac{1}{Z(\theta)} \cdot q(x)$

일반적인 확률분포 $= 1 \rightarrow q$ 은 확률 분포 1로 만들어주는 역할을 $Z(\theta)$

그러나 $Z(\theta) \Rightarrow$ 구하기 어렵음 \rightarrow 따라서 근사하는 다양한 solution 존재.

denoising Score matching

가장 큰 장점 $\rightarrow \text{tr}(\nabla_{\theta} S_{\theta}(x))$ 를 알면 충분 대체 가능

이때 \Rightarrow 대체 가능한 $\text{tr}(\nabla_{\theta} S_{\theta}(x))$

noise $\sim \mathcal{N}(0, I)$ Gaussian, $\tilde{x} = x + \text{noise}$

이때 \rightarrow noise가 Gaussian을 따르므로 $\text{tr}(\nabla_{\theta} S_{\theta}(x)) \rightarrow \frac{\|\nabla_{\theta} S_{\theta}(\tilde{x})\|^2}{2}$

objective $\rightarrow \frac{1}{2} \int q_{\theta}(\tilde{x}|x) P_{data}(x) [\|\nabla_{\theta} S_{\theta}(\tilde{x}) - \nabla_{\theta} \log q_{\theta}(\tilde{x}|x)\|^2]$

Sliced Score Matching

\rightarrow 장점 \Rightarrow Add noise 필요 없음.

\rightarrow 이분 간변을 해결하는 다른 방법들의 Score Matching 장점 \rightarrow 계산량 $\times 4$

$$\mathbb{E}_{P_{data}} \mathbb{E}_{P_v} \left[v^T \nabla_{\theta} S_{\theta}(x) v + \frac{1}{N} \|\nabla_{\theta} S_{\theta}(x)\|^2 \right]$$

P_v : simple distribution of Random Vector
e.g. multivariate normal

Sampling with Langevin dynamics

Langevin dynamics $\rightarrow \nabla_{\mathbf{x}} \log p(\mathbf{x})$ 사용 sampling 가능

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\epsilon}{\sqrt{t}} \nabla_{\mathbf{x}} \log p(\mathbf{x}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t$$

① $\mathbf{x}_0 \approx \pi(\mathbf{x})$, $\pi \rightarrow$ prior distribution

② $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$

③ $\mathbf{x}_T = p(\mathbf{z})$ ($\epsilon \rightarrow 0$ & $T \rightarrow \infty$)

④ Metropolis-Hastings 필요 ($\epsilon > 0$ & $T < \infty$)

\rightarrow error 수정해 필요 (그러나 skip하는 경우로 존재)

$\rightarrow \epsilon$ 이 충분히 작고, T 가 충분히 크다면 해당 error 무시 가능

\rightarrow 불 중요하게 무시

\rightarrow Diffusion
forwarding에 사용

Metropolis-Hastings 2.

base - 1 mcmc

(rotation은 대부분 해당 알고리즘에서 제공)

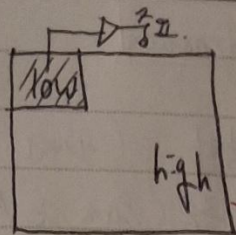
sampling 하려는 분포 $\pi(\mathbf{x})$ ($p_{\text{data}}(\mathbf{x})$)

\mathbf{x}_t 생성하기 생성할 분포 $q(\mathbf{y} | \mathbf{x}_t)$

$\mathbf{x}_t \rightarrow \mathbf{y}$ 생성 (편향은 $\pi(\mathbf{x})$ 가 크면 적음)

The manifold hypothesis

real world data. \rightarrow



2. \rightarrow high-dimensional

low-dimensional space

Score-based generative model

문제점 1. $\nabla_{\theta} \log p_{\text{data}}(x) = \text{high Dimensional}$

따라서 x 가 low dimensional 한 경우 표현이 어렵다

문제점 2. General 한 Score Matching objective function의 경우

$\rightarrow d \in \mathbb{R}^n$ 인데 d 의 값이 x 에 의존한다

2. $d \in \mathbb{R}^n$ 이 아닌 $d \in \mathbb{R}^k$ low-dimension 경우 문제가능

예) ResNet with CIFAR

Sliced Score Matching (상관 \Rightarrow noise x $\Rightarrow d \in \mathbb{R}^k$ 사용)

$\rightarrow x$ input 입력값은 x = x + noise

$\rightarrow x + \text{noise} = x$

문제점

\rightarrow low-dimensional 이므로 x 가 d 에 비해 d 가 더 크다

Data

(+) Score \rightarrow vector

\rightarrow data와 $S(x)$ 사이의 비교를 용이하게

\rightarrow high dimension 부공간에서만 비슷한 분포를 갖는다

Slow mixing of Langevin dynamics

문제 \rightarrow low dimension 에 취약함. (고차원 공간이 나뉘짐)

결과 \rightarrow Langevin dynamics 샘플링 ($p_{data} \approx p_g$) 어렵다.

$$\rightarrow p_{data} = \pi p_1(x) + (1-\pi)p_2(x) \quad \text{이때 } \pi \sim N(0, I)$$

Score function $\rightarrow \log p_{data} = \log p_1(x) + \log p_2(x)$

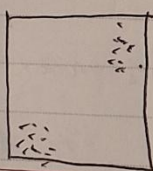
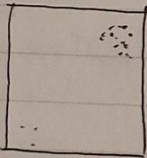
$$\rightarrow \nabla \log p_{data} = \nabla \log p_1(x) + \nabla \log p_2(x)$$

score functions 두개의 distribution으로 표현했으므로

\rightarrow π 에 해당하는 부분이 사라지는 결과를 피할 수 있음

~~이동하면~~

결과적으로 \rightarrow mixing 하도록



가변 반쪽

\rightarrow Langevin 한계 문제

Δ

3:7

\rightarrow

π 가 생각되면

Langevin 따지면

}

\rightarrow

5:5로 결과 나오는

결과

Learning VCSAs via Score Matching

→ $\sigma_{\text{iced}} \times \frac{1}{x} = \text{denoise}$

이런 Score Matching : denoising Score Matching → 중요!

$$\text{denoising Score Matching} \rightarrow \frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{x \sim \mathcal{N}(x, \sigma^2)} \left[\left\| s_{\theta}(x, \sigma) + \frac{x - \mu}{\sigma^2} \right\|^2 \right]$$

sigma (분산/변동) 제거됨

$$\mathcal{L}(\theta; \{\sigma_i\}_{i=1}^L) \triangleq \frac{1}{2} \sum_{i=1}^L \lambda(\sigma_i) \ell(\theta; \sigma_i)$$

$$s_{\theta}(x, \sigma) \approx s_{\theta}^*(x, \sigma)$$

정규화된 가중

Coefficient Function
→ $\lambda(\sigma) = \sigma^{-2}$

Langevin dynamics inference

$$\tilde{x}_t \leftarrow \tilde{x}_{t-1} + \left[\frac{\alpha_i}{\nu} \right] s_{\theta}(\tilde{x}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \epsilon_t$$

$$\alpha_i = \sigma_i^2 / \sigma_L^2 \quad (\text{정확도 감소}) \rightarrow \sigma_L \approx 0$$

σ_1 일수록 잘 샘플링했다면 → low dimensional 문제

(noise)

★ σ_1 은 가장 큰 σ_{L-1} 으로

즉, σ_1 은 σ_{L-1} 이므로 high dimensional 문제

(L=10 step)

VAE, GAN 의 성능을 평가할 때 → KL Divergence / KL Divergence / KL Divergence