

Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
 - A. 이름: Deep Unsupervised Learning using Nonequilibrium Thermodynamics
 - B. 저널: ICML
 - C. 도메인: Diffusion
 - D. 출판연도: 2015
 - E. 저자: Jascha Sohl-Dickstein et al.
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. 논문 Figure를 그대로 따라 그리면 안됩니다.
 - A. 이 논문은 VAE의 장점을 흡수하되, 단점을 보완하는 논문임.
 - B. [선행연구] VAE전, Generative model들의 경우, Inference($X \rightarrow Z$) 과정과, Model의 Probability 설계는 독립된 구조였음
 - C. [선행연구] 그러나, VAE의 경우, Inference - z를 통한 Generate 과정을 결합하여 해당 구조를 깨뜨리되, 우수한 성능을 보였음.
 - D. 다만, VAE의 단점은 Learning 과정에서 Inference + Generate과정이 결합되었다는 점임(학습이 한꺼번에 진행)
 - E. Diffusion에서는 이런 한계를 해결하고자 output 까지의 Process를 Forward + Reverse process로 구분함.
 - F. [Forward Process] 해당 과정에 대해서는 학습을 하지 않으며, Input: x_0 를 Gaussian Kernel 에 통과시키는 과정으로 구성됨.
 - G. 이는, $q(x_t) * \pi(t | t - 1) = q(x_{t+1})$, t는 0, T(T = 1000)의 값. 이때, Kernel은 $\pi \sim N(0, I)$ 를 따르며, $q(x_{0,1,2,...,T})$ 가 Gaussian Distribution을 따르도록 유도함. Markov Chain을 따르므로 해당 값은 $q(x_{0,1,2,...,T}) = q(x_0) \prod_1^T q(x_t | x_{t-1})$, 이때, $q(x_t | x_{t-1})$ 은 $N(x_t; x_{t-1} \sqrt{1 - \beta_t I \beta_t})$, β_t 는 Noise

- H. [Reverse Process] Notation: $p(x_0)$:Sample 결과라고 할 수 있음. 해당 값을 구하기 위한 과정은, Gaussian Distribution을 따른다는 점에 의거하여(β 가 충분히 작다면 = T가 충분히 크다면 Gaussian Distribution을 따름-Kernel이 Gaussian을 따르므로) μ, Σ 만으로 표현할 수 있다는 장점이 있음.
- I. 그러나 해당 과정을 위해서는 우선, Model Probability를 구해야 하며, 이는 다음과 같은 과정에 의해 구할 수 있음.
- J. $p(x_0) = \int dx_{0,1,2,\dots,T} * p_{0,1,2,\dots,T}$ 그러나, p에 대한 Joint Distribution은 계산이 불가능하므로 Tractable한 q에 대한 식으로 변환시키되, p와 q에 대해서, 동일한 Task (t시점일 때, $q(x_t | x_{t-1})/p(x_{t-1} | x_t)$, 심지어 Diffusion의 경우 Dimension이 다른 Generative Model과 달리 감소하지도 않아서 훨씬 어려움)
- K. Training(Loss Function)의 경우 Log Likelihood를 기반으로 학습하며, 이때 Intractable한 부분을 처리하기 위해서 KL-Divergence로 변환해주는 과정을 거침. (이때, 직접 Loss를 구하는 것이 불가능하다는 점 때문에, Jensen's Inequality를 통한, ELBO를 구하는 형태로 진행)
3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? 논문 **Contribution bullet**을 그대로 따라 적으면 안됩니다.
- A. **Diffusion 모델의 초기 과정 (Forward + Reverse)를 제안함.**
- B. Sampling(AIS) + Nonequilibrium Theory = 이미 존재했던 이론. 그러나 여러 이론 들을 결합해 Generative model에 적용했다는 점이 이 논문의 가장 큰 장점이라고 생각함.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
- A. 매번 읽을 때마다, 그렇구나 하고 읽었던 부분에 대해서 깊게 공부하고, 실제 모델의 의미에 대해서 추가적으로 파악해 나가는 과정속에서 부족했던 개념을 채워 나간다는 느낌이 듭.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. 복잡한 Loss Function (Lower Bound)의 존재는 학습을 어렵게 만들.
- B. Sampling까지 필요한 시간 (T=1000)이 너무 길

6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.

A. Diffusion 선행연구

7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?

A. 아니요

날짜: 2025-06-28

이름: 신준원

6.28 Deep Unsupervised Learning using Nonequilibrium Thermodynamics

AIS, Folker-Blanch, Holmognor

문제 0 제기

기존의 Generative model \rightarrow flexibility & tractability

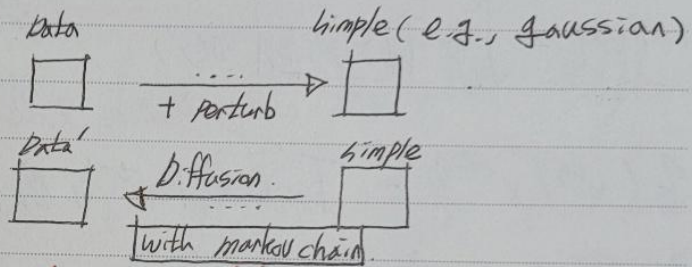
- tractability : 계산가능성
- flexibility : 미친 분포와 여러가지 표현가능성 \rightarrow trade-off

계산이 쉬우면 \rightarrow 분포가 편
분포 복잡하면 \rightarrow 계산 비용 \rightarrow 두가지 모두 만족하는 Generative model?

\rightarrow 고려할 다양한 방안 존재 \Rightarrow trade-off 라인을 선택

제한하는 새로운 노점 (가정사항)

- 모든 시이리 분포 A
- 샘플링 속도 A
- 다른 분포와 결합이 쉬움
- log likelihood 계산 쉬움

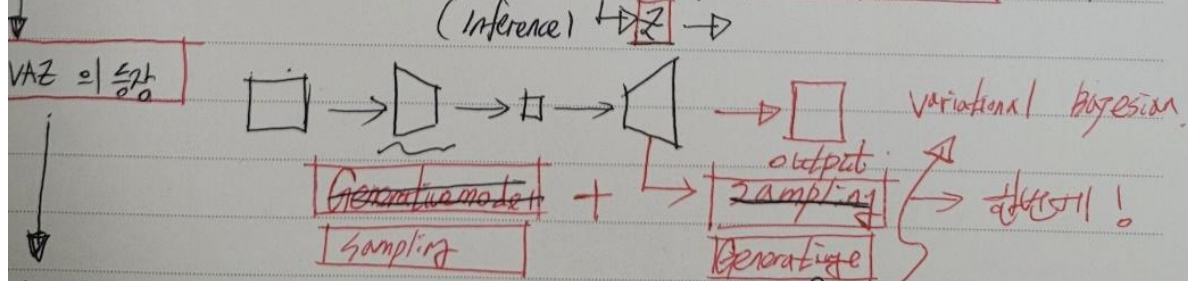


(non-equilibrium statistical physics)
+ Sequential MC

- \rightarrow step 여러개로 구성, 샘플링을 한번에 광-유
- \rightarrow Noise \rightarrow 여러 step으로 나눠 하여 \Rightarrow 다루기 쉬움

다른 방법의 과제

The Wake-sleep algorithm \rightarrow Probability model의 학습과 Sampling 과정은 무관함 \rightarrow 이어서는 새 방법



분 노점까지 증명하기 "한번에" 복잡도 (VAE와 자이로 점)

1. Wake-sleep (quasi-static process) 기반 설계 + annealed importance sampling
2. 다른 분포와 계산 쉬움
3. VAE: generator + sampler 사이의 상호 작용을 통해
4. lower (수렴)
5. upper + lower bound

Algorithm

Forward Process

Goal: Complex data distribution \rightarrow simple distribution

data distribution: $q(x^0)$ input \leadsto $\pi(y)$ sample \rightarrow Gaussian

with Markov chain kernel

$$\pi(y) = \int dy' T(y|y'; \beta) \pi(y') \quad / \quad y_{+1} = y$$

Diffusion Rate

$$q(x^t | x^{t-1}) = T_\pi(x^t | x^{t-1}; \beta_t)$$

$$\pi(x) = q(x^0) \cdot q(x^1 | x^0) \cdots q(x^T | x^{T-1})$$

$$q(x^{(0, \dots, T)}) = q(x^0) \cdot \prod_{t=1}^T q(x^t | x^{t-1}) \quad : \text{whole forward process}$$

$$q(x^t | x^{t-1}) = \mathcal{N}(x^t; x^{t-1} \sqrt{1 - \beta_t} \cdot I \cdot \beta_t) \quad : \text{Gaussian kernel}$$

- $\frac{1}{\beta_t} x$

$x^0 \rightarrow$ forward $f_T / p^T f_T \rightarrow$ Reverse

Reverse Process \rightarrow training target

$$p(x^T) = \pi(x^T)$$

$$p(x^{0 \dots T}) = p(x^T) \prod_{t=1}^T p(x^{t-1} | x^t)$$

정확점 $p(x^0)$

$$p(x^{t-1} | x^t) = \mathcal{N}(x^{t-1}; f_{\mu}(x^t, t), f_{\Sigma}(x^t, t))$$

계산량 $\propto T$

~~with $f_{\Sigma}(x^t, t)$~~ 계산량

model probability $\rightarrow p(x^0) = \int dx^{(1 \dots T)} \cdot p(x^{(1 \dots T)})$

$$p(x^0) = \int dx^{(1 \dots T)} p(x^{(1 \dots T)}) \frac{q(x^{(1 \dots T)} | x^0)}{q(x^{(1 \dots T)}) | x^0} \quad \Delta \text{Jensen's inequality}$$

$$p(x^0) = \int dx^{(1 \dots T)} \cdot q(x^{(1 \dots T)} | x_0) \cdot \frac{p(x^{(1 \dots T)})}{q(x^{(1 \dots T)} | x_0)}$$

$$p(x^0) = \int dx^{(1 \dots T)} \cdot q(x^{(1 \dots T)} | x_0) \left[p(x^T) \prod_{t=1}^T \frac{p(x^{t-1} | x^t)}{q(x^{t-1} | x^t)} \right] \quad \Delta$$

$(\beta \rightarrow \infty) \rightarrow$ 만약 서로 다른 종횡이면 $p(x^t)$ 은 \rightarrow 음이 아닌 계산가능

즉, noise를 step당 noise가 \rightarrow 무한으로, 서로 동일.
(이러한 경우...)

Training (log likelihood) $\rightarrow \int dx^0 p_{data}(x^0) \cdot \log p_\theta(x^0)$) $\frac{1}{N}$ \rightarrow 평균

$$L = \int dx^0 q(x^0) \cdot \log p(x^0)$$

$p(x^0)$ 은 양의 정의 \rightarrow 적용

$$L = \int dx^0 q(x^0) \cdot \log \left[\int dx^{(1 \dots T)} q(x^{(1 \dots T)} | x^0) \cdot p(x^T) \prod_{t=1}^T \frac{p(x^{t+1} | x^t)}{q(x^t | x^{t-1})} \right]$$

Jensen's inequality $f(E[x]) \leq E[f(x)]$ - 볼록
 $f(E[x]) \geq E[f(x)]$ - 오목 (log 오목)
 $\therefore \log E[x] \geq E[\log x]$

\rightarrow log를 integral 안으로 들어 올릴 수 있다!

\rightarrow Fixed $q(x^{(1 \dots T)} | x^0) \geq \frac{q(x^{(0 \dots T)})}{q(x^0)}$
 \rightarrow $\int dx^{(1 \dots T)}$ in $dx^{(1 \dots T)}$

$$L \geq \int dx^{(1 \dots T)} q(x^{(1 \dots T)} | x^0) \cdot \int dx^0 q(x^0) \log \left[p(x^T) \prod_{t=1}^T \frac{p(x^{t+1} | x^t)}{q(x^t | x^{t-1})} \right]$$

$$L \geq \int dx^{(0 \dots T)} q(x^{(0 \dots T)}) \cdot \log \left[p(x^T) \prod_{t=1}^T \frac{p(x^{t+1} | x^t)}{q(x^t | x^{t-1})} \right]$$

\rightarrow log \rightarrow 곱 \rightarrow 합으로 바뀐다.

$$L = \int \sum_{t=1}^T \frac{p(x^{t+1} | x^t)}{q(x^t | x^{t-1})} + \int dx^T \cdot q(x^T) \log p(x^T)$$

$p(x^T) = \pi(x^T) = q''(x^T)$

$-H_p(x^T) \rightarrow H(x) = - \int p(x) \log p(x) dx$

$$-H_p(x^T) + \int \sum_{t=1}^T \frac{p(x^{t+1} | x^t)}{q(x^t | x^{t-1})}$$

→ 가장자리 비정규화인 점 (해답률) $\beta(x'|x^0)$
edge effect $\Rightarrow z=0$ → forward 와 동일하게 설정

$$p(x^0|x') = \beta(x'|x^0) \cdot \frac{\pi(x^0)}{\pi(x')} = T_{\pi}(x^0|x'; \beta)$$

$$\left\{ L \geq \sum_{t=2}^T \int dx^{(0..T)} g^{(0..T)} \log \left[\frac{p(x^{t-1}|x^t)}{g(x^{t-1}|x^t)} \right] \right\} + \int dx^0 \cdot dx' g(x^0, x') \log \left[\frac{\beta(x'|x^0) \pi(x^0)}{\beta(x'|x^0) \pi(x')} \right]$$

ZLBO

\rightarrow ~~sec~~, $t=1$ → \rightarrow \rightarrow

$$p(x^0|x') = \frac{\beta(x'|x^0) \cdot \frac{\pi(x^0)}{\pi(x')}}{\beta(x'|x^0)} \rightarrow \text{sec.}$$

2eq2

Posterior with markov chain

$\rightarrow p(x_0)$ → x_0 부터 시작점 (input)

$p(x_{t-1}|x_t)$ 를 해답률 \Rightarrow 해답률 \Rightarrow x_0 Condition

markov chain $\rightarrow g(x_t|x_{t-1}) = g(x_t|x_{t-1}, x_0)$
 (forward)

$$\Rightarrow L \geq \sum_{t=2}^T \int dx^{(0..T)} g^{(0..T)} \log \left[\frac{p(x^{t-1}|x^t)}{g(x^{t-1}|x^t, x_0)} \right] - H_p(x^T)$$

[\rightarrow 해답률]

$$\frac{g(x^t|x^{t-1}, x_0)}{g(x^t|x^{t-1}, x_0)} = \frac{g(x^t, x^{t-1}, x_0)}{g(x^{t-1}, x_0)} \rightarrow \frac{g(x^{t-1}|x^t, x_0)}{g(x^{t-1}, x_0)} \cdot \frac{1}{g(x^t|x_0)}$$

$$= \frac{g(x^{t-1}|x^t, x_0)}{g(x^{t-1}|x_0)} \cdot \frac{g(x_0)}{g(x_0)} \cdot \frac{1}{g(x^t|x_0)}$$

$$g(x^t, x^{t-1}|x_0) = g(x^{t-1}|x^t, x_0) \cdot g(x^t|x_0)$$

$$g(x^t|x^{t-1}, x_0) = \frac{g(x^{t-1}|x^t, x_0) \cdot g(x^t|x_0)}{g(x_0)}$$

(Posterior) \star

$$\rightarrow g(x_{t-1}|x_t, x_0) = \frac{g(x_t|x_{t-1}, x_0) \cdot g(x_{t-1}|x_0)}{g(x_t|x_0)}$$

λ_0 를 Condition으로 제시 \rightarrow 해당 값 p 함수를 위해 사용 ($q \rightarrow p \rightarrow (q-p)$ 차이)

$$L \geq k = \sum_{t=2}^T \int dx^{(0...T)} q(x^{0...T}) \log \left[\frac{p(x^{t-1}|x^t)}{q(x^{t-1}|x^t, x^0)} \right] - H_p(x^T)$$

$$L \geq k = \sum_{t=2}^T \int dx^{(0...T)} q(x^{0...T}) \log \left[\frac{p(x^{t-1}|x^t)}{q(x^{t-1}|x^t, x^0)} \right] + \sum_{t=2}^T [H_q(x^t|x^0) - H_q(x^{t-1}|x^0)]$$

항에 대하여 정리

$D_{KL}(q||p) \rightarrow \mathbb{E}(\log q - \log p)$

$H_q(x^t|x^0) - H_q(x^{t-1}|x^0) + H_q(x^t|x^0) - H_q(x^t|x^0)$
 $\rightarrow H_q(x^t|x^0) - H_q(x^{t-1}|x^0)$ 안 씌울 것
 $H_q(x^T|x^0) - H_q(x^1|x^0) - H_p(x^T)$

$$L \geq k = - \sum_{t=2}^T \int dx^0 \cdot dx^t q(x^0, x^t) \cdot D_{KL}(q(x^{t-1}|x^t, x^0) || p(x^{t-1}|x^t)) +$$

\rightarrow KLD를 포함한 항 "k" \rightarrow Variational Bayesian Method (VBM) 사용

① T 충분히 큰 경우 부등호 \rightarrow 등호로 가능 ($p \rightarrow q$)

① 결과적으로 중복한 T \rightarrow Gaussian distribution

\rightarrow m.s.e에 대해 regression - L2 loss 동일.

β Scheduling. (noise scheduling)

$\beta_1 \rightarrow$ over fitting 예방 (과적합 예방) - 고정

$\beta_2 \dots T \rightarrow$ 증가함.

↳ Denoising on inputting < | ~~필요한~~ 정보.

Multiplying Distributions and Computing Posteriors

$$\tilde{p}(x_0) \propto p(x_0) \cdot \boxed{\gamma(x_0)}$$

→ Diffusion process!

intermediate distribution function → noise

$$\tilde{p}(x^{0 \dots T}) = \frac{1}{Z_T} \cdot p(x^T) \cdot \gamma(x^T)$$

$$\boxed{\tilde{p}(x^t) = \frac{1}{Z_t} \cdot p(x^t) \cdot \gamma(x^t)}$$

↓ $\gamma(x^t)$

1 step at a time

$$\boxed{\tilde{p}(x^t) = \int dx^{t+1} \tilde{p}(x^t | x^{t+1}) \cdot \tilde{p}(x^{t+1})}$$

$$\frac{p(x^t) \cdot \gamma(x^t)}{Z_t} = \int dx^{t+1} \tilde{p}(x^t | x^{t+1}) \cdot \frac{p(x^{t+1}) \cdot \gamma(x^{t+1})}{Z_{t+1}}$$

양쪽

$$p(x^t) = \int dx^{t+1} \tilde{p}(x^t | x^{t+1}) \cdot \frac{Z_t \gamma(x^{t+1})}{Z_{t+1} \gamma(x^t)} \cdot p(x^{t+1})$$

$$\boxed{p(x^t) = \int dx^{t+1} \tilde{p}(x^t | x^{t+1}) \cdot p(x^{t+1})} \quad \text{이므로 } p, \tilde{p} \text{의 관계가 성립}$$

$$\tilde{p}(x^t | x^{t+1}) = p(x^t | x^{t+1}) \cdot \boxed{\frac{Z_{t+1} \gamma(x^t)}{Z_t \gamma(x^{t+1})}}$$

weight ratio

→ 이 값을 곱하여

각에 곱함

$$\tilde{p}(x^t | x^{t+1}) = \boxed{\frac{1}{Z_t \gamma(x^{t+1})}} \cdot p(x^t | x^{t+1}) \cdot \gamma(x^t)$$

(normalization)

normalization

Forward process / Reverse process algorithm

↳ weight, γ

based AIS algorithm