

Score-based + Annealing Langevi Dynamics.

* Introduction.

likelihood-based or GAN은 성능적이지만 불안정한 문제가 있다.

likelihood는 surrogate loss를 쓰거나, flow-based의 경우 reversible한 arch이 필요. (or Autoregressive)

GAN은 unstable, model compare가 힘들 (직접적인 data 분포의 NLL과 달리 낮은 Loss가 훨씬 보잘것없다)

논문은 pdf의 gradient를 score ($\nabla \log p_\theta(x)$)로 삼고, model은 score를 추정하고,

이를 바탕으로 $p_\theta(x)$ 를 높은 것으로 이동.

Challenge.

1. low-dimension의 dataset의 경우에, 공간의 space에서 score 추정이 힘들.

2. certain noise에 respect sampling 시에, 정확하지 않음.

↳ perturbated noise로 대체

0) 반복은 tractable 하기, objective function은 model에 따라 사용 가능.

* Model

pdf: $P(x)$

score: $\nabla_x \log P(x)$

network: $S_\theta(x) \rightarrow P_{\text{data}}$ 의 case score

* Score Matching

$$\Theta^* = \arg \min_{\theta} E_{P(x)} \left[\frac{1}{2} \left\| \nabla_x \log P(x) - S_\theta(x) \right\|_2^2 \right]$$

P(x) 예상 x
Score matching

$$\Theta^* = \arg \min_{\theta} E_{P(x)} \left[\frac{1}{2} \|S_\theta(x)\|_2^2 + \text{tr}(\nabla_x S_\theta(x)) \right]$$

Jacobian은 이미지가 x에 걸 수록 예상 가능 ↑

- Denoising Score Matching

perturbated x score

$$\Theta^* = \arg \min_{\theta} E_{P(x)} \left[\frac{1}{2} \|S_\theta(x)\|_2^2 + \text{tr}(\nabla_x S_\theta(x)) \right]$$

Denoising Score Matching

$$\frac{1}{I} \mathbb{E}_{q_\phi(\tilde{x}|x) P_{\text{data}}(x)} \left[\|S_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log q_\phi(\tilde{x}|x)\|_2^2 \right]$$

Gaussian 모형: $\frac{1}{\sigma^2} (\tilde{x}-x), \Sigma^{-1} (\tilde{x}-x)$

- Sliced Score Matching

projection

$$\mathbb{E}_{p_v} \mathbb{E}_{P_{\text{data}}} \left[v^\top \nabla_x s_\theta(x) v + \frac{1}{2} \|s_\theta(x)\|_2^2 \right],$$

DSM 보다 예상 가능성이 많아서 좋음.

* Sampling with Langevin dynamics

Score function을 PCN으로부터 sample 하는 방법

$$\tilde{x}_t = \tilde{x}_{t-1} + \frac{\epsilon}{2} D_\phi \log P(\tilde{x}_{t-1}) + \sqrt{\epsilon} z_t$$

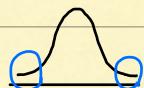
* Challenges of Score-based GM.

* Manifold hypothesis

All data in high-dimensional ambient spaces are linear

문제점.

1. ambient manifold is low dimension while score manifold is high.



2. score matching & data are linearly independent 가능.

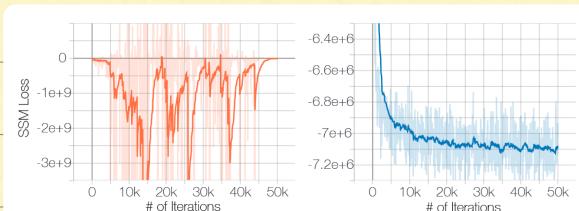


Figure 1: Left: Sliced score matching (SSM) loss w.r.t. iterations. No noise is added to data. Right: Same but data are perturbed with $\mathcal{N}(0, 0.0001)$.

이미지의 속성이 알아보기 힘든 작은 깊은 노이즈로 높은 확률로 혼란.

* Low data density regions

Dataset이 부족할 때 일어나는 현상.

데이터가 빠른 곳은 예상이 잘 안됨.

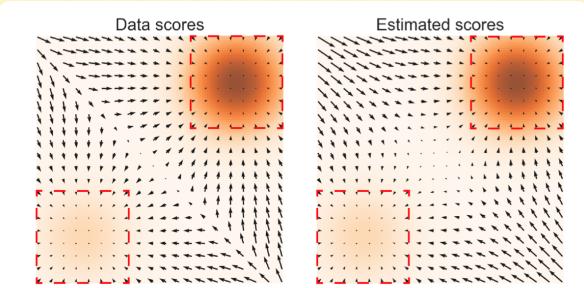


Figure 2: **Left:** $\nabla_x \log p_{\text{data}}(x)$; **Right:** $s_\theta(x)$. The data density $p_{\text{data}}(x)$ is encoded using an orange colormap: darker color implies higher density. Red rectangles highlight regions where $\nabla_x \log p_{\text{data}}(x) \approx s_\theta(x)$.

- Slow Mixing of Langevin dynamics

Langevin dynamics의 느雠로 사이를 잘 modeling 못하고, 상대적 분포를 깨짐.

$P_{\text{data}}(x) = \pi P_1(x) + (1-\pi)P_2(x)$ 일 때, score는 π 의 영향을 안받아서, 두 분포의 상대비율도 흐짐이 있음.

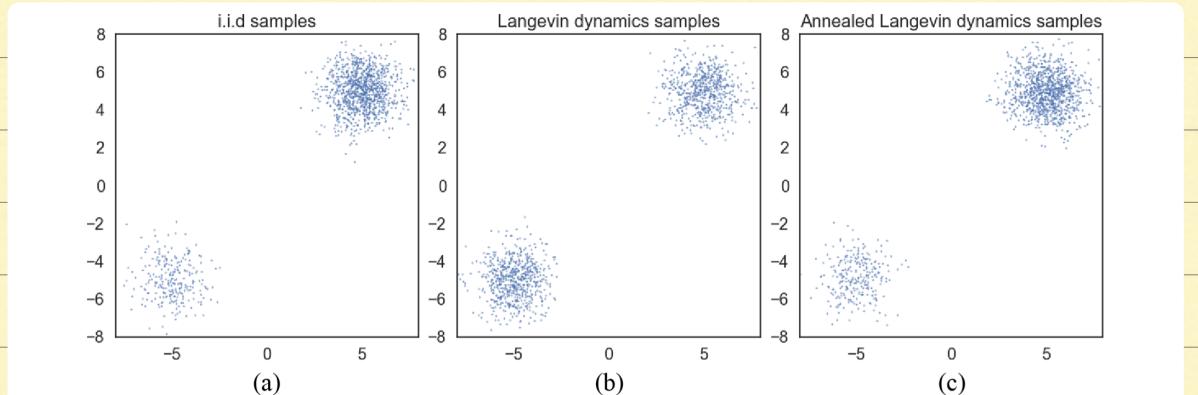


Figure 3: Samples from a mixture of Gaussian with different methods. (a) Exact sampling. (b) Sampling using Langevin dynamics with the exact scores. (c) Sampling using annealed Langevin dynamics with the exact scores. Clearly Langevin dynamics estimate the relative weights between the two modes incorrectly, while annealed Langevin dynamics recover the relative weights faithfully.

* Noise condition Score Network : learning & Infer.

Gaussian은 whole space에서 있을 때 low-dimension manifold에 고정됨. → score 흐름 ↑

" origin data from low-density 영역에 고정되는 경향이 있다.

∴ 다양성은 noise-level or data modeling.

[2] Langevin MCMC는 noise level Σ.

* NCSN

Score Net의 noise를 condition으로 험.

$t \rightarrow 0$ 일 때 σ_t 가 작아짐. $\frac{\sigma_t}{\sigma_{t+1}} > 1$

$$s_\theta(x, \sigma) \approx \nabla_x \log p_\theta(x)$$

Network는 일단 segmentation의 유형은 UNet 사용

* Learning NCSNs

$$q_\sigma(\tilde{x}|x) = \mathcal{N}(\tilde{x}|x, \sigma^2 I)$$

$$\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) = -\frac{1}{\sigma^2}(\tilde{x}-x)$$

$$DSM: \ell(\theta; \sigma) \triangleq \frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 I)} \left[\left\| s_\theta(\tilde{x}, \sigma) + \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right].$$

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

$$\text{정규화되었을 때, } \|s_\theta(x, \sigma)\| \propto 1/\sigma \text{ 일 때 } \lambda \propto \sigma^{-2}. \therefore \lambda(\sigma) = \sigma^{-2}$$

• NCFN inference via annealed Langevin Dynamics

Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

- 1: Initialize $\tilde{\mathbf{x}}_0$
- 2: **for** $i \leftarrow 1$ to L **do**
- 3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ ▷ α_i is the step size.
- 4: **for** $t \leftarrow 1$ to T **do**
- 5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$
- 6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7: **end for**
- 8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$
- 9: **end for**
- return $\tilde{\mathbf{x}}_T$

* Experiments

Model	Inception	FID
CIFAR-10 Unconditional		
PixelCNN [59]	4.60	65.93
PixelIQN [42]	5.29	49.46
EBM [12]	6.02	40.58
WGAN-GP [18]	7.86 ± .07	36.4
MoLM [45]	7.90 ± .10	18.9
SNGAN [36]	8.22 ± .05	21.7
ProgressiveGAN [25]	8.80 ± .05	-
NCSN (Ours)	8.87 ± .12	25.32
CIFAR-10 Conditional		
EBM [12]	8.30	37.9
SNGAN [36]	8.60 ± .08	25.5
BigGAN [6]	9.22	14.73

Table 1: Inception and FID scores for CIFAR-10

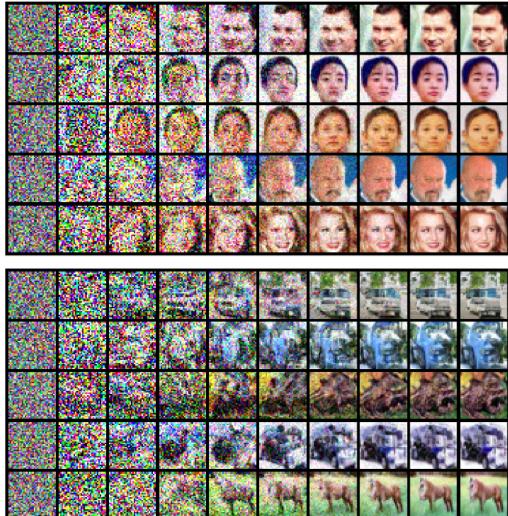


Figure 4: Intermediate samples of annealed Langevin dynamics.



(a) MNIST



(b) CelebA



(c) CIFAR-10

Figure 5: Uncurated samples on MNIST, CelebA, and CIFAR-10 datasets.



Figure 6: Image inpainting on CelebA (**left**) and CIFAR-10 (**right**). The leftmost column of each figure shows the occluded images, while the rightmost column shows the original images.

Algorithm 2 Inpainting with annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$ $\triangleright \epsilon$ is smallest step size; T is the number of iteration for each noise level.
Require: \mathbf{m}, \mathbf{x} $\triangleright \mathbf{m}$ is a mask to indicate regions not occluded; \mathbf{x} is the given image.

```

1: Initialize  $\tilde{\mathbf{x}}_0$ 
2: for  $i \leftarrow 1$  to  $L$  do
3:    $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ 
4:   Draw  $\tilde{\mathbf{z}} \sim \mathcal{N}(0, \sigma_i^2)$   $\triangleright \alpha_i$  is the step size.
5:    $\mathbf{y} \leftarrow \mathbf{x} + \tilde{\mathbf{z}}$ 
6:   for  $t \leftarrow 1$  to  $T$  do
7:     Draw  $\mathbf{z}_t \sim \mathcal{N}(0, I)$ 
8:      $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$ 
9:      $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t \odot (1 - \mathbf{m}) + \mathbf{y} \odot \mathbf{m}$ 
10:    end for
11:     $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$ 
12: end for
return  $\tilde{\mathbf{x}}_T$ 

```

* NCSN ++

* Introduction

NCSN은 score estin으로 Langevin 등은 noise를 제거한다.

DDPM은 noise와 process를 modeling 한다.

Continous의 경우 DDPM의 loss는 각 noise scale에 score를 implicit compute 한다.

Score-based의 학습을 위해, SDE를 통해, 위의 방법을 통합

유한한 data에 pertubing 하는 대신.

Diffusion process의 따라, 시간이 지나면서 넓어지는 continuum of distribution을 제안한다.

↳ data는 random noise로 diffusion

↳ 이 과정은 trainable한 param을 data에 의존하지 않는, 고정된 SDE로 의해 주어짐.

Reverse process의 noise → sample 가능.

↳ Reverse-time SDE를 만족하고 시간의 흐름으로 marginal probability density의 score

는 아니면 forward SDE로도 가능.

∴ Reverse-time SDE의 score의 대신 t 조건의 NN을 훈련하여 사용 가능.

↳ numerical SDE solver로 sample 가능

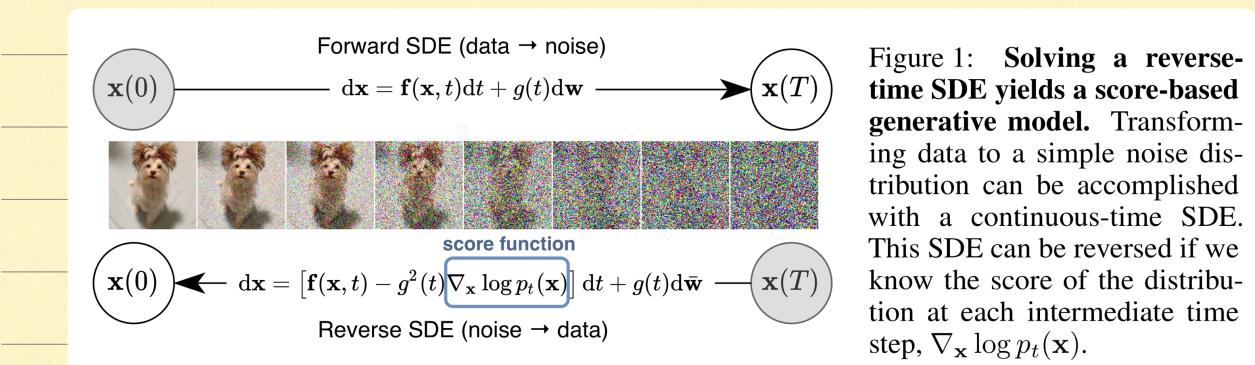


Figure 1: Solving a reverse-time SDE yields a score-based generative model. Transforming data to a simple noise distribution can be accomplished with a continuous-time SDE. This SDE can be reversed if we know the score of the distribution at each intermediate time step, $\nabla_x \log p_t(x)$.

* contribution.

1. Flexible sampling & likelihood compute.

reverse-time SDE or general-purpose SDE solvers 더 넓은 범위.

특히

1. Predictor-Corrector (PC) : SDE-solver or score-based MCMC 포함.

↳ 같은 방법 동작.

2. ODE-based deterministic sampler

↳ black-box ODE solver를 통한更具 adaptive sampling.

↳ latent를 통한 flexible manipulation.

↳ likelihood compute

↳ unique identifiable encoding

3. Controllable generation

Conditional reverse-time SDE는 unconditional score/mim estim 가능성이 있다.

훈련 데이터가 있는지 없는지 condition information으로 generation process 조절 가능.

↳ inpainting, class condition 등을 모두 재현하지 않고 생성 가능.

4. Unified framework

여러 score-based model은 unity

* Background

* Denoising Score Matching with Langevin Dynamics (SMLD)

= NCSN

* DDPM

* Score-based generative modeling with SDEs

Multi-level noise의 perturbing of 학습이다.

한번에서는 noise scale은 infinite로 generalize 가능,

perturbated data의 distribution이 SDE의 흐름을 update 되도록 한다

* Perturbing data with SDEs

$X(0) \sim P_0, X(T) \sim P_T$ or $t \in [0, T]$ 의 diffusion process $\{X(t)\}_{t=0}^T$

$\nearrow X(0)$ 이 tractable. \nwarrow continuous

Diffusion process는 SDE로 modeling 가능.

$$dx = f(x,t)dt + g(t)dw$$

$\hookrightarrow w$: standard Wiener process (우연성이 미지인가)가 움직이는 모양

$\hookrightarrow f(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, drift coefficient of $x(t)$

$\hookrightarrow g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, diffusion coefficient of $x(t)$

coefficient가 state에 대해 Lipschitz 일 때, SDE는 unique strong solution.

$X(t)$ 의 확률은 $P_t(x)$. $P_{s,t}(x(t) | x(s))$, $0 \leq s < t \leq 1$

Distribution이 고정된 사건 분포로 놓은 $x(0) \rightarrow x(T)$

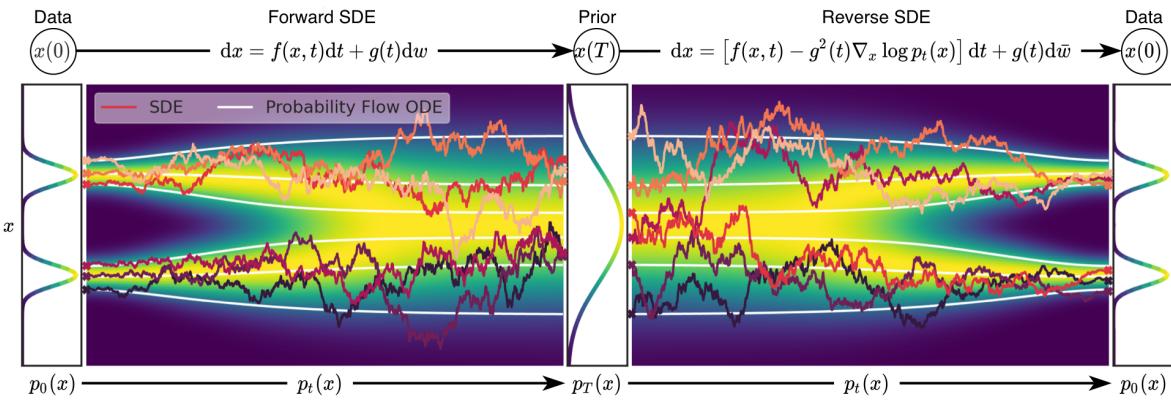


Figure 2: **Overview of score-based generative modeling through SDEs.** We can map data to a noise distribution (the prior) with an SDE (Section 3.1), and reverse this SDE for generative modeling (Section 3.2). We can also reverse the associated probability flow ODE (Section 4.3), which yields a deterministic process that samples from the same distribution as the SDE. Both the reverse-time SDE and probability flow ODE can be obtained by estimating the score $\nabla_x \log p_t(x)$ (Section 3.3).

* Generating samples by reversing the SDE

Forward \rightarrow diffusion step, reverse \leftarrow diffusion.

$$dx = [f(x, t) - g(t)^2 \nabla_x \log p_t(x)]dt + g(t)d\bar{w}, \quad t = T \rightarrow 0$$

↙ 모든 텝이 대수 일때, p_0 가능.

* Estimating score for SDE

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{x(0)} \mathbb{E}_{x(t)|x(0)} \left[\| s_{\theta}(x(t), t) - \nabla_{x(t)} \log p_{0t}(x(t) | x(0)) \|_2^2 \right] \right\}.$$

↳ ADVIM DSM 사용 (Sliced SMC 사용)

↳ 확장으로 풀기위해 transition kernel $p_{0t}(x(t) | x(0))$ 을 알아야 한다

↳ $f(\cdot, t) \sim$ affine 확장, transition kernel은 항상 Gaussian distribution이다. (mean, var는 closed form)

* Example VE, VP SDEs and beyond

SMLD와 DDPM은 SDE의 각각 다른 이산화.

Noise scale in SMLD와

$$X_i = X_{i-1} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2}, Z_{i-1}, i=1, \dots, N$$

$N \rightarrow \infty$ 이면, $\sigma_i \approx \sigma(t)$, $Z_i \sim Z(t)$ 가 됨 (continuous, $t=[0, T]$)

$$\{X(t)\}_{t=0}^T \text{의 SDE : } dx = \sqrt{\frac{d[\sigma^2(t)]}{dt}} dw \quad \dots \textcircled{1}$$

DDPM의 M은

$$X_i = \sqrt{1-\beta_i} X_{i-1} + \sqrt{\beta_i} Z_{i-1}, i=1, \dots, N$$

$N \rightarrow \infty$ 이면,

$$\{X(t)\}_{t=0}^T \text{의 SDE : } dx = -\frac{1}{2}\beta(t)x dt + \sqrt{\beta(t)} dw \quad \dots \textcircled{2}$$

$t \rightarrow \infty$ 일 때 $\textcircled{1}$ 은 variance가 exploding 된 CL. \Rightarrow VE SDE

$\textcircled{2}$ 는 fixed variance. \Rightarrow VP SDE

VP로부터 likelihood 잘 맞는 SDE 찾기.

$$dx = -\frac{1}{2}\beta(t)x dt + \sqrt{\beta(t)(1 - e^{-2 \int_0^t \beta(s) ds})} dw.$$

L const $\beta(t)$ 를 뜻 때, 충분히 짧은 time step 만 VP SDE의 variance bound 찾기.

L :: sub-VP SDE

VE, VP, sub-VP 은 affine drift coefficient 를 갖는다. $P_{tt}(Y(t)|X(0))$ 은 가우시안으로 closed form으로

* Solving the Reverse SDE

time dependent score-based S_θ 를 훈련하고 낸 후,

reverse-time SDE를 구성하고, numerical approach로 Point sample을 생성할 수 있다.

* General-purpose numerical SDE solver

Numerical solver SDE는 trajectory를 계산한다.

Euler-maruyama or runge-kutta 같은 general-purpose numerical methods

학습률이 다른 stochastic dynamics의 discretization이며, 이를 sample 생성 가능.

All SDE의 경우 reverse diffusion sampler를 제안한다.

↳ DDPM or SMLD의 ancestral sampling과 유사하다.

Table 1: Comparing different reverse-time SDE solvers on CIFAR-10. Shaded regions are obtained with the same computation (number of score function evaluations). Mean and standard deviation are reported over five sampling runs. “P1000” or “P2000”: predictor-only samplers using 1000 or 2000 steps. “C2000”: corrector-only samplers using 2000 steps. “PC1000”: Predictor-Corrector (PC) samplers using 1000 predictor and 1000 corrector steps.

FID↓ Predictor	Variance Exploding SDE (SMLD)				Variance Preserving SDE (DDPM)			
	P1000	P2000	C2000	PC1000	P1000	P2000	C2000	PC1000
ancestral sampling	4.98 ± .06	4.88 ± .06		3.62 ± .03	3.24 ± .02	3.24 ± .02		3.21 ± .02
reverse diffusion	4.79 ± .07	4.74 ± .08	20.43 ± .07	3.60 ± .02	3.21 ± .02	3.19 ± .02	19.06 ± .06	3.18 ± .01
probability flow	15.41 ± .15	10.54 ± .08		3.51 ± .04	3.59 ± .04	3.23 ± .03		3.06 ± .03

→ Predictor - Corrector Samplers

MCMC은 sampling은 numerical SDE solver로 막을 끌어.

각 time step마다 SDE solver는 Cholesky step을 sample은 estimate predictor.

Score-based MCMC는 Corrector는, estimate sample은 marginal distribution의 모양을 예상

↳ Predictor & Corrector (PC)

↳ SMULD와 DDPN을 활용한

SMULD: identity function은 predictor, LD는 corrector.

DDPM: ancestral sampling은 predictor, identity function은 corrector

PC가 다른 것보다 성능↑

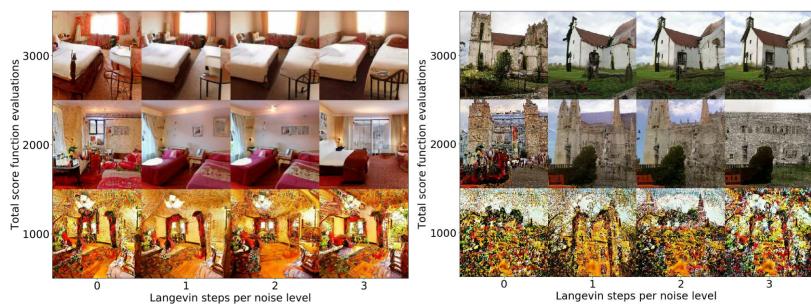


Figure 9: PC sampling for LSUN bedroom and church. The vertical axis corresponds to the total computation, and the horizontal axis represents the amount of computation allocated to the corrector. Samples are the best when computation is split between the predictor and corrector.

⇒ Probability flow and connection to neural ODE.

Score-based는 reverse-time SDE의 numerical method를 계승한다

모든 diffusion process는 trajectory가 동일한 marginal probability distribution $\{P_t(x)\}_{t=0}^T$ 을 SDE로써

공유하는 deterministic process) NCL. ⇒ SDE 뿐만 아니라 결론은 같다라고 말하는듯?

or deterministic process는 다음 만족

$$dx = \left[f(x, t) - \frac{1}{2} g(t)^2 \nabla_x \log p_t(x) \right] dt, \Rightarrow \text{probability flow ODE}$$

- Exact likelihood computation.

기존 DDPM의 경우 가능성이 훨씬 낫다.

- Manipulating latent representation.

$x_0 \in X_T$ encoding의 경우.

latent x_T 의 manipulating 가능

- Uniquely identifiable encoding

unique but invertible SCL.

동일한 data, capability, optimization accuracy가 있으면, x_0 와 x_T 가 unique does not happen

$dx = f(x, t) dt + g(t) dw \rightarrow$ trainable parameter

ODE의 normalize flow $dx = [f(x, t) - \frac{1}{2} g(t)^2 \nabla_x \log P_t(x)] dt$ 로 얻어진 esim score가 필요

같은 계산을 고려

- Efficient sampling.

* Architecture Improvement.

★ Controllable generation.