

CS 59300 – Algorithms for Data Science

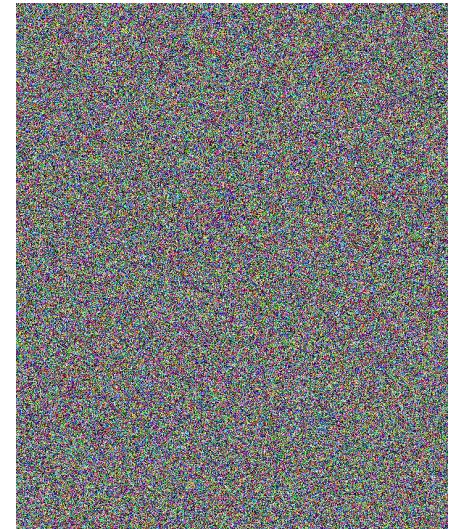
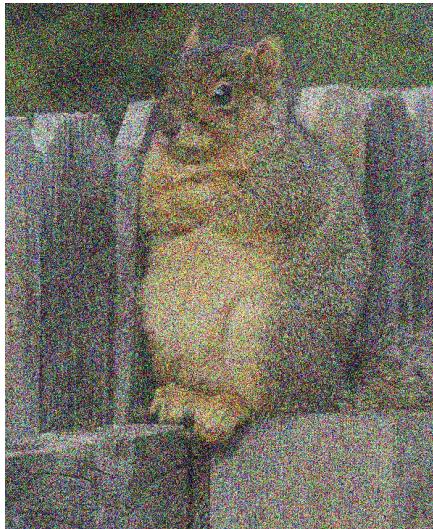
Classical and Quantum approaches

Lecture 13 (10/23)
Diffusion Models

https://ruizhezhang.com/course_fall_2025.html

Diffusion models

Forward process: Ornstein-Uhlenbeck



$t = 0$

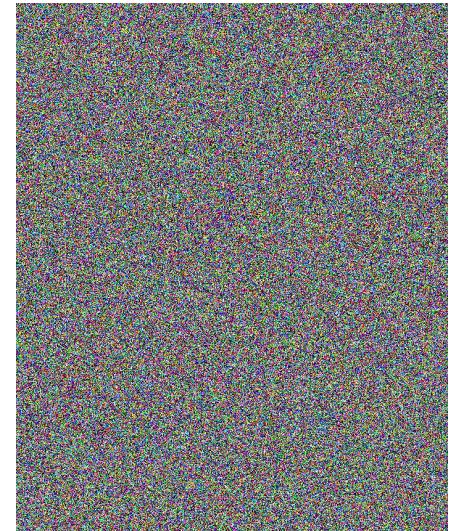


$t = \infty$

$$\begin{aligned} dX_t &= -X_t + \sqrt{2}dB_t \\ X_0 &\sim q \text{ (data distribution)} \end{aligned}$$

Diffusion models

Forward process: Ornstein-Uhlenbeck



$t = 0$



$t = \infty$

$$X_t = e^{-t} \cdot X_0 + \sqrt{1 - e^{-2t}} \cdot g \quad \text{for } g \sim N(0, I)$$

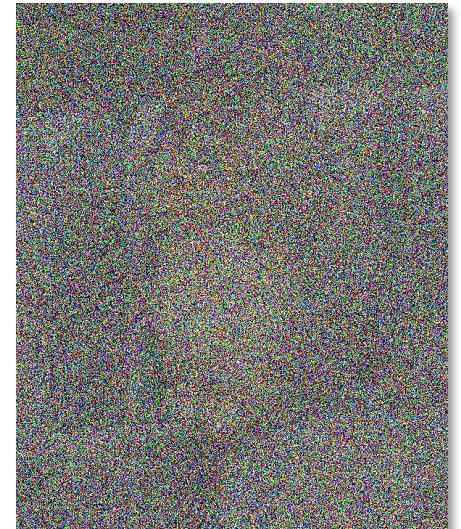
$X_0 \sim q$ (data distribution)

(prove this latter)

Diffusion models

Forward process:

$$q_0 = q$$



$t = 0$



$t = T$

$$\begin{aligned} dX_t &= -X_t + \sqrt{2}dB_t \\ X_0 &\sim q \text{ (data distribution)} \end{aligned}$$

in practice,
can only run
for finite time

Diffusion models

Reverse process:



$t = 0$

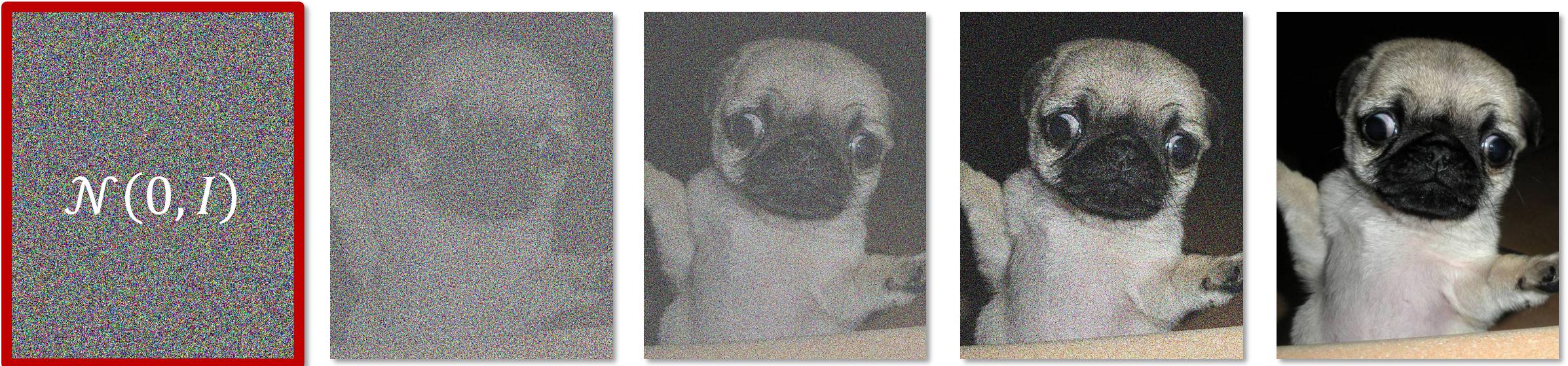
$$dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$$

(prove this latter)

$$X_0^\leftarrow \sim q_T \text{ (forward process at time } T\text{)}$$

Diffusion models

To sample fresh images, run reverse process with **Gaussian initialization**



$t = 0$

$t = T$

$$dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$$

$X_0^\leftarrow \sim q_T$ (forward process at time T)

“score function”
How to estimate it?

Score matching

Tweedie's formula. Given $\tilde{x} = x + e$ for $x \sim p$ and $e \sim \mathcal{N}(0, \sigma^2 I)$,

$$\mathbb{E}[x | \tilde{x}] = \tilde{x} + \sigma^2 \cdot \nabla \ln \tilde{p}(\tilde{x})$$

where \tilde{p} is the density for \tilde{x}

$$\Leftrightarrow \hat{e}_{\text{Bayes}} = -\sigma^2 \nabla \ln \tilde{p}(\tilde{x})$$

Song-Ermon '19: reduce estimating the score function to a **supervised learning task**:

Given noisy image X_t , predict noise γ that was added

$$X_t = \underbrace{e^{-t} \cdot X_0}_x + \underbrace{\sqrt{1 - e^{-2t}} \cdot g}_{e \sim \mathcal{N}(0, (1 - e^{-2t})I)}$$

Score matching

Tweedie's formula. Given $\tilde{x} = x + e$ for $x \sim p$ and $e \sim \mathcal{N}(0, \sigma^2 I)$,

$$\mathbb{E}[x | \tilde{x}] = \tilde{x} + \sigma^2 \cdot \nabla \ln \tilde{p}(\tilde{x})$$

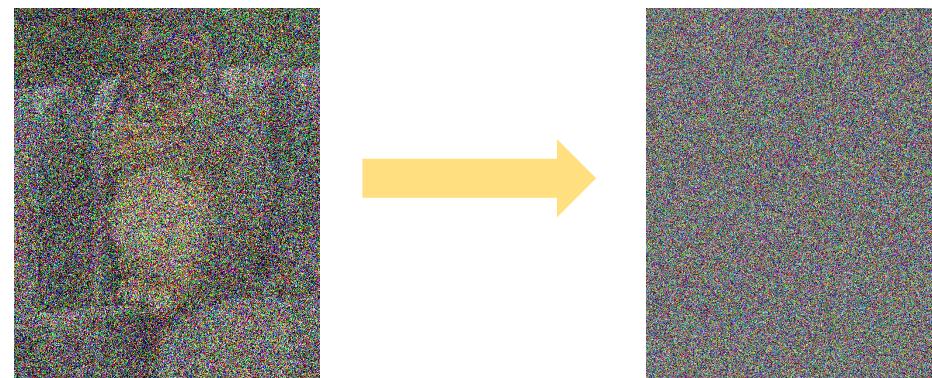
where \tilde{p} is the density for \tilde{x}

$$\Leftrightarrow \hat{e}_{\text{Bayes}} = -\sigma^2 \nabla \ln \tilde{p}(\tilde{x})$$

Song-Ermon '19: reduce estimating the score function to a **supervised learning task**:

Given noisy image X_t , predict noise γ that was added

Fit a neural net to training examples to drawn from q_t



Score matching

Tweedie's formula. Given $\tilde{x} = x + e$ for $x \sim p$ and $e \sim \mathcal{N}(0, \sigma^2 I)$,

$$\mathbb{E}[x | \tilde{x}] = \tilde{x} + \sigma^2 \cdot \nabla \ln \tilde{p}(\tilde{x})$$

where \tilde{p} is the density for \tilde{x}

$$\Leftrightarrow \hat{e}_{\text{Bayes}} = -\sigma^2 \nabla \ln \tilde{p}(\tilde{x})$$

Song-Ermon '19: reduce estimating the score function to a **supervised learning task**:

Given noisy image X_t , predict noise γ that was added

Fit a neural net to training examples to drawn from q_t

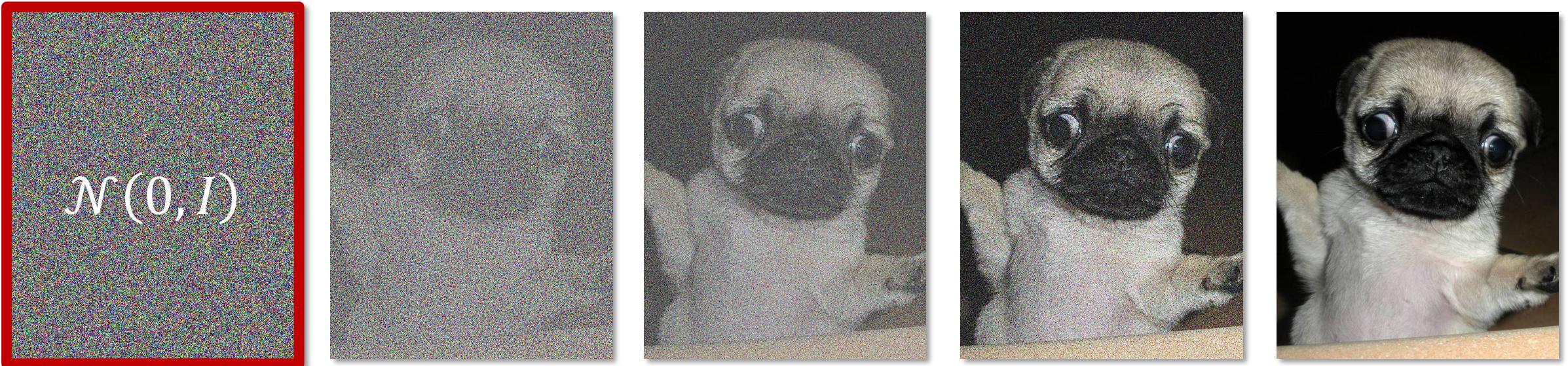
$$s_t := \arg \min_{\text{NN} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left\| \frac{1}{\sqrt{1-e^{-2t}}} g^{(i)} + \text{NN}\left(e^{-t} X^{(i)} + \sqrt{1-e^{-2t}} g^{(i)}\right) \right\|^2$$

$g^{(i)} \sim \mathcal{N}(0, I)$

Assumption: $\mathbb{E}_{q_t}[\|s_t(X_t) - \nabla \ln q_t(X_t)\|^2] \leq \epsilon_{sc}^2 \quad \forall t$

Diffusion models

To sample fresh images, run reverse process with **Gaussian initialization**



$t = 0$

$t = T$

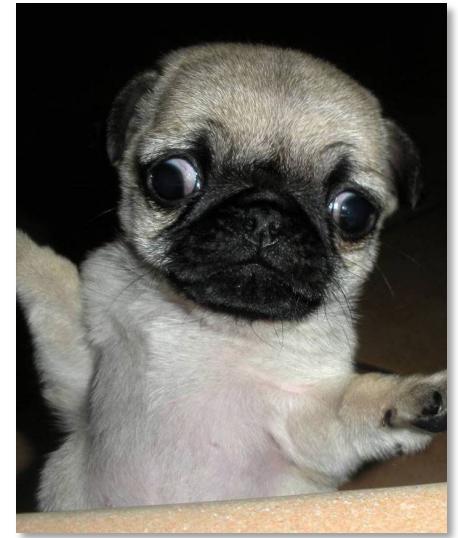
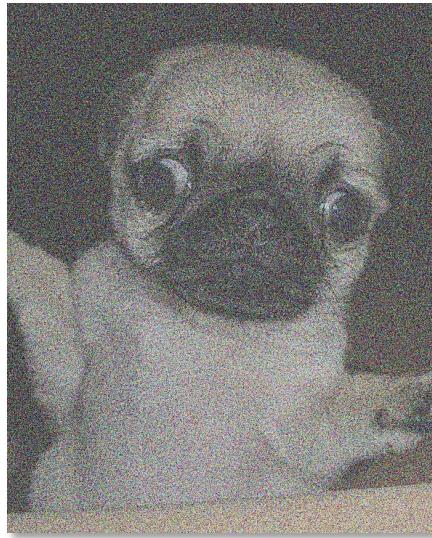
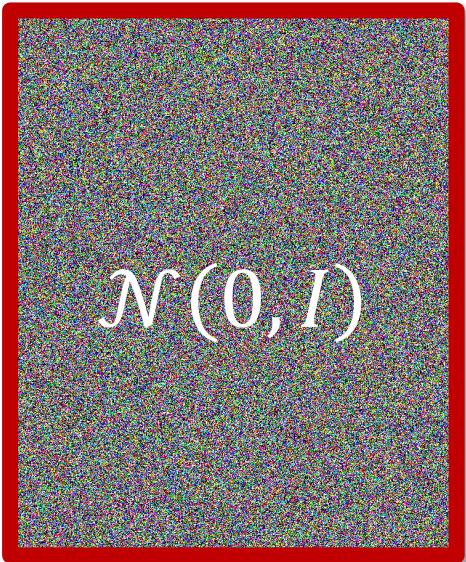
$$dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$$

$X_0^\leftarrow \sim q_T$ (forward process at time T)

“score function”
How to estimate it?

Diffusion models

To sample fresh images, run reverse process with **Gaussian initialization**



$t = 0$



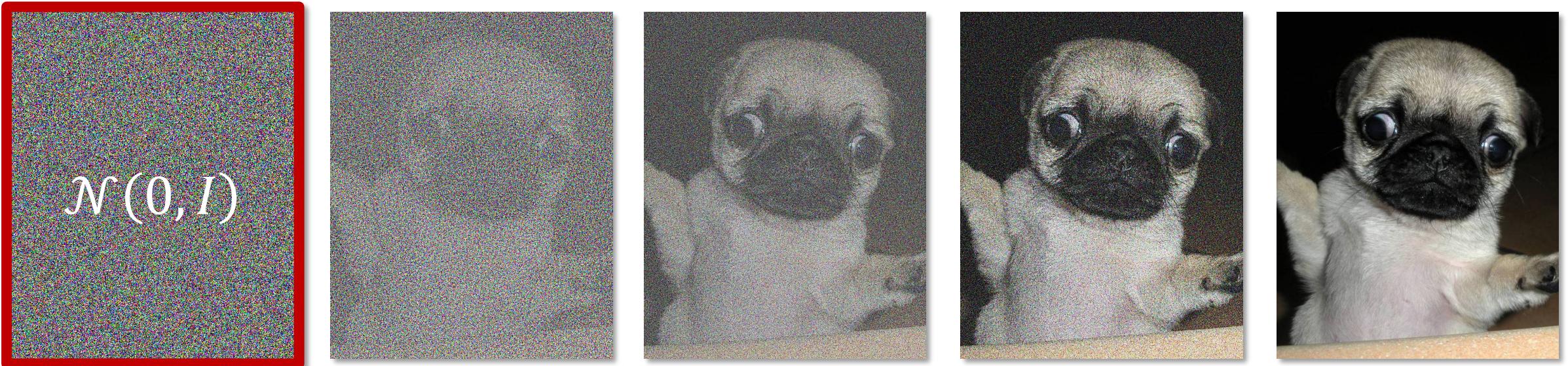
$t = T$

$$\begin{aligned} dX_t^\leftarrow &= \{X_t^\leftarrow + 2s_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t \\ X_0^\leftarrow &\sim q_T \text{ (forward process at time } T\text{)} \end{aligned}$$

can't run in continuous time, need to **discretize**

Diffusion models

To sample fresh images, run reverse process with **Gaussian initialization**



$t = 0$

$t = h, 2h, 3h, \dots$

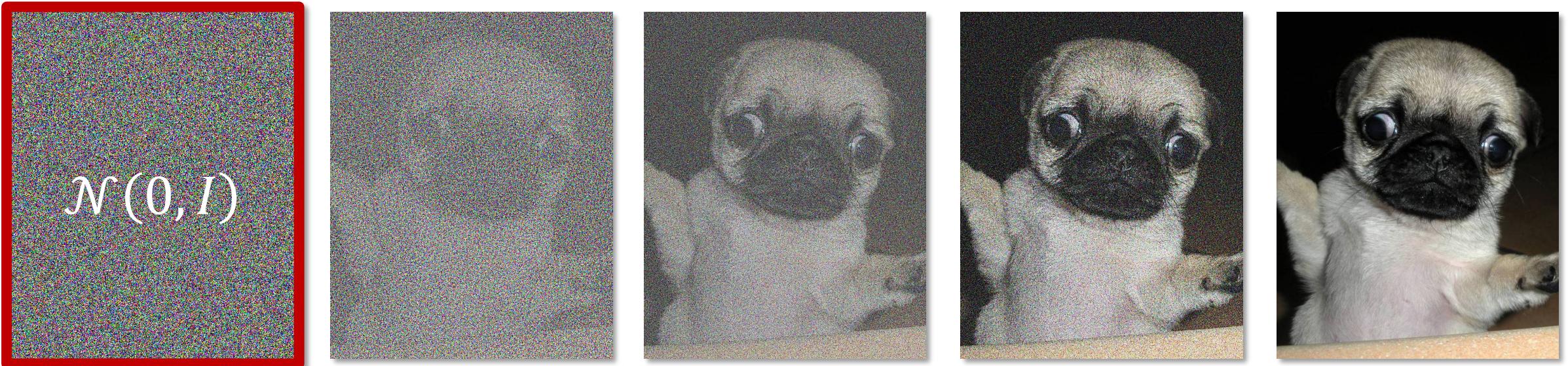
$t = T$

$$dX_t^\leftarrow = \{X_t^\leftarrow + 2s_{T-kh}(X_{kh}^\leftarrow)\} dt + \sqrt{2} dB_t \quad \text{for } t \in [kh, (k+1)h]$$

$X_0^\leftarrow \sim q_T$ (forward process at time T)

Diffusion models

To sample fresh images, run reverse process with **Gaussian initialization**



$$t = 0 \quad \text{---} \quad t = h, 2h, 3h, \dots \quad \text{---} \quad t = T$$

$$\begin{aligned} X_{(k+1)h}^{\leftarrow} &= e^h X_{kh}^{\leftarrow} + (e^h - 1) 2s_{T-kh}(X_{kh}^{\leftarrow}) + \mathcal{N}(0, (e^{2h} - 1)I) \\ X_0^{\leftarrow} &\sim q_T \text{ (forward process at time } T\text{)} \end{aligned}$$

Deferred proofs

1. Integral solution for the Ornstein-Uhlenbeck process

$$dX_t = -X_t + \sqrt{2}dB_t \implies X_t = e^{-t} \cdot X_0 + \mathcal{N}(0, (1 - e^{-2t})I)$$

2. Tweedie's formula

$$\mathbb{E}[x | \tilde{x}] = \tilde{x} + \sigma^2 \cdot \nabla \ln \tilde{p}(\tilde{x})$$

3. Integral solution for the discretized reverse process:

$$X_{(k+1)h}^\leftarrow = e^h X_{kh}^\leftarrow + (e^h - 1) 2s_{T-kh}(X_{kh}^\leftarrow) + \mathcal{N}(0, (e^{2h}-1)I)$$

4. Reverse-time SDE

$$dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$$

Solve linear SDE

$$dX_t = -X_t + \sqrt{2}dB_t$$

- $e^t(X_t + dX_t) = \sqrt{2}e^t dB_t$, which implies that

$$d(e^t X_t) = \sqrt{2}e^t dB_t$$

- Integrate from 0 to t on both sides:

$$e^t X_t - X_0 = \sqrt{2} \int_0^t e^s dB_s \quad \Rightarrow \quad X_t = e^{-t} X_0 + \boxed{\sqrt{2} \int_0^t e^{s-t} dB_s}$$

- By Itô's isometry,

$$\sqrt{2} \int_0^t e^{s-t} dB_s \equiv \mathcal{N} \left(0, 2 \int_0^t (e^{s-t})^2 ds \right) = \mathcal{N}(0, (1 - e^{-2t})I)$$

Proof of Tweedie's formula

$$\tilde{x} = x + e \text{ for } x \sim p \text{ and } e \sim \mathcal{N}(0, \sigma^2 I)$$

By Bayes' rule,

$$\mathbb{P}[x | \tilde{x}] = \frac{\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\tilde{x}-x)^2}{2\sigma^2}\right) \cdot p(x)}{\tilde{p}(\tilde{x})}$$

so

$$\mathbb{E}\left[\frac{x - \tilde{x}}{\sigma^2} \mid \tilde{x}\right] = \tilde{p}(\tilde{x})^{-1} \underbrace{\frac{1}{\sigma\sqrt{2\pi}} \int \exp\left(-\frac{(\tilde{x}-x)^2}{2\sigma^2}\right) \cdot \frac{x - \tilde{x}}{\sigma^2} \cdot p(x) dx}_{\nabla \tilde{p}(\tilde{x})}$$

Observe that

$$\tilde{p}(\tilde{x}) = \int \exp\left(-\frac{(\tilde{x}-x)^2}{2\sigma^2}\right) \cdot p(x) dx$$

Proof of Tweedie's formula

$$\tilde{x} = x + e \text{ for } x \sim p \text{ and } e \sim \mathcal{N}(0, \sigma^2 I)$$

By Bayes' rule,

$$\mathbb{P}[x \mid \tilde{x}] = \frac{\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\tilde{x}-x)^2}{2\sigma^2}\right) \cdot p(x)}{\tilde{p}(\tilde{x})}$$

so

$$\mathbb{E}\left[\frac{x - \tilde{x}}{\sigma^2} \mid \tilde{x}\right] = \frac{\nabla \tilde{p}(\tilde{x})}{\tilde{p}(\tilde{x})} = \nabla \ln \tilde{p}(\tilde{x})$$

■

Reverse-time SDE

- Forward process: $dX_t = -X_t + \sqrt{2}dB_t$

$$\frac{\partial}{\partial t} q_t(\mathbf{x}) = \nabla \cdot (\mathbf{x} q_t(\mathbf{x})) + \Delta q_t$$

- Reverse process: $dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$

$$\frac{\partial}{\partial t} q_t^\leftarrow(\mathbf{x}) = -\nabla \cdot ((\mathbf{x} + 2\nabla \ln q_t^\leftarrow(\mathbf{x})) q_t^\leftarrow(\mathbf{x})) + \Delta q_t^\leftarrow$$

Fokker-Planck equation. Let $\{x_t\}_{t \geq 0}$ follows $d\mathbf{x}_t = \boldsymbol{\mu}_t(x_t)dt + \boldsymbol{\sigma}_t d\mathbf{B}_t$ and $\mathbf{x}_0 \sim \pi_0$. Then for all $t \geq 0$, denoting the law of x_t by π_t , we have

$$\frac{\partial}{\partial t} \pi_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\mu}_t(\mathbf{x}) \pi_t(\mathbf{x})) + \frac{1}{2} \sum_{i,j \in [d]} \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (\boldsymbol{\sigma}_t \boldsymbol{\sigma}_t^\top(\mathbf{x})_{ij} \pi_t(\mathbf{x}))$$

Reverse-time SDE

- Forward process: $dX_t = -X_t + \sqrt{2}dB_t$

$$\frac{\partial}{\partial t} q_t(\mathbf{x}) = \nabla \cdot (\mathbf{x} q_t(\mathbf{x})) + \Delta q_t$$

- Reverse process: $dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$

$$\begin{aligned}\frac{\partial}{\partial t} q_t^\leftarrow(\mathbf{x}) &= -\nabla \cdot ((\mathbf{x} + 2\nabla \ln q_t^\leftarrow(\mathbf{x})) q_t^\leftarrow(\mathbf{x})) + \Delta q_t^\leftarrow \\ &= -\nabla \cdot (\mathbf{x} q_t^\leftarrow(\mathbf{x})) - \Delta q_t^\leftarrow\end{aligned}$$

$$\frac{\partial}{\partial t} q_t(\mathbf{x}) = \nabla \cdot (\mathbf{x} q_t(\mathbf{x})) + \Delta q_t \quad \xleftrightarrow{t \Rightarrow T-t} \quad \frac{\partial}{\partial t} q_t^\leftarrow(\mathbf{x}) = -\nabla \cdot (\mathbf{x} q_t^\leftarrow(\mathbf{x})) - \Delta q_t^\leftarrow$$

Discretization analysis

- Let p_T be the law of the output of the algorithm after $N = T/h$ steps
- Our goal is to show that $\text{TV}(p_T, q)$ is small

We need to bound three sources of error:

- 1) the initialization of the algorithm at pure Gaussian noise rather than at q_T
- 2) the estimation of the score function
- 3) the discretization of the SDE with step size $h > 0$

Assumption I (L^2 -accurate score estimate): For all $t = 0, h, 2h, \dots, T$, the score estimate $s_t(\cdot)$ satisfies

$$\mathbb{E}_{q_t}[\|s_t(X_t) - \nabla \ln q_t(X_t)\|^2] \leq \epsilon_{\text{sc}}^2$$

Assumption II (Smoothness): For all $t \geq 0$, $\nabla \ln q_t(\cdot)$ is L -Lipschitz

Assumption III (Bounded second moment): $\mathfrak{m}_2^2 := \mathbb{E}_q[\|x\|^2] < \infty$

Convergence guarantee

Theorem (Chen et al. '23). Under **Assumptions I-III**, if p_T is the law of the output of the algorithm after $N = T/h$ iterations with step size h ,

$$\text{TV}(p_T, q) \lesssim \underbrace{\sqrt{\text{KL}(q\|\gamma)} \cdot \exp(-T)}_{\text{error from initializing at } \gamma = \mathcal{N}(0, I) \text{ instead of } q_T} + \underbrace{(L\sqrt{dh} + L\mathfrak{m}_2 h)\sqrt{T}}_{\text{discretization error}} + \underbrace{\epsilon_{sc}\sqrt{T}}_{\text{score error}}$$

- Suppose $\text{KL}(q\|\gamma) \leq \text{poly}(d)$ and $\mathfrak{m}_2^2 \lesssim d$. Choose $T = \log(\text{KL}(q\|\gamma)/\epsilon)$ and $h = \frac{\epsilon^2}{L^2 d}$ gives
$$\text{TV}(p_T, q) = \tilde{\mathcal{O}}(\epsilon + \epsilon_{sc})$$
 with $N = \tilde{\mathcal{O}}(L^2 d / \epsilon^2)$ steps
- To get ϵ -close in TV, it suffices to estimate the score function to within $\epsilon_{sc} \leq \tilde{\mathcal{O}}(\epsilon)$ accuracy

Bound initialization error

1. Forward process converges exponentially quickly to Gaussian

$$\text{KL}(q_T \parallel \gamma) \leq \exp(-2T) \text{KL}(q_0 \parallel \gamma) = \exp(-2T) \text{KL}(q \parallel \gamma)$$

2. Running reverse SDE starting from Gaussian vs. from q_T cannot increase distance between them (**data processing inequality**)

$$\text{KL}(q \parallel p_T) = \text{KL}(\text{reverse}(q_T) \parallel \text{reverse}(\gamma)) \leq \text{KL}(q_T \parallel \gamma)$$

3. **Pinsker's inequality:** $\text{TV}(p_T, q) \leq \sqrt{\frac{1}{2} \text{KL}(q \parallel p_T)}$

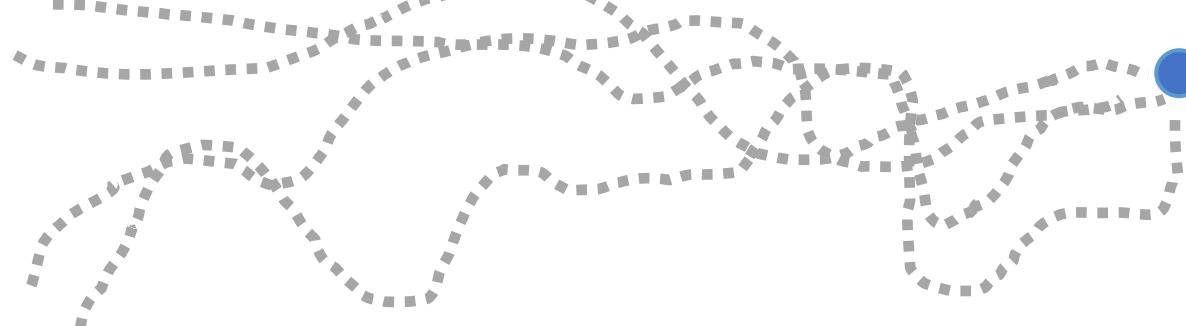
Discretization argument

Consider the **ideal reverse process** (continuous, perfect score) and **algorithm** (discrete, estimated score), with **both initialized at q_T**

- $dX_t^\leftarrow = \{X_t^\leftarrow + 2\nabla \ln q_{T-t}(X_t^\leftarrow)\} dt + \sqrt{2} dB_t$
- $dX_t^\leftarrow = \{X_t^\leftarrow + 2s_{T-kh}(X_{kh}^\leftarrow)\} dt + \sqrt{2} dB_t \quad X_{(k+1)h}^\leftarrow = e^h X_{kh}^\leftarrow + (e^h - 1) 2s_{T-kh}(X_{kh}^\leftarrow) + \mathcal{N}(0, e^{2h} - 1)$

To control the distance between these processes, we use **Girsanov's theorem**, a powerful tool for bounding distance between laws of processes driven by the same Brownian motion

Intuition: distribution over **last iterate** of reverse process is hard to characterize, but distribution over **trajectory** just given by a bunch of Gaussian samples



Girsanov's theorem

Consider the SDE's

$$dX_t = b_t dt + \sqrt{2} dB_t$$

$$dX_t = b'_t dt + \sqrt{2} dB_t$$

with the same initial distribution.

Let Q_T, P_T denote the laws of the **trajectories** over time T , i.e. $\mathcal{C}([0, T]; \mathbb{R}^d)$. Then

$$\frac{dQ_T}{dP_T} = \exp\left(\frac{1}{\sqrt{2}} \int_0^T (b_t - b'_t) dB_t - \frac{1}{4} \int_0^T \|b_t - b'_t\|^2 dt\right)$$

where B_t is a standard Brownian motion w.r.t. P_T

Heuristic proof of Girsanov's theorem

For the SDE's:

$$dx_t = \mathbf{b}_t dt + \sqrt{2} dB_t \quad (\text{in } Q\text{-world})$$

$$dx_t = \mathbf{b}'_t dt + \sqrt{2} dB_t \quad (\text{in } P\text{-world})$$

consider the discrete-time approximation:

$$\hat{x}_{(k+1)h} = \hat{x}_{kh} + h\mathbf{b}_{kh}(\hat{x}_{kh}) + \sqrt{2h}g_{kh} \quad (\text{in } Q\text{-world})$$

$$\hat{x}_{(k+1)h} = \hat{x}_{kh} + h\mathbf{b}'_{kh}(\hat{x}_{kh}) + \sqrt{2h}g_{kh} \quad (\text{in } P\text{-world})$$

For a trajectory $(\hat{x}_0, \hat{x}_h, \hat{x}_{2h}, \dots, \hat{x}_{Nh})$, what are the likelihoods in Q -world and P -world?

- Q -world: $L_Q \propto \prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} \|\hat{x}_{(k+1)h} - \hat{x}_{kh} - h\mathbf{b}_{kh}(\hat{x}_{kh})\|^2\right)$
- P -world: $L_P \propto \prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} \|\hat{x}_{(k+1)h} - \hat{x}_{kh} - h\mathbf{b}'_{kh}(\hat{x}_{kh})\|^2\right)$

Heuristic proof of Girsanov's theorem

For a trajectory $(\hat{x}_0, \hat{x}_h, \hat{x}_{2h}, \dots, \hat{x}_{Nh})$, what are the **likelihoods** in Q -world and P -world?

$$\begin{aligned}
 \frac{L_Q}{L_P} &= \frac{\prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} \|\hat{x}_{(k+1)h} - \hat{x}_{kh} - h b_{kh}(\hat{x}_{kh})\|^2\right)}{\prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} \|\hat{x}_{(k+1)h} - \hat{x}_{kh} - h b'_{kh}(\hat{x}_{kh})\|^2\right)} \\
 &= \prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} (h^2 \|b_{kh}(\hat{x}_{kh})\|^2 - h^2 \|b'_{kh}(\hat{x}_{kh})\|^2 - 2h \underbrace{\langle \hat{x}_{(k+1)h} - \hat{x}_{kh}, b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh}) \rangle}_{hb'_{kh}(\hat{x}_{kh}) + \sqrt{2h}g_{kh} \text{ in } P\text{-world}})\right) \\
 &= \prod_{k=0}^{N-1} \exp\left(-\frac{1}{4h} (h^2 \|b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh})\|^2 - 2\sqrt{2h} \langle \sqrt{h}g_{kh}, b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh}) \rangle)\right) \\
 &= \exp\left(-\frac{1}{4} \sum_{k=0}^{N-1} h \|b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh})\|^2 + \frac{1}{\sqrt{2}} \sum_{k=0}^{N-1} \langle \sqrt{h}g_{kh}, b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh}) \rangle\right)
 \end{aligned}$$

Heuristic proof of Girsanov's theorem

For a trajectory $(\hat{x}_0, \hat{x}_h, \hat{x}_{2h}, \dots, \hat{x}_{Nh})$, what are the **likelihoods** in Q -world and P -world?

$$\frac{L_Q}{L_P} = \exp \left(-\frac{1}{4} \sum_{k=0}^{N-1} h \| b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh}) \|^2 + \frac{1}{\sqrt{2}} \sum_{k=0}^{N-1} \underbrace{\langle \sqrt{h} g_{kh}, b_{kh}(\hat{x}_{kh}) - b'_{kh}(\hat{x}_{kh}) \rangle}_{\text{"dB}_{kh}\text{"}} \right)$$

$$\xrightarrow{(h \rightarrow 0)} \exp \left(-\frac{1}{4} \int_0^T \| b_t - b'_t \|^2 dt + \frac{1}{\sqrt{2}} \underbrace{\int_0^T (b_t - b'_t) dB_t}_{\text{martingale in } P\text{-world}} \right) = \frac{dQ_T}{dP_T}$$

KL divergence bound from Girsanov's theorem

$$\text{KL}(Q_T \| P_T) = \mathbb{E}_Q \left[-\frac{1}{4} \int_0^T \|b_t - b'_t\|^2 dt + \frac{1}{\sqrt{2}} \int_0^T (b_t - b'_t) dB_t \right]$$

BM in P -world

By Girsanov's theorem, we can relate the Brownian motion in P -world to Q -world:

$$b_t dt + \sqrt{2} d\tilde{B}_t = b'_t dt + \sqrt{2} dB_t \quad \Rightarrow \quad d\tilde{B}_t = \frac{1}{\sqrt{2}} (b'_t - b_t) dt + dB_t$$

a new BM in Q -world

$$\begin{aligned} \text{KL}(Q_T \| P_T) &= \mathbb{E}_Q \left[-\frac{1}{4} \int_0^T \|b_t - b'_t\|^2 dt + \frac{1}{\sqrt{2}} \int_0^T \left\langle b_t - b'_t, \frac{1}{\sqrt{2}} (b_t - b'_t) dt + d\tilde{B}_t \right\rangle \right] \\ &= \frac{1}{4} \mathbb{E}_Q \left[\int_0^T \|b_t - b'_t\|^2 dt \right] + \frac{1}{\sqrt{2}} \mathbb{E}_Q [(b_t - b'_t) d\tilde{B}_t] \end{aligned}$$

Warmup: Girsanov analysis for Langevin dynamics

Suppose the target distribution q is strongly log-concave: $\alpha \cdot I \preccurlyeq -\nabla^2 \ln q \preccurlyeq L \cdot I$

ULA: $d\hat{x}_t = -\nabla \ln \hat{q}(\hat{x}_{kh}) dt + \sqrt{2} dB_t \quad \text{for } t \in [kh, (k+1)h]$

LD: $dx_t = -\nabla \ln q(x_t) dt + \sqrt{2} dB_t$

$$\text{KL}(\text{LD} \parallel \text{ULA}) = \frac{1}{4} \mathbb{E}_{\text{LD}} \left[\sum_{k=0}^{T/h-1} \int_{kh}^{(k+1)h} \|\nabla \ln q(x_t) - \nabla \ln \hat{q}(x_{kh})\|^2 dt \right]$$

$$\leq \frac{L^2}{4} \mathbb{E}_{\text{LD}} \left[\sum_{k=0}^{T/h-1} \int_{kh}^{(k+1)h} \|x_t - x_{kh}\|^2 dt \right]$$

$$\begin{aligned} \|x_t - x_{kh}\|^2 &= \left\| \int_0^{t-kh} \nabla \ln q(x_{kh+s}) ds + \sqrt{2}(B_t - B_{kh}) \right\|^2 \\ &\leq 2h \int_0^{t-kh} \|\nabla \ln q(x_{kh+s})\|^2 ds + 4\|B_t - B_{kh}\|^2 \end{aligned}$$

Warmup: Girsanov analysis for Langevin dynamics

Suppose the target distribution q is strongly log-concave: $\alpha \cdot I \leq -\nabla^2 \ln q \leq L \cdot I$

ULA: $d\hat{x}_t = -\nabla \ln \hat{q}(\hat{x}_{kh}) dt + \sqrt{2} dB_t \quad \text{for } t \in [kh, (k+1)h]$

LD: $dx_t = -\nabla \ln q(x_t) dt + \sqrt{2} dB_t$

$$\begin{aligned} \text{KL}(\text{LD} \parallel \text{ULA}) &= \frac{1}{4} \mathbb{E}_{\text{LD}} \left[\sum_{k=0}^{T/h-1} \int_{kh}^{(k+1)h} \|\nabla \ln q(x_t) - \nabla \ln \hat{q}(x_{kh})\|^2 dt \right] \\ &\leq \frac{L^2}{4} \mathbb{E}_{\text{LD}} \left[\sum_{k=0}^{T/h-1} \int_{kh}^{(k+1)h} \|x_t - x_{kh}\|^2 dt \right] \\ &\lesssim L^2 h^2 \int_0^T \mathbb{E}_{\text{LD}} [\|\nabla \ln q(x_t)\|^2] dt + L^2 d h T \end{aligned}$$

$$\mathbb{E}_{\text{LD}} [\|\nabla \ln q(x_t)\|^2] \lesssim Ld + \frac{L^2}{\alpha} \text{KL}(\text{Law}(x_0) \parallel q)$$

Warmup: Girsanov analysis for Langevin dynamics

Suppose the target distribution q is strongly log-concave: $\alpha \cdot I \leq -\nabla^2 \ln q \leq L \cdot I$

ULA: $d\hat{x}_t = -\nabla \ln \hat{q}(\hat{x}_{kh}) dt + \sqrt{2} dB_t \quad \text{for } t \in [kh, (k+1)h]$

LD: $dx_t = -\nabla \ln q(x_t) dt + \sqrt{2} dB_t$

$$\text{KL}(\text{LD} \parallel \text{ULA}) \lesssim L^2 h^2 T \left(Ld + \frac{L^2}{\alpha} \text{KL}(\text{Law}(x_0) \parallel q) \right) + L^2 dhT$$

- For sufficiently small ϵ , if we take $h = \frac{\epsilon^2}{L^2 d T}$, then

$$\text{KL}(\text{Law}(x_T) \parallel \text{Law}(\hat{x}_T)) \leq \text{KL}(\text{LD} \parallel \text{ULA}) \lesssim \epsilon^2$$

- Since q satisfies α -LSI, if we take $T = \frac{1}{\alpha} \log(\text{KL}(\text{Law}(x_0) \parallel q) / \epsilon)$, then $\text{KL}(\text{Law}(x_T) \parallel q) \leq \epsilon^2$
- Finally, by Pinsker and triangle inequality for TV distance, we get that

$$\text{TV}(\text{Law}(\hat{x}_T), q) \lesssim \sqrt{\text{KL}(\text{Law}(x_T) \parallel \text{Law}(\hat{x}_T))} + \sqrt{\text{KL}(\text{Law}(x_T) \parallel q)} \leq \epsilon$$

Warmup: Girsanov analysis for Langevin dynamics: bound

$$\mathbb{E}_{\text{LD}}[\|\nabla \ln q(x_t)\|^2]$$

Suppose the target distribution q is strongly log-concave: $\alpha \cdot I \preccurlyeq -\nabla^2 \ln q \preccurlyeq L \cdot I$

By the definition of W_2 ,

$$\mathbb{E}_{\text{LD}}[\|\nabla \ln q(x_t)\|^2] \lesssim \mathbb{E}_q[\|\nabla \ln q(x)\|^2] + L^2 W_2^2(\text{Law}(x_t), q)$$

- For the first term,

$$\mathbb{E}_q[\|\nabla \ln q(x)\|^2] = \int \langle \nabla q, \nabla \ln q \rangle dx \stackrel{\text{i.b.p.}}{=} - \int (\Delta \ln q) q(x) dx \leq Ld$$

- For the second term,

$$L^2 W_2^2(\text{Law}(x_t), q) \leq \frac{L^2}{\alpha} \text{KL}(\text{Law}(x_t) \| q) \leq \frac{L^2}{\alpha} \text{KL}(\text{Law}(x_0) \| q)$$

Talagrand's T₂ inequality data processing inequality

Girsanov analysis for diffusion model

ALG: $dX_t = (\hat{X}_t + 2s_{T-kh}(\hat{X}_{kh}))dt + \sqrt{2}dB_t \quad \text{for } t \in [kh, (k+1)h]$

TRUE: $dX_t = (X_t + 2\nabla \ln q_{T-t}(X_t))dt + \sqrt{2}dB_t$

Following the same argument,

$$\text{KL}(\text{TRUE} \parallel \text{ALG}) = \mathbb{E}_{\text{TRUE}} \left[\sum_{k=0}^{T/h-1} \int_{kh}^{(k+1)h} \|\nabla \ln q_{T-t}(X_t) - s_{T-kh}(X_{kh})\|^2 dt \right]$$

Girsanov analysis for diffusion model

We can decompose $\mathbb{E}_{\text{TRUE}}[\|\nabla \ln q_{T-t}(X_t) - \nabla \ln s_{T-kh}(X_{kh})\|^2]$ into:

1. $\mathbb{E}[\|\nabla \ln q_{T-kh}(X_{kh}) - s_{T-kh}(X_{kh})\|^2] \leq \epsilon_{sc}^2$

score estimation error, bounded by assumption

2. $\mathbb{E}[\|\nabla \ln q_{T-t}(X_{kh}) - \nabla \ln q_{T-kh}(X_{kh})\|^2]$

bounded because score doesn't change much over short period

3. $\mathbb{E}[\|\nabla \ln q_{T-t}(X_t) - \nabla \ln q_{T-t}(X_{kh})\|^2]$

bounded because score is Lipschitz, and process doesn't move too much over short time

Space discretization

- By Lipschitzness of the score function,

$$\mathbb{E}[\|\nabla \ln q_{T-t}(X_t) - \nabla \ln q_{T-t}(X_{kh})\|^2] \lesssim L^2 \mathbb{E}[\|X_t - X_{kh}\|^2]$$

- The joint distribution of $(X_t, X_{kh}) \equiv \left(X_t, e^{-(t-kh)} X_t + \mathcal{N}(0, (1 - e^{-2(t-kh)})I) \right)$

$$\mathbb{E}[\|X_t - X_{kh}\|^2] \lesssim (1 - e^{-(t-kh)})^2 \mathbb{E}[\|X_t\|^2] + (1 - e^{-2(t-kh)})d$$

- $X_t \equiv e^{-(T-t)}X + \mathcal{N}(0, (1 - e^{-2(T-t)})I)$ for $X \sim q$

$$\mathbb{E}[\|X_t\|^2] \lesssim e^{-2(T-t)} \mathbb{E}_q[\|X\|^2] + (1 - e^{-2(T-t)})d \leq \mathfrak{m}_2^2 + d$$

- Therefore,

$$\mathbb{E}[\|\nabla \ln q_{T-t}(X_t) - \nabla \ln q_{T-t}(X_{kh})\|^2] \leq L^2 h^2 \mathfrak{m}_2^2 + L^2 h d$$

Time discretization

$$\mathbb{E}[\|\nabla \ln q_{\textcolor{red}{T-t}}(\textcolor{blue}{X}_{kh}) - \nabla \ln q_{T-kh}(\textcolor{blue}{X}_{kh})\|^2]$$

- Consider $p \propto \exp(-V)$, $\nabla^2 V \geq L \cdot I$, and $p' = p \star \mathcal{N}(0, \sigma^2 I)$, i.e. Gaussian convolution
- Our goal: $\mathbb{E}_{p'}[\|\nabla \ln p - \nabla \ln p'\|^2]$
- Notice that

$$\nabla \ln p'(x) = -\mathbb{E}_{p_{x,\sigma}}[\nabla V(y)] \quad \text{where } p_{x,\sigma} = \text{Law}(y|y + \sigma g = x) \text{ for } y \sim p$$

- $\mathbb{E}_{x \sim p'}[\|\nabla \ln p - \nabla \ln p'\|^2] = \mathbb{E}_{x \sim p'} \left[\left\| \mathbb{E}_{y \sim p_{x,\sigma}}[\nabla V(y) - \nabla V(x)] \right\|^2 \right] \leq L^2 \mathbb{E}_{x \sim p'} \mathbb{E}_{y \sim p_{x,\sigma}}[\|x - y\|^2]$
- $\text{Law}(x, y) \equiv \text{Law}(\tilde{y} + \sigma g, \tilde{y})$ for $\tilde{y} \sim p$
- Thus, $\mathbb{E}_{x \sim p'}[\|\nabla \ln p - \nabla \ln p'\|^2] \leq L^2 \sigma^2 d$

See Lemma C.12 in (Lee-Lu-Tan '22) for the full proof

Convergence guarantee

Assumption I (L^2 -accurate score estimate): For all $t = 0, h, 2h, \dots, T$, the score estimate $s_t(\cdot)$ satisfies

$$\mathbb{E}_{q_t}[\|s_t(X_t) - \nabla \ln q_t(X_t)\|^2] \leq \epsilon_{sc}^2$$

Assumption II (Smoothness): For all $t \geq 0$, $\nabla \ln q_t(\cdot)$ is L -Lipschitz

“early stopping” (Chen-Lee-Lu ’23)

Assumption III (Bounded second moment): $m_2^2 := \mathbb{E}_q[\|x\|^2] < \infty$

Theorem (Chen et al. ’23). Under Assumptions I-III, if p_T is the law of the output of the algorithm after $N = T/h$ iterations with step size h ,

$$TV(p_T, q) \lesssim \underbrace{\sqrt{KL(q\|\gamma)} \cdot \exp(-T)}_{\text{initialization error}} + \underbrace{(L\sqrt{dh} + Lm_2h)\sqrt{T}}_{\text{discretization error}} + \underbrace{\epsilon_{sc}\sqrt{T}}_{\text{score error}}$$

- Choose $T = \log(KL(q\|\gamma)/\epsilon)$ and $h = \frac{\epsilon^2}{L^2 d}$ gives $TV(p_T, q) = \tilde{\mathcal{O}}(\epsilon + \epsilon_{sc})$ with $N = \tilde{\mathcal{O}}(L^2 d/\epsilon^2)$ steps

Using ODE flow + Langevin corrector can achieve \sqrt{d} steps (Chen-Chewi-Lee-Li-Lu-Salim ’23)

SOTA DDPM convergence bound: $N = d/\epsilon$ (Li-Yan ’25)

Stochastic localization

Let π_0 be a probability measure over \mathbb{R}^d with mean \mathbf{m}_0

Define a stochastic process $\{\mathbf{c}_t\}_{t \geq 0}$ as follows:

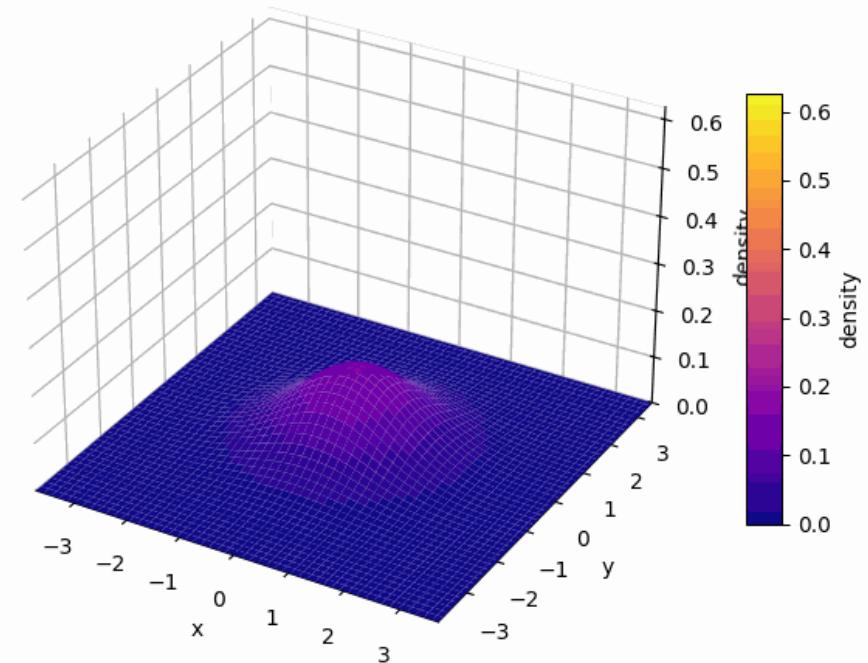
$$\mathbf{c}_0 = 0, \quad d\mathbf{c}_t = \mathbf{m}_t dt + dB_t$$

where $\mathbf{m}_t := \mathbb{E}_{\pi_t}[\mathbf{x}]$ and $\pi_t(\mathbf{x}) \propto \exp\left(\langle \mathbf{c}_t, \mathbf{x} \rangle - \frac{t}{2} \|\mathbf{x}\|^2\right) \pi_0(\mathbf{x})$

We call the random induced measures $\{\pi_t\}_{t \geq 0}$ the stochastic localization

- As $t \rightarrow \infty$, $\pi_t \rightarrow \delta_{\mathbf{x}}$ i.e. π_t localizes towards a **delta-measure** at son
- $\{\pi_t\}_{t \geq 0}$ is measure-valued martingale i.e. $\mathbb{E}[\pi_t(\mathbf{x})] = \pi_0(x)$ for all

Stochastic Localization ($t=0.00$)



SL \iff Diffusion model

Consider the reverse process of diffusion model (re-parameterized):

$$t \leftrightarrow \frac{1}{2} \log \frac{t+1}{t}$$

$$dX_t^\leftarrow = \left\{ \frac{X_t^\leftarrow}{2t(t+1)} + \frac{1}{t(t+1)} \nabla \ln \pi_t^\leftarrow(X_t^\leftarrow) \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

with $X_0^\leftarrow \sim \mathcal{N}(0, I)$

Then the processes $\{\mathbf{c}_t\}_{t \geq 0}$ and $\{X_t^\leftarrow\}_{t \geq 0}$ satisfy $\sqrt{t(t+1)}X_t^\leftarrow = \mathbf{c}_t$

- $d\mathbf{c}_t = \mathbf{m}_t dt + dB_t$ with the change of variable $\sqrt{t(t+1)}X_t^\leftarrow = \mathbf{c}_t$ gives

$$dX_t^\leftarrow = -\frac{2t+1}{2t(t+1)} X_t^\leftarrow dt + \frac{1}{\sqrt{t(t+1)}} \mathbf{m}_t dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

SL \iff Diffusion model

Consider the reverse process of diffusion model (re-parameterized):

$$t \leftrightarrow \frac{1}{2} \log \frac{t+1}{t}$$

$$dX_t^\leftarrow = \left\{ \frac{X_t^\leftarrow}{2t(t+1)} + \frac{1}{t(t+1)} \nabla \ln \pi_t^\leftarrow(X_t^\leftarrow) \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

with $X_0^\leftarrow \sim \mathcal{N}(0, I)$

Then the processes $\{\mathbf{c}_t\}_{t \geq 0}$ and $\{X_t^\leftarrow\}_{t \geq 0}$ satisfy $\sqrt{t(t+1)} X_t^\leftarrow = \mathbf{c}_t$

- $X_t^\leftarrow = X_{t'}$ for $t' = \frac{1}{2} \log \frac{t+1}{t}$ in the forward process, which distributed as $e^{-t'} X_0 + (1 - e^{-2t'}) g$
- By Tweedie's formula, $\nabla \ln \pi_t^\leftarrow(X_t^\leftarrow) = \sqrt{t(t+1)} \mathbb{E}_{\pi_t^\leftarrow}[X_0] - (t+1) X_t^\leftarrow$

$$dX_t^\leftarrow = \left\{ -\frac{2t+1}{2t(t+1)} X_t^\leftarrow + \frac{1}{\sqrt{t(t+1)}} \mathbb{E}_{\pi_t^\leftarrow}[X_0] \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

SL \iff Diffusion model

Consider the reverse process of diffusion model (re-parameterized):

$$t \leftrightarrow \frac{1}{2} \log \frac{t+1}{t}$$

$$dX_t^\leftarrow = \left\{ \frac{X_t^\leftarrow}{2t(t+1)} + \frac{1}{t(t+1)} \nabla \ln \pi_t^\leftarrow(X_t^\leftarrow) \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

with $X_0^\leftarrow \sim \mathcal{N}(0, I)$

Then the processes $\{\mathbf{c}_t\}_{t \geq 0}$ and $\{X_t^\leftarrow\}_{t \geq 0}$ satisfy $\sqrt{t(t+1)} X_t^\leftarrow = \mathbf{c}_t$

$$dX_t^\leftarrow = -\frac{2t+1}{2t(t+1)} X_t^\leftarrow dt + \frac{1}{\sqrt{t(t+1)}} \mathbf{m}_t dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

$$dX_t^\leftarrow = \left\{ -\frac{2t+1}{2t(t+1)} X_t^\leftarrow + \frac{1}{\sqrt{t(t+1)}} \mathbb{E}_{\pi_t^\leftarrow}[X_0] \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

SL \iff Diffusion model

Consider the reverse process of diffusion model (re-parameterized):

$$t \leftrightarrow \frac{1}{2} \log \frac{t+1}{t}$$

$$dX_t^\leftarrow = \left\{ \frac{X_t^\leftarrow}{2t(t+1)} + \frac{1}{t(t+1)} \nabla \ln \pi_t^\leftarrow(X_t^\leftarrow) \right\} dt + \frac{1}{\sqrt{t(t+1)}} dB_t$$

with $X_0^\leftarrow \sim \mathcal{N}(0, I)$

Then the processes $\{\mathbf{c}_t\}_{t \geq 0}$ and $\{X_t^\leftarrow\}_{t \geq 0}$ satisfy $\sqrt{t(t+1)} X_t^\leftarrow = \mathbf{c}_t$

- π_t^\leftarrow is the density of $e^{-t'} X_0 + (1 - e^{-2t'}) g = \sqrt{\frac{t}{t+1}} X_0 + \frac{1}{\sqrt{t+1}} g$

$$\begin{aligned} \pi_t^\leftarrow(X_0 | X_t^\leftarrow) &\propto \exp\left(-\frac{t+1}{2} \left\| \sqrt{t/t+1} X_0 - X_t^\leftarrow \right\|^2\right) \pi_0(X_0) \\ &\propto \exp\left(\sqrt{t(t+1)} \langle X_t^\leftarrow, X_0 \rangle - \frac{t}{2} \|X_0\|^2\right) \pi_0(X_0) = \underbrace{\exp\left(\langle \mathbf{c}_t, X_0 \rangle - \frac{t}{2} \|X_0\|^2\right)}_{\pi_t(X_0)} \pi_0(X_0) \end{aligned}$$

Provable score estimation

- El Alaoui-Montanari-Sellke '22; Celentano '24: Sampling from the Sherrington–Kirkpatrick model where $\pi_W(x) \propto \exp(-\beta x^\top W x)$
 $W \in \mathbb{R}^{d \times d}$ is a **random matrix** with i.i.d. Gaussian entries
- El Alaoui-Montanari-Sellke '23; Huang-Montanari-Pham '24; Huang-Mohanty-Rajaraman-Wu '24:
Sampling from the p -spin spherical spin glass
- Montanari-Wu '23; Montanari-Wu '24: Bayesian posterior sampling: observing $A = \frac{\beta}{d} \theta^\top \theta + W$ where $\theta \sim_u \{-1,1\}^d$, sample from the posterior distribution
$$P(\theta|A) \propto \exp\left(-\frac{\beta}{2} \theta^\top A \theta\right)$$
- Shah-Chen-Klivans '23; Chen-Kontonis-Shah '24; Gatmiry-Kelner-Lee '24: learning mixtures of Gaussians
- Chewi-Kalavasis-Mehrotra-Montasser '25: Statistical theory for score estimation (**and a comprehensive literature review on diffusion model theory**)