

[AI2613] Convergence of Langevin Diffusion, DDPM

June 4, 2024

1 Convergence of Langevin diffusion

In this section, we analyse the convergence rate of Langevin diffusion.

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a μ -strongly convex function¹. Recall that when analysing the gradient flow, we calculate $\frac{d\|X_t - Y_t\|^2}{dt}$, where X_t and Y_t come from the ODE $dX_t = -\nabla f(X_t) dt$ and $dY_t = -\nabla f(Y_t) dt$ and start from some $X_0, Y_0 \in \mathbb{R}^d$ respectively.²

For the Langevin diffusion, let $\{X_t\}$ and $\{Y_t\}$ be two processes generated by

$$\begin{cases} dX_t = -\nabla f(X_t) dt + dB_t \\ dY_t = -\nabla f(Y_t) dt + dB_t \end{cases} \quad (1)$$

Assume Y_0 is drawn from the stationary distribution $\pi \propto e^{-f}$ and the distribution of X_0 can be arbitrary. We use the Wasserstein metric to measure the distance between the distribution of X_t and Y_t .

We couple X_t and Y_t with the same Brownian motion. From Equation (1), we have

$$\frac{d(X_t - Y_t)}{dt} = \nabla f(Y_t) - \nabla f(X_t).$$

Consequently,

$$\begin{aligned} \frac{d\mathbb{E}[\|X_t - Y_t\|^2]}{dt} &= 2\mathbb{E}\left[\left\langle \frac{d(X_t - Y_t)}{dt}, X_t - Y_t \right\rangle\right] \\ &= 2\mathbb{E}[\langle \nabla f(Y_t) - \nabla f(X_t), X_t - Y_t \rangle] \\ &\stackrel{(\heartsuit)}{\leq} -2\mu\mathbb{E}[\|X_t - Y_t\|^2], \end{aligned} \quad (2)$$

where (\heartsuit) follows from

$$\begin{cases} f(y) - f(x) \geq \langle \nabla f(x), y - x \rangle + \frac{\mu}{2}\|x - y\|_2^2 \\ f(x) - f(y) \geq \langle \nabla f(y), x - y \rangle + \frac{\mu}{2}\|x - y\|_2^2 \end{cases}$$

which is due to the strong convexity of f . Equation (2) indicates the convergence of Langevin diffusion

$$\mathbb{E}[\|X_t - Y_t\|] \leq e^{-2\mu t} \mathbb{E}[\|X_0 - Y_0\|].$$

For the discretized Langevin algorithm, we left its analysis in homework.

2 DDPM

Note that the above analysis relies heavily on the strong log-concavity of π (or equivalently, the strong convexity of f), which does not hold for most distributions in practice. One recent popular method to sample a general

¹ A differentiable function f is μ -strongly convex if

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2}\|x - y\|_2^2.$$

² Unless otherwise stated, $\|\cdot\|$ refers to the Euclidean norm in this note.

For two distributions μ and ν , let ζ be the set of all couplings between μ and ν . The Wasserstein distance $W_2(\mu, \nu) = \inf_{\omega \in \zeta} \mathbb{E}_{(X,Y) \sim \omega} [\|X - Y\|^2]^{\frac{1}{2}}$.

distribution π is the so-called *denoising diffusion probabilistic modeling* (DDPM).

The idea of DDPM is to add Gaussian noise to an initial sample $X_0 \sim \pi$ until the distribution is close enough to a Gaussian. Then it starts with a sample from the Gaussian distribution and executes a reverse process to generate a sample from the target distribution.

To be specific, in the forward procedure, we run an Ornstein-Uhlenbeck process (OU process) starting from $X_0 \sim \pi$ and get $\{\bar{X}_t\}$. Recall that the OU process follows the stochastic differential equation $d\bar{X}_t = -\bar{X}_t dt + \sqrt{2} dB_t$. Therefore we have

$$\bar{X}_t = e^{-t} X_0 + \sqrt{1 - e^{-2t}} Z_t$$

where $Z_t \sim \mathcal{N}(0, 1)$ is a standard Gaussian random variable. Running this for a time of T , we get a series of random variables $\{\bar{X}_t\}$ and \bar{X}_T is close to a Gaussian random variable.

Let q_t be the law of \bar{X}_t . In the reversed process, we first sample \bar{X}_T from a Gaussian distribution and then calculate $\{\bar{X}_t\}$ according to the reversed equation of OU process:

$$d\bar{X}_t = \left(\bar{X}_t + 2\nabla \log q_{T-t}(\bar{X}_t) \right) dt + \sqrt{2} dB_t. \quad (3)$$

If law of \bar{X}_0 is $\mathcal{N}(0, 1)$, then the law of \bar{X}_0 will converge to π . In the following, we will prove that Equation (3) is indeed the reverse of the OU process.

2.1 The Reverse of OU Process

Recall that for a diffusion $dX_t = \mu(X_t) dt + \sigma(X_t) dB_t$, the generator \mathcal{L} satisfies

$$\mathcal{L}f(x) = \mu(x) \cdot \partial_x f + \frac{1}{2} \sigma^2(x) \cdot \partial_x^2 f \quad (4)$$

for any function f . Also,

$$\mathcal{L}^* f(x) = -\partial_x (\mu \cdot f) + \frac{1}{2} \partial_x^2 (\sigma^2 \cdot f). \quad (5)$$

Recall the Kolmogorov Forward equation $\frac{\partial P_t}{\partial t} = \mathcal{L}^* P_t$. This is equivalent to say

$$\forall s < t, \partial_t p[X_t = x | X_s = y] = \mathcal{L}^* p[X_t = x | X_s = y]. \quad (6)$$

With fixed y , for any $s > \tau$, let $f(x) = \delta(y - x)$ and $f_\tau(x) = Q_\tau f(x) = p[X_s = y | X_{s-\tau} = x]$. Then the Kolmogorov backward equation indicates that

$$\partial_\tau p[X_s = x | X_{s-\tau} = y] = \mathcal{L} p[X_s = x | X_{s-\tau} = y],$$

or equivalently,

$$\forall t < s, -\partial_t p[X_s = x | X_t = y] = \mathcal{L} p[X_s = x | X_t = y]. \quad (7)$$

Here we slightly abuse the notation and let $p[X_s = y | X_{s-\tau} = x]$ represent the density of X_s at y given $X_{s-\tau} = x$. Similarly define $p[X_t = x]$ and $p[X_s = y, X_t = x]$.

In the following part, we abbreviate $p[X_s = x_s, X_t = x_t]$ as $p(x_s, x_t)$ and $p[X_s = x_s | X_t = x_t]$ as $p(x_s | x_t)$ (here $s > t$). We have

$$\begin{aligned} -\partial_t p(x_s, x_t) &= -\partial_t (p(x_s | x_t) \cdot p(x_t)) \\ &= -p(x_t) \cdot \partial_t p(x_s | x_t) - p(x_s | x_t) \cdot \partial_t p(x_t). \end{aligned} \quad (8)$$

From Equations (5) and (6),

$$\begin{aligned} \partial_t p(x_t) &= \mathcal{L}^* p(x_t) \\ &= -\partial_{x_t} (\mu(x_t) \cdot p(x_t)) + \frac{1}{2} \cdot \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_t)). \end{aligned} \quad (9)$$

From Equations (4) and (7),

$$\begin{aligned} \partial_t (p(x_s | x_t)) &= -\mathcal{L} p(x_s | x_t) \\ &= -\mu(x_t) \cdot \partial_{x_t} p(x_s | x_t) - \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t). \end{aligned} \quad (10)$$

Plugging Equations (9) and (10) into Equation (8), we have

$$\begin{aligned} -\partial_t p(x_s, x_t) &= p(x_t) \cdot \left(\mu(x_t) \cdot \partial_{x_t} p(x_s | x_t) + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) \right) \\ &\quad + p(x_s | x_t) \cdot \left(\partial_{x_t} (\mu(x_t) \cdot p(x_t)) - \frac{1}{2} \cdot \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_t)) \right). \end{aligned} \quad (11)$$

Note that

$$\partial_{x_t} p(x_s | x_t) = \partial_{x_t} \left(\frac{p(x_s, x_t)}{p(x_t)} \right) = \frac{\partial_{x_t} p(x_s, x_t)}{p(x_t)} - \frac{p(x_s, x_t) \cdot \partial_{x_t} p(x_t)}{p^2(x_t)}$$

and

$$\partial_{x_t} (\mu(x_t) \cdot p(x_t)) = p(x_t) \cdot \partial_{x_t} (\mu(x_t)) + \mu(x_t) \cdot \partial_{x_t} (p(x_t)).$$

Then by direct calculation, we can further write Equation (11) as

$$\begin{aligned} -\partial_t p(x_s, x_t) &= \partial_{x_t} (p(x_s, x_t) \cdot \mu(x_t)) + \frac{1}{2} p(x_t) \sigma^2(x_t) \cdot \partial_{x_t}^2 p(x_s | x_t) \\ &\quad - \frac{1}{2} p(x_s | x_t) \cdot \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_t)) \end{aligned} \quad (12)$$

We have

$$\begin{aligned} \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_s, x_t)) &= \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_s | x_t) p(x_t)) \\ &= \frac{1}{2} \sigma^2(x_t) p(x_t) \cdot \partial_{x_t}^2 p(x_s | x_t) + \frac{1}{2} p(x_s | x_t) \cdot \partial_{x_t}^2 (\sigma^2(x_t) p(x_t)) \\ &\quad + \partial_{x_t} p(x_s | x_t) \partial_{x_t} (\sigma^2(x_t) p(x_t)). \end{aligned}$$

Plugging this into Equation (12),

$$\begin{aligned}
-\partial_t p(x_s, x_t) &= \partial_{x_t} (p(x_s, x_t) \cdot \mu(x_t)) + \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_s, x_t)) \\
&\quad - p(x_s | x_t) \cdot \partial_{x_t}^2 (\sigma^2(x_t) p(x_t)) - \partial_{x_t} p(x_s | x_t) \partial_{x_t} (\sigma^2(x_t) p(x_t)) \\
&= \partial_{x_t} (p(x_s, x_t) \cdot \mu(x_t)) + \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_s, x_t)) \\
&\quad - \partial_{x_t} [p(x_s | x_t) \partial_{x_t} (\sigma^2(x_t) p(x_t))] \\
&= \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_s, x_t)) + \partial_{x_t} \left[p(x_s, x_t) \cdot \left(\mu(x_t) - \frac{1}{p(x_t)} \partial_{x_t} (\sigma^2(x_t) p(x_t)) \right) \right].
\end{aligned}$$

Since $p(x_s, x_t) = p(x_t | x_s) p(x_s)$, we have

$$-\partial_t p(x_t | x_s) = \frac{1}{2} \partial_{x_t}^2 (\sigma^2(x_t) \cdot p(x_t | x_s)) + \partial_{x_t} \left[p(x_t | x_s) \cdot \left(\mu(x_t) - \frac{1}{p(x_t)} \partial_{x_t} (\sigma^2(x_t) p(x_t)) \right) \right].$$

In the OU process, we have $\mu(x_t) = -x_t$ and $\sigma(x_t) = \sqrt{2}$. Therefore,

$$\begin{aligned}
-\partial_t p(x_t | x_s) &= \partial_{x_t}^2 p(x_t | x_s) + \partial_{x_t} \left[p(x_t | x_s) \cdot \left(-x_t - \frac{2}{p(x_t)} \partial_{x_t} p(x_t) \right) \right] \\
&= \partial_{x_t}^2 p(x_t | x_s) + \partial_{x_t} [p(x_t | x_s) \cdot (-x_t - 2 \partial_{x_t} \log p(x_t))].
\end{aligned}$$

Let $\tau = T - t$. Note that $p(x_t | x_T)$ is exactly the density of \overleftarrow{X}_τ , which is generated by the reversed OU process. From the Kolmogorov forward equation, the reversed process $\{\overleftarrow{X}_t\}$ satisfies

$$d\overleftarrow{X}_\tau = \left(\overleftarrow{X}_\tau + 2 \nabla \log q_{T-\tau}(\overleftarrow{X}_\tau) \right) dt + \sqrt{2} dB_\tau,$$

where $q_{T-\tau}$ is the law of $\overleftarrow{X}_{T-\tau}$.

2.2 Score Matching

To execute the reversed process, one of the most challenging problem is to estimate the score function $\nabla \log q_t(x)$. Let \mathcal{F} be a family of candidate functions, for example, the functions can be represented by neural networks. Our goal is to find

$$\arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} [\|S_t(x) - \nabla \log q_t(x)\|^2].$$

In this subsection, we introduce the idea of score matching.

Since $\nabla \log q_t(x)$ is independent with S_t ,

$$\arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} [\|S_t(x) - \nabla \log q_t(x)\|^2] = \arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} [\|S_t(x)\|^2 - 2 \langle S_t(x), \nabla \log q_t(x) \rangle].$$

Let $\gamma(\cdot)$ be the density of the standard Gaussian distribution. We then

show that estimating $\nabla \log q_t(x)$ is equivalent to find an S_t minimizing

$$\mathbb{E}_{x \sim q_t} \left[\left\| S_t(x) - \frac{Z_t}{\sqrt{1-e^{-2t}}} \right\|^2 \right].$$

We give a proof for the one dimensional case.

The proof is similar in high-dimension space.

Lemma 1.

$$\arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} [\|S_t(x) - \nabla \log q_t(x)\|^2] = \arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} \left[\left\| S_t(x) - \frac{Z_t}{\sqrt{1 - e^{-2t}}} \right\|^2 \right].$$

Proof. From direct calculation,

$$\begin{aligned} \mathbb{E}_{x \sim q_t} [\langle S_t(x), \nabla \log q_t(x) \rangle] &= \int_{\mathbb{R}} S_t(x) (\log q_t(x))' q_t(x) dx \\ &\stackrel{(\spadesuit)}{=} - \int_{\mathbb{R}} q_t(x) S_t'(x) dx \\ &= - \int_{\mathbb{R}} \int_{\mathbb{R}} S_t'(e^{-t}x_0 + \sqrt{1 - e^{-2t}}z_t) q_0(x_0) \gamma(z_t) dx_0 dz_t \\ &\stackrel{(\clubsuit)}{=} - \frac{1}{\sqrt{1 - e^{-2t}}} \int_{\mathbb{R}} q_0(x_0) \int_{\mathbb{R}} S_t(e^{-t}x_0 + \sqrt{1 - e^{-2t}}z_t) \gamma(z_t) \cdot z_t dx_0 dz_t \\ &= \mathbb{E}_{x \sim q_t} \left[\left\langle S_t(x), \frac{Z_t}{\sqrt{1 - e^{-2t}}} \right\rangle \right], \end{aligned}$$

where (\spadesuit) and (\clubsuit) is derived through integrating by parts. Therefore,

$$\begin{aligned} \arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} [\|S_t(x) - \nabla \log q_t(x)\|^2] &= \arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} \left[\|S_t(x)\|^2 - 2 \left\langle S_t(x), \frac{Z_t}{\sqrt{1 - e^{-2t}}} \right\rangle \right] \\ &= \arg \min_{S_t \in \mathcal{F}} \mathbb{E}_{x \sim q_t} \left[\left\| S_t(x) - \frac{Z_t}{\sqrt{1 - e^{-2t}}} \right\|^2 \right]. \end{aligned}$$

□

References