

§7. Diffusion Models.

Goal: Given iid samples $\{x_i\}_{i=1}^n$ from the target distribution p .
generate more samples from p (approximately).

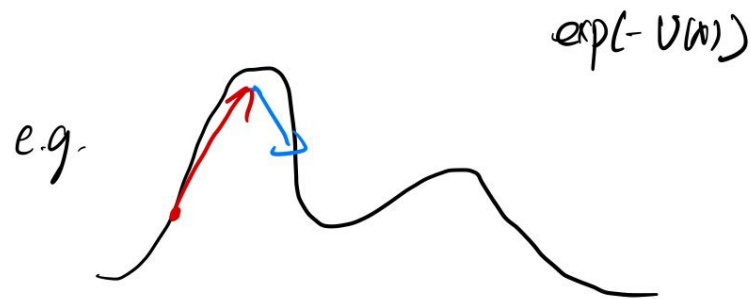
§7.1. Langevin and score-matching.

Suppose $p(x) \propto \exp(-V(x))$. $V: \mathbb{R}^d \rightarrow \mathbb{R}$

Langevin diffusion. $dX_t = -\nabla V(X_t) dt + \sqrt{2} dB_t$

↳ move toward
the modes

↳ some random noise
preventing $X_t \rightarrow$ (local) max.



Thm. If p is log-concave. (i.e., V is convex) then $\text{Law}(X_t) \rightarrow p$.

Question. How to estimate the score function. $\nabla \log p(x) = -\nabla V(x)$. ?

Let $S(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be our score estimator. (e.g. some deep network).
↳ parameter.

Goal: minimize $J(\theta) := \mathbb{E}_{x \sim p} \| S(x; \theta) - \underbrace{(-\nabla U(x))}_{\substack{\text{↳ we don't have access to this.}}}\|_2^2$

Lemma. Assume some boundary conditions. We have

$$J(\theta) = \frac{1}{2} \mathbb{E}_{x \sim p} \| S(x; \theta) \|^2 + \mathbb{E}_{x \sim p} \nabla \cdot S(x; \theta) + C$$

↳ can be estimated using $\{x_i\}_{i=1}^n$

↳ doesn't depend on θ .

Pf. We write.

$$J(\theta) = \frac{1}{2} \mathbb{E}_{x \sim p} \| S(x; \theta) \|^2 + \underbrace{\frac{1}{2} \mathbb{E}_{x \sim p} \| \nabla U(x) \|^2 + \mathbb{E}_{x \sim p} \langle S(x; \theta), \nabla U(x) \rangle}_{=: C}$$

For the third term, we compute.

$$\mathbb{E}_p \langle S(x; \theta), \nabla U(x) \rangle = \sum_{k=1}^d \int S_k(x; \theta) \partial_k U(x) p(x) dx$$

$$= - \sum_{k=1}^d \int S_k(x; \theta) \partial_k \log p(x) p(x) dx$$

$$= - \sum_{k=1}^d \int S_k(x; \theta) \partial_k p(x) dx$$

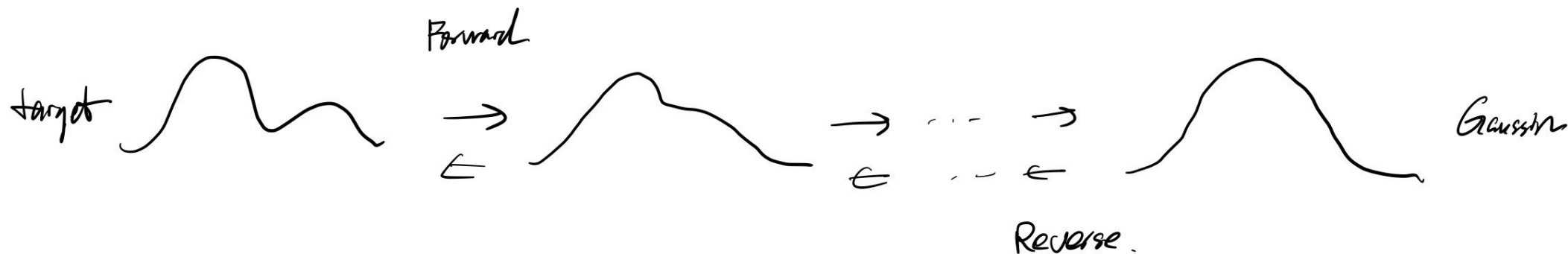
$$= - \sum_{k=1}^d \left(\cancel{S_k(x; \theta) p(x)} \Big|_{\|x\|=\infty}^{\theta} - \int \partial_k S_k(x; \theta) p(x) dx \right)$$

the boundary condition was needed

$$= \int \sum_{k=1}^d \partial_k S_k(x; \theta) p(x) dx = \mathbb{E}_{x \sim p} \nabla \cdot S(x; \theta).$$



§ 7.2 Diffusion models.



Forward process: $dX_t = -X_t dt + \sqrt{2} dB_t$ $X_0 \sim p$.

Reverse process: $X_t^{\leftarrow} := X_{T-t}$.

$$q_t := \text{Law}(X_t)$$

We have $dX_t^{\leftarrow} = \left(X_t^{\leftarrow} + 2 \nabla \log q_{T-t}(X_t^{\leftarrow}) \right) dt + \sqrt{2} dB_t$ $X_0^{\leftarrow} \sim q_T$.

\hookrightarrow score fn. can be learned in the forward process.

Facts. 1) $\text{Law}(X_t) \rightarrow \gamma^d$ (d-dim std. Gaussian).

2) $\text{Law}(X_t^{\leftarrow}) = q_{T-t}$.

Discretized reverse process. (step size $h > 0$, score estimator $S_t \approx \nabla \log q_t$).

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

$X_0^{\leftarrow} \sim \gamma^d$  not q_T since we can't sample from q_T .

(No need to replace all X_t^{\leftarrow} with X_{kt}^{\leftarrow} . The above (1-step) linear SDE can be solved in closed form.)

$$P_t := \text{Law}(X_t^{\leftarrow}).$$

$$\text{Goal: } P_T \approx q_0 = p.$$

★ No log-concavity conditions or isoperimetric inequalities.

Assumptions a). $\forall t \geq 0$. $\text{supp } q_t = \mathbb{R}^d$, $\nabla \log q_t$ is L -Lipschitz.

b). For some $\eta > 0$. $\mathbb{E}_p \|X\|^{2+\eta} < \infty$. Define $m_2^2 := \mathbb{E}_p \|X\|^2$

c). $\forall k \in \mathbb{N}$. $\mathbb{E}_{q_{kh}} \|S_{kh} - \nabla \log q_{kh}\|^2 \leq \epsilon_{\text{score}}^2$.

- Error sources:
- 1) $X_0^{\leftarrow} \sim \gamma^d$ instead of q_T / forward process not converging
 - 2) Using the time-discretized SDE.
 - 3) non-exact score estimator s .

ϕ Kendall.

Thm. 2 of (Chen et al. 2023). Under the previous assumptions, choose $h = T/N \leq \min\{1/L, \epsilon\}$.

$$\Rightarrow TV(P_T, p) \leq \underbrace{\sqrt{KL(p \parallel \gamma^d)} \exp(-T)}_{\text{error 1}} + \underbrace{(2\sqrt{ah} + Lm_2h) \sqrt{T}}_2 + \underbrace{\epsilon_{\text{score}} \sqrt{T}}_3.$$

1). \tilde{P}_T : distribution of X_T^{\leftarrow} if we run the continuous reverse SDE with init = γ^d .

$$TV(\tilde{P}_T, p) \leq TV(\gamma^d, q_T) \leq \sqrt{KL(q_T \parallel \gamma^d)} \leq \exp(-T) \sqrt{KL(p \parallel \gamma^d)}$$

↓

Data processing inequality
along the reverse process

↓

Pinsker's inequality

↓

exp. convergence of
the forward process

To handle errors 2) and 3). we need Girsanov's theorem.

Thm 1 (Girsanov diffusion). Consider $dX_t = b(X_t) dt + \sigma(X_t) dB_t$, $t \leq T$.
and $dY_t = [b(X_t) + \gamma(t, \omega)] dt + \sigma(X_t) dB_t$, $t \leq T$.

Assume some regularity conditions (and Novikov's condition). Define.

$$u(t, \omega) = \gamma(t, \omega) / \sigma(X_t). \quad \longrightarrow \text{normalized error}$$

$$M_t(\omega) = \exp\left(-\int_0^t u(s, \omega) dB_s - \frac{1}{2} \int_0^t u^2(s, \omega) ds\right).$$

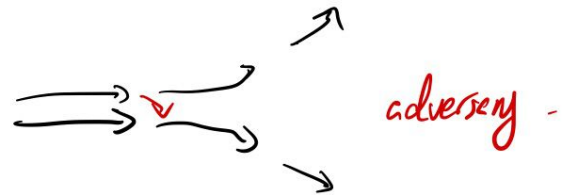
then, we have $\text{Law}(Y_t) / \text{Law}(X_t) = M_t$ The statement here is wrong. TODO: fix this

Corollary. Put $P = \text{Law}(X_t)$, $Q = \text{Law}(Y_t)$. Then we have L² error of the drift term.

$$KL(P \| Q) = \mathbb{E}_P \log \frac{P}{Q} = -\mathbb{E}_P \log M_t = \cancel{\mathbb{E}_P \int_0^T u(s, \omega) dB_s} + \frac{1}{2} \mathbb{E}_P \int_0^T u^2(s, \omega) ds.$$

Remark. Novikov's condition doesn't hold in our setting, and we actually need some bounding argument.

Remark. For ODEs, small L^2 error $\not\Rightarrow$ small worst case error

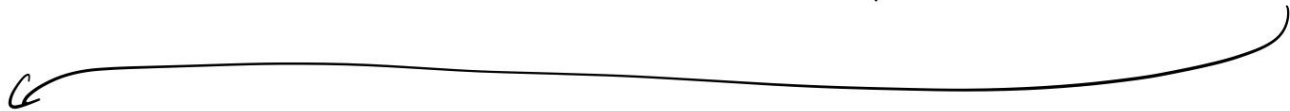
e.g.  adversary.

For SDEs, we don't have "worst case"!

Bounding the discretization error. Under the assumptions of Thm 1, assuming Novikov's condition, let Q_T^{\leftarrow} and $P_T^{q_T}$ denote the distribution of the continuous and discretized reverse process initialized at q_T , respectively.

Goal: Bound $\text{TVC}(P_T^{q_T}, Q_T^{\leftarrow})$.

Pf. Girsanov $\Rightarrow \text{KL}(Q_T^{\leftarrow} \parallel P_T^{q_T}) \leq \sum_{k=0}^{N-1} \mathbb{E}_{Q_T^{\leftarrow}} \int_{t_k}^{t_{k+1}} \mathbb{1}_{\text{detach}} \|S_{T-t_k}(X_{t_k}) - \nabla \log q_{T-t}(X_t)\|^2 dt$.



$$\mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|\mathcal{S}_{T-kh}(X_{kh}) - \nabla \log q_{T-t}(X_t)\|^2$$

$$\lesssim \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|\mathcal{S}_{T-kh}(X_{kh}) - \nabla \log q_{T-kh}(X_{kh})\|^2 + \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|\nabla \log q_{T-kh}(X_{kh}) - \nabla \log q_{T-t}(X_{kh})\|^2$$

$$+ \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|\nabla \log q_{T-t}(X_{kh}) - \nabla \log q_{T-t}(X_t)\|^2$$

$$\lesssim \epsilon_{\text{score}}^2 + \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \left\| \nabla \log \frac{q_{T-kh}}{q_{T-t}}(X_{kh}) \right\|^2 + L^2 \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|X_{kh} - X_t\|^2 \rightarrow \text{movement of } X.$$

↳ change of the score function within each step along the forward process

Define $S(t) = \exp(-(t-kh))\pi$. We have $q_{T-kh} = S \# q_{T-t} * N(0, 1 - \exp(-2(t-kh)))$.

Lemma (Score perturbation Lemma). $M_0, M_1 \in \mathbb{R}^{d \times d}$, $\|M_0 - I\|_{\text{op}} \leq \zeta < 1$, M_1 symmetric.

$q = \exp(-H) \in \mathcal{P}(\mathbb{R}^d)$, ∇H L -Lipschitz, $L \leq \frac{1}{4\|M_1\|_{\text{op}}}$.

$$\Rightarrow \left\| \nabla \log \frac{M_0 \# q * N(0, M_1)}{q}(\theta) \right\| \lesssim L \sqrt{\|M_1\|_{\text{op}}} + L\zeta \|\theta\| + (\zeta + 2\|M_1\|_{\text{op}}) \|\nabla H(\theta)\|.$$

$$\Rightarrow \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \left\| \nabla \log \frac{q_{T+h}}{q_{T-1}}(X_{T+h}) \right\|^2 = \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \left\| \nabla \log \frac{\sum_{\#} q_{T-t} * N(0, 1 - \exp(-2(t-h)))}{q_{T-t}}(X_{T+h}) \right\|^2$$

$$\lesssim L^2 d h + L^2 h^1 \mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \|X_{T+h}\|^2 + L^2 h^1 \underbrace{\mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \left\| \nabla \log q_{T-t}(X_{T+h}) \right\|^2}_{\rightarrow \lesssim \left\| \nabla \log q_{T-t}(X_t) \right\|^2 + L^2 \|X_{T+h} - X_t\|^2}$$

$$\rightarrow \lesssim \left\| \nabla \log q_{T-t}(X_t) \right\|^2 + L^2 \|X_{T+h} - X_t\|^2$$

Claim. $\mathbb{E} \|X_{t+1}\|^2 \leq d \vee m_2^2$. $\mathbb{E} \left\| \nabla \log q_t(X_t) \right\|^2 \leq Ld$. $\mathbb{E} \|X_{t+1} - X_t\|^2 \lesssim (t-s)^1 m_2^2 + (t-s)d$.

Thus. $\mathbb{E}_{\mathcal{Q}_T^{\epsilon}} \left\| S_{T+h}(X_{T+h}) - \nabla \log q_{T+h}(X_t) \right\|^2 \leq \epsilon_{\text{score}}^2 + L^2 d h + L^2 h^2 (d + m_2^2) + L^3 h^2 d$
 $+ L^4 h^2 (h^2 m_2^2 + d h)$
 $\lesssim \epsilon_{\text{score}}^2 + L^2 d h + L^2 m_2^2 h^2$.

