

# Maximum-Likelihood Training of Score-Based Diffusion Models

Sequentially corrupt data with increasing noise

Learn to reverse corruption to create generative model

Score-Matching with Langevin Dynamics (Song + Ermon 19)

Estimate "score" at each noise scale

grad of log probability density

Sample from sequence of decreasing noise scales by Langevin Dynamics

Denoising Diffusion Probabilistic Modeling DDPM (Sohl-Dickstein 15, Ho 15)

Can be shown to implicitly compute a score

$$x_0 \rightarrow x_t$$
$$dx = f(x, t) dt + g(t) dw \quad [\text{data} \rightarrow \text{noise}]$$

forward SDE

$$x_0 \leftarrow x_t$$
$$dx = \left[ f(x, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x) \right] dt + g(t) dw \quad [\text{noise} \rightarrow \text{data}]$$

reverse SDE

Perturb data with a continuum of distributions that evolve in time according to diffusion process

Diffuse data to noise with non-parametric SDE

The reverse process can be ~~beaten~~ described by reverse-time SDE that can be derived from forward SDE given  $\nabla_x \log f_t(x)$

Score-based diffusion can be converted to flow for numerical likelihood computation

Score-based models

Let  $p(x)$  unknown distribution of 0-dim iid samples

Score-based drift models diffuse  $p(x) \rightarrow$  noise by SDE

$$dx = f(x, t) dt + g(t) dw \leftarrow \text{infinitesimal Gaussian noise}$$

where  $f(\cdot, t)$  drift coefficient

$g(t) \in \mathbb{R}$  diffusion coefficient

$w \in \mathbb{R}^0$  Wiener process (Brownian Motion) aka diffusion by Langevin eq

$x(t)$  then is a diffusion process

$p_t(x)$  marginal distribution

$p_{0t}(x'|x)$  transition distribution from  $x_0$  to  $x_t$

Intuitively: smoothen the data distribution in time (by adding noise)

We smooth into an analytically tractable prior  $\pi(x)$  at  $t=T$

$$\text{i.e. } p_T(x) \propto \pi(x)$$

this can be done by  
a family of SDEs

Now we can reverse the SDE

$$\begin{aligned} dx &= f(x, t) dt + g(t) dw \\ \text{data } x(0) &\uparrow \quad \quad \quad x(T) \text{ noise} \\ dx &= \left[ f(x, t) - g(t)^2 \nabla_x \log p_t(x) \right] dt + g(t) dw \end{aligned}$$

↑  
standard Wiener process

this is the learnable piece

Choose  $f$  and  $g$  so we have reversal by

Fit a "score-based model" on  $s_\theta(x, t) \approx \nabla_x \log p_t(x)$  [notice not normalized]

$$J_{\text{SM}}(\theta; \lambda(\cdot)) := \frac{1}{2} \int_0^T \mathbb{E}_{p_t(x)} \left[ \lambda(t) \left\| \nabla_x \log p_t(x) - s_\theta(x, t) \right\|_2^2 \right] dt$$

loss function

$$\lambda: [0, T] \rightarrow \mathbb{R}^+$$

the "score matching objective"

We want to transform this objective to something computable

By denising score-matching [Vincent NIPS 11] we can get  $J$  up to a constant ind of  $\Theta$

$$J_{PSM}(\Theta; \lambda(\cdot)) := \frac{1}{2} \int_0^T \mathbb{E}_{p(x)} p_{\theta+}(x'|x) [\lambda(+)] \|\nabla_{x'} \log p_{\theta+}(x'|x) - s_\theta(x', t)\|_2^2 dt$$

when  $f_\theta(x, t)$  linear in  $x$  [loose]

$p_{\theta+}(x'|x)$  is Gaussian

To estimate draw  $(t, x, x')$

$$t \sim [0, T]$$

$$x \sim p(x) \text{ data sample}$$

$$x' \sim p_{\theta+}(x'|x) \Rightarrow \nabla_{x'} \log p_{\theta+}(x'|x) \text{ since } p \text{ gaussian}$$

Once we have  $s_\theta(x, t)$  we will numerical SDE solvers to sample

Likelihood of score-based diffusion Models

There are 2 probabilistic models for which we can define a likelihood

1.  $p_\theta^{SDE}(x)$

Let  $\{\hat{x}_\theta(t)\}_{t \in [0, T]}$  a stochastic process given by

$$d\hat{x} = [f(\hat{x}, t) - g(t)^T s_\theta(\hat{x}, t)] dt + g(t) d\bar{w} \quad \hat{x}_\theta(T) \sim \pi$$

Then  $p_\theta^{SDE}(x)$  is marginal distribution of  $\hat{x}_\theta(0)$

we sample by

$$\hat{x}_\theta(0) \sim p_\theta^{SDE} \text{ numerically solving reverse-time SDE w/ noise } x_\theta(T) \sim \pi$$

2.  $p_\theta^{ODE}(x)$  derived from SDE's probability flow ODE

Just ODE with the same  $p_+(x)$  given by

$$\frac{dx}{dt} = f(x, t) - \frac{1}{2} g(t)^T \nabla_x \log p_+(x) \quad \text{ANL score function}$$

$$\approx f(x, t) - \frac{1}{2} g(t)^T s_\theta(x, t)$$

this is actually a CNF!

So given prior  $\pi(\cdot)$  and trajectory functions  $\tilde{x}_0: [0, T] \rightarrow \mathbb{R}^d$

satisfying ODE then  $\rho_0^{\text{ODE}}$  is the marginal of  $\tilde{x}_0(0)$  when  $\tilde{x}_0(T) \sim \pi$

Can evaluate exactly with ODE solver

Computing this ODE  $\log \rho_0^{\text{ODE}}(x)$  is tractable

But solving requires running a solver at every step

Similarly we can't ~~exactly~~ evaluate  $\rho_0^{\text{SDE}}$  but we can solve lower bound.

Upper  
Lower Bounds

With correct choice of  $\lambda$

$J_{\text{SM}}(\theta; \lambda(\cdot))$  becomes a <sup>upper</sup> bound  $D_{\text{KL}}(\rho \parallel \rho_0^{\text{SDE}})$

$$\text{when } \lambda(t) = g(t)^2$$

$$D_{\text{KL}}(\rho \parallel \rho_0^{\text{SDE}}) \leq J_{\text{SM}}(\theta; g(\cdot)^2) + D_{\text{KL}}(\rho_T \parallel \pi)$$

AM)

$$-\mathbb{E}_{\rho(x)} \left[ \log \rho_0^{\text{SDE}}(x) \right] \leq J_{\text{SM}}(\theta; g(\cdot)^2) + C_1 = J_{\text{SM}}(\theta; g(\cdot)^2) + C_2$$

AM) when  $s_\theta(x, t) = \nabla_x \log g_t(x) \quad \forall t \in [0, T]$  then

$$\rho_0^{\text{SDE}} = \rho_0^{\text{ODE}} = g \text{ and}$$

$$D_{\text{KL}}(\rho \parallel \rho_0^{\text{SDE}}) = J_{\text{SM}}(\theta; g(\cdot)^2) + D_{\text{KL}}(\rho_T \parallel \pi)$$