

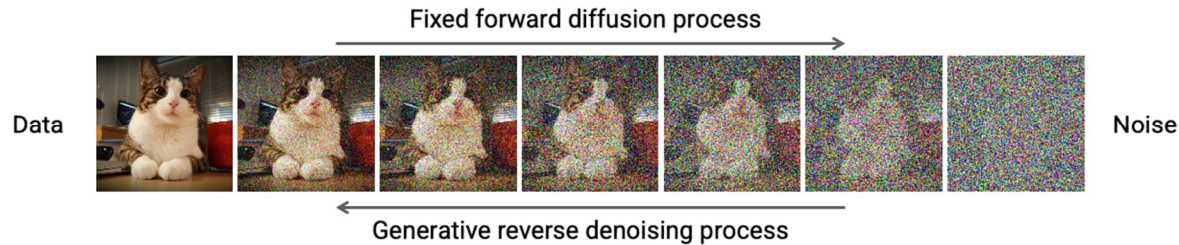
Restricted Boltzmann machines and generative diffusion models III

Eric Paquet

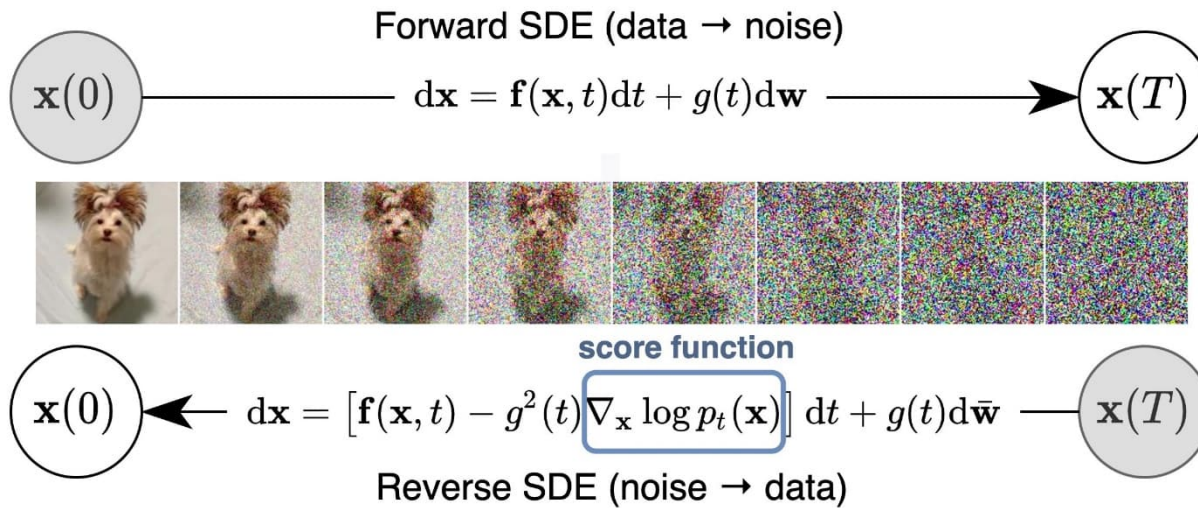
National Research Council

28 November 2024

Denosing diffusion probabilistic models:



Score-based generative modelling with stochastic differential equations:



Forward process; noising process; Wiener process; drift; diffusion coefficient:

$$dx = f(x, t)dt + g(t)dw \quad (1)$$

Reverse process; denoising process; score function; exact results:

$$dx = \left[f(x, t) - g(t)^2 \nabla_x \ln p(x) \right] dt + g(t)d\bar{w} \quad (2)$$

Variance preserving stochastic differential equation:

$$f(x, t) = -\frac{1}{2}\beta(t)x \quad (3)$$

$$g(t) = \sqrt{\beta(t)} \quad (4)$$

Noise scheduling: linear, sinusoidal, signal-to-noise ratio e.g.

$$\beta(t) = \beta_{\min} + \frac{t}{T}(\beta_{\max} - \beta_{\min}) \quad (5)$$

<https://arxiv.org/abs/2011.13456>

Restricted Boltzmann machine; energy:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{v}^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} \quad (6)$$

Probability associated with the visible units; marginalisation over the hidden units:

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \equiv \frac{e^{-F(\mathbf{v})}}{Z} \quad (7)$$
$$F = \sum_{\mathbf{v}} F(\mathbf{v})$$

Free energy associated with the visible units:

$$F_{\boldsymbol{\theta}}(\mathbf{v}) = -\mathbf{b}^\top \mathbf{v} - \sum_{j=1}^H \ln \left(1 + e^{c_j + \mathbf{W}_j^\top \mathbf{v}} \right), \quad \boldsymbol{\theta} = \{ \mathbf{b}, \mathbf{c}, \mathbf{W} \} \quad (8)$$

The connection is made through the score function associated with the model:

$$\boxed{\nabla_{\mathbf{v}} \ln p_{\theta}(\mathbf{v}) = -\nabla_{\mathbf{v}} F_{\theta}(\mathbf{v})} \quad (9)$$

$$\frac{\partial F_{\theta}(\mathbf{v})}{\partial v_i} = \frac{v_i - b_i}{\sigma_i^2} - \sum_j \frac{W_{ij}}{\sigma_i^2} \sigma \left(c_j + \sum_k \frac{v_k}{\sigma_k^2} W_{kj} \right) \quad (10)$$

$$v_i \in \{0,1\} \Rightarrow \nabla? \quad (11)$$

The gradient may be evaluated with the Gumbel trick.

Alternatively, with the mean field approximation:

$$\boxed{v_i \rightarrow \langle v_i \rangle} \quad (12)$$

$$\boxed{v_i \in \{0,1\} \rightarrow \langle v_i \rangle \in [0,1]} \quad (13)$$

$$\langle v_i \rangle = 1p(v_i = 1) + 0p(v_i = 0) \Rightarrow$$
$$\boxed{\langle v_i \rangle = p(v_i = 1)} \quad (14)$$

Logit; activation function; real number:

$$a_i = \text{logit}(\langle v_i \rangle) = \ln\left(\frac{\langle v_i \rangle}{1 - \langle v_i \rangle}\right) \in \mathbb{R} :]0,1[\rightarrow \mathbb{R} \quad (15)$$

This transformation is invertible:

$$\langle v_i \rangle = \sigma(a_i) = \frac{1}{1 + \exp(-a_i)} \quad (16)$$

Free energy of the activation function:

$$v_i \rightarrow \langle v_i \rangle = p(v_i = 1) = \sigma(a_i) \Rightarrow$$
$$\boxed{F_{\theta}(\mathbf{a}) = -\mathbf{b}^{\top} \sigma(\mathbf{a}) - \ln \left(1 + e^{c_j + \mathbf{W}_j^{\top} \sigma(\mathbf{a})} \right)} \quad (17)$$

Score of the activation function associated with the model:

$$\boxed{s_{\theta}(\mathbf{a}) = -\nabla_{\mathbf{a}} F_{\theta}(\mathbf{a})} \quad (18)$$

$$s_{\theta}(\mathbf{a}) = \left(\mathbf{b} + \sum_{j=1}^H \sigma(c_j + \mathbf{W}_j^{\top} \sigma(\mathbf{a})) \mathbf{W}_j \right) \odot \sigma(\mathbf{a}) \odot (1 - \sigma(\mathbf{a})) \quad (19)$$

Denosing score matching; non-tractable; expectation:

$$\mathcal{L}(\theta) = \mathbb{E}_{t, a_0, a_t} \left[\lambda(t) \left\| s_{\theta}(a_t, t) - \nabla_{a_t} \ln p(a_t) \right\|_2^2 \right] \quad (20)$$

Denosing score matching; conditional score; closed-form; expectation:

$$\mathcal{L}(\theta) = \mathbb{E}_{t, a_0, a_t} \left[\lambda(t) \left\| s_{\theta}(a_t, t) - \nabla_{a_t} \ln p(a_t | a_0) \right\|_2^2 \right] \quad (21)$$

Noising process; generation of noisy data; variance preserving SDE; location-scale; reparametrisation trick:

$$da = -\frac{1}{2} \beta(t) a dt + \sqrt{\beta(t)} dw_t \quad (22)$$

$$\begin{aligned}
& \overbrace{\mu(a_0, t)} \\
& a_t \simeq e^{-\frac{1}{2} \int_0^t \beta(s) ds} a_0 + \sigma_t \zeta, \quad \zeta \sim \mathcal{N}(0, \mathbf{I}) \Rightarrow \\
& a_t \sim p(a_t | a_0) = \mathcal{N}(a_t; \mu(a_0, t), \Sigma(t)) \quad (23)
\end{aligned}$$

$$\mu(a_0, t) = e^{-\frac{1}{2} \int_0^t \beta(s) ds} a_0 \quad (24)$$

$$\Sigma(t) = \sigma_t^2 \mathbf{I}, \quad \sigma_t^2 = \int_0^t e^{-\int_s^t \beta(r) dr} \beta(s) ds \simeq 1 - e^{-\int_0^t \beta(s) ds} \quad (25)$$

$$\nabla_{a_t} \ln p(a_t | a_0) = -\Sigma^{-1}(t) (a_t - \mu(a_0, t)) \quad (26)$$

$$\nabla_{a_t} \ln p(a_t | a_0) = -\frac{1}{\sigma_t^2} \left(a_t - e^{-\frac{1}{2} \int_0^t \beta(s) ds} a_0 \right) \quad (27)$$

Score matching:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{t, a_0, a_t} \left[\lambda(t) \left\| s_{\boldsymbol{\theta}}(a_t, t) + \frac{1}{\sigma_t^2} \left(a_t - e^{-\frac{1}{2} \int_0^t \beta(s) ds} a_0 \right) \right\|_2^2 \right] \quad (28)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{t, a_0, a_t} \left[\lambda(t) \left\| \left(\mathbf{b} + \sum_{j=1}^H \sigma(c_j + \mathbf{W}_j^\top \sigma(a_t)) \mathbf{W}_j \right) \odot \sigma(a_t) \odot (1 - \sigma(a_t)) + \frac{1}{\sigma_t^2} \left(a_t - e^{-\frac{1}{2} \int_0^t \beta(s) ds} a_0 \right) \right\|_2^2 \right] \quad (29)$$

Stochastic optimisation:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \quad (30)$$

The activation function associated with the data may be evaluated from a mini-batch; sampling with replacement:

$$a_{0,i}^{\mathcal{B}} = \text{logit} \left(\frac{1}{|\mathcal{B}|} \sum_{j \in \mathcal{B}} v_{0,i}^{(j)} \right), \quad \mathcal{B} \ll \quad (31)$$

Pros: simplicity; deterministic, differentiable

Cons: ignore stochasticity, which may lead to a biased gradient estimate; underestimation of variance

Gumbel trick

Pros: unbiased gradient estimate

Cons: high variance in gradient estimates; temperature sensitivity

Generation of new samples from random noise; reverse process; Euler – Maruyama integration:

$$a_{t-1} = a_t + \left[f(a_t, t) - g(t)^2 s_\theta(a_t, t) \right] \Delta t + g(t) \sqrt{\Delta t} \zeta$$
$$\zeta \sim \mathcal{N}(0, \mathbf{I}) \quad (32)$$

$$a_T \rightarrow a_{T-1} \rightarrow \dots \rightarrow a_0 \quad (33)$$

Inverse of the activation function:

$$\boxed{p(v_i = 1) = \langle v_i \rangle = \sigma(a_{0,i}), \quad \mathbb{R} \rightarrow]0, 1[} \quad (34)$$

Sampling:

$$\begin{aligned} (p(v_i = 1) < u) &\Rightarrow v_i = 1, & v_i = 0 \\ u &\sim \mathcal{U}(0,1) \\]0,1[&\rightarrow \{0,1\} \end{aligned} \tag{35}$$