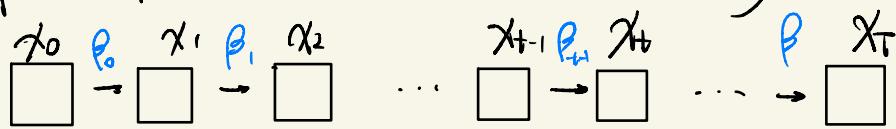


Diffusion models

1. Forward process (deterministic) $g(x_t | x_{t-1})$



$$g(x_t | x_{t-1}) \sim N(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

Markov process: $g(x_{1:T} | x_0) = \prod_{t=1}^T g(x_t | x_{t-1})$

One step process: $x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$

$$\alpha_t = 1 - \beta_t$$

$$\beta_T = 1 - \alpha_T$$

$$\begin{aligned} &= \sqrt{1-\beta_t} (\sqrt{1-\beta_{t-1}} x_{t-2} + \sqrt{\beta_{t-1}} \varepsilon_{t-2}) + \beta_t \varepsilon_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_{t-1}} \varepsilon_{t-1} \\ &= \boxed{\sqrt{\alpha_t} x_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon} \quad \delta^2 = \sigma_1^2 + \sigma_2^2 \end{aligned}$$

continuous form (Stochastic diffusion equation)

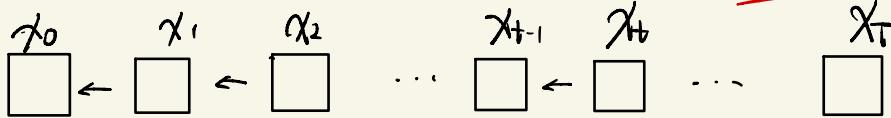
$$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$$

$$= \left(1 - \frac{1}{2}\beta_t\right) x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$$

$$\Delta x_t = -\frac{1}{2}\beta_t x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$$

$$\frac{dx}{dt} = -\frac{1}{2}\beta_t x + \sqrt{\beta_t} dW \leq \frac{dx}{dt} = f(x, t) + g(x, t) W$$

Reverse Process $P(X_{t-1}|X_t)$ is unknown $\rightarrow P_\theta(X_{t-1}|X_t)$



When $\beta_t \ll 1$, $P(X_{t-1}|X_t) \sim N(X_{t-1}; \mu_t(X_t, t), \Sigma_t(X_t, t))$

Markov Process: $P(X_{0:T}) = P(X_T) \prod_{t=1}^T P(X_{t-1}|X_t)$

Idea: Use a neural network to approximate $\hat{g}(X_{t-1}|X_t)$

Sadly: $\hat{g}(X_{t-1}|X_t)$ is not tractable, trajectories arriving X_{t-1}

$$\hat{g}(X_{t-1}|X_t) = \frac{\hat{g}(X_t|X_{t-1}) \hat{g}(X_{t-1})}{\hat{g}(X_t)} \rightarrow x$$

But: $\hat{g}(X_{t-1}|X_t, X_0)$ is deterministic!

$$\hat{g}(X_{t-1}|X_t, X_0) = \frac{\hat{g}(X_t|X_{t-1}, X_0) \hat{g}(X_{t-1})}{\hat{g}(X_t|X_0)} \xrightarrow{\text{known}} \hat{g}(X_t|X_0) \xrightarrow{\text{known}}$$

$$\text{Also: } X_0 = \frac{1}{\sqrt{1-\alpha_t}}(X_t - \sqrt{1-\alpha_t}\varepsilon_t) \sim N(X_t; \tilde{\mu}_t, \tilde{\Sigma}_t)$$

$$\tilde{\mu}_t = \frac{1}{\sqrt{1-\alpha_t}}(X_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}}\varepsilon_t)$$

$$P_\theta \xrightarrow{\text{expectation}} \xrightarrow{\text{fixed}}$$

Now P_θ can learn $\hat{g}(X_{t-1}|X_t)$ in the following way:

$$P_\theta(X_t, t) \rightarrow E_{X_0 \sim P(X_0)} \hat{g}(X_{t-1}|X_t, X_0)$$

$$P_\theta(X_t, t) \rightarrow E_{X_t \sim P(X_t|X_0), X_0 \sim P(X_0), t \sim U(0, T)} \hat{g}(X_{t-1}|X_t, X_0)$$

More on \hat{p}_t : $\hat{p}_t(x_t, t) = a_t(x_t - b_t \epsilon_t)$

$$\epsilon_t = -\frac{\hat{p}_t}{ab} + \frac{at}{ab} \delta_t$$

$$p_\theta \rightarrow \epsilon_\theta \rightarrow \epsilon_t$$

\uparrow Δ input \uparrow

objective: $E_{x_t, x_0, t} \| \epsilon_\theta - \epsilon_t \|^2$

Idea: Essentially, the network is learning to remove noise.
denoising process

To derive: use the variational lower bound for $P(x_0)$

Question: How to relate to $\nabla_x \log g(x)$, the score?

Again: $\nabla_x \log g(x)$ is not tractable. But: $\nabla_x \log g(x|x_0)$ is.

$$\nabla_{x_t} \log g(x_t|x_0) = -\nabla_{x_t} \frac{(x_t - \mu)^2}{2\sigma^2} = -\frac{2\delta_t \epsilon_t}{2\sigma^2} = -\frac{\epsilon_t}{\sigma^2}$$

So $\rightarrow E_{x_t, x_0, t} \nabla_{x_t} \log g(x_t|x_0) \rightarrow \nabla_x \log g(x_t)$

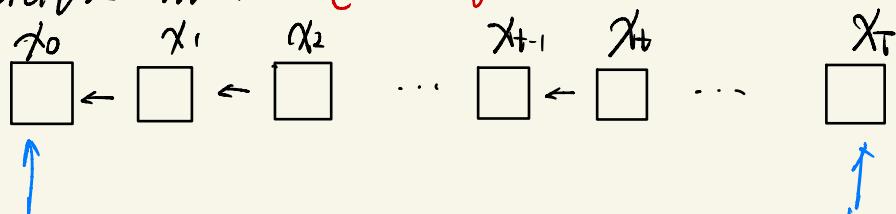
Other Score matching method:

$$\min E_{x_t, t} \| s_\theta(x_t, t) \|_2^2 + \text{tr} \| \nabla_{x_t} s_\theta(x_t, t) \| \rightarrow \boxed{\begin{array}{l} \text{sliced} \\ \text{projected 1D} \end{array}}$$

continuous form: (RDE)

$$dx = -\frac{1}{2} \beta_t' x - \beta_t' \boxed{\nabla_x \log g(x_t)} dt + \sqrt{\beta_t'} dW_t$$

Generative model (Denoising diffusion)

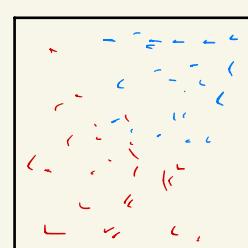
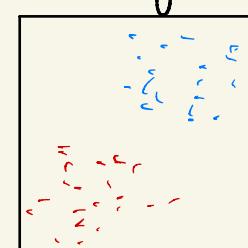
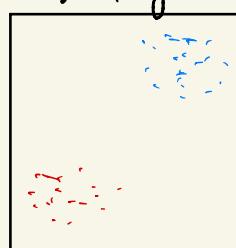
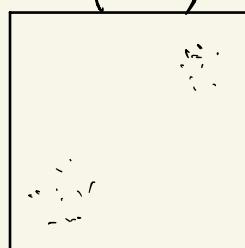


$$x_{t-1} = a_t [x_t - b_t \epsilon_0(x_t, t)] + \sqrt{b_t} \epsilon$$

$$\vdots$$

$$x_0 = a_1 [x_1 - b_1 \epsilon_0(x_1, t=1)]$$

score-based generative model



$$x_{t-1} = x_t - \boxed{\lambda \nabla \log p(x_t)} + \bar{r} \epsilon$$

etc

avoid local minima.

To define the objective function: (Maximum Likelihood)

$$\max_{\theta} P_{\theta}(x_0) \rightarrow \max_{\theta} \log P_{\theta}(x_0)$$

variational lower bound.

$$\log P_{\theta}(x_0) \geq \log P_{\theta}(x_0) - KL \left[q(x_{1:T}|x_0) \| p_{\theta}(x_{1:T}|x_0) \right]$$

$$-\log P_{\theta}(x_0) \leq -\log P_{\theta}(x_0) + \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{P_{\theta}(x_{0:T})} \right]$$

$$-\log P_{\theta}(x_0) \leq \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{P_{\theta}(x_{0:T})} \right] P_{\theta}(x_0)$$

X

$$-\mathbb{E}_{x_0 \sim q(x_0)} \left[\log P_{\theta}(x_0) \right] \leq \mathbb{E}_{x_{0:T} \sim q(x_{0:T})} \left[\log \frac{q(x_{0:T}|x_0)}{P_{\theta}(x_{0:T})} \right]$$

$$\min L_0 \rightarrow \min \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(x_t|x_{t-1})}{P_{\theta}(x_T) \prod_{t=1}^T P_{\theta}(x_{t+1}|x_t)} \right]$$

$$= \min \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(x_{t-1}|x_t)}{P_{\theta}(x_T) \prod_{t=1}^T P_{\theta}(x_{t+1}|x_t)} \frac{q(x_t)}{q(x_{t-1})} \right]$$

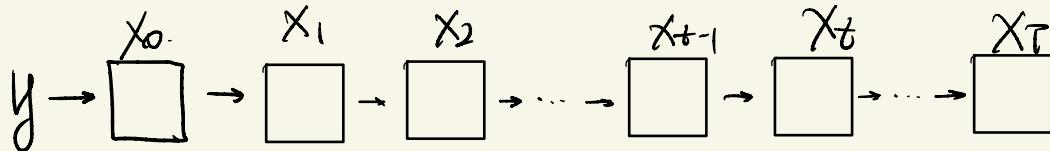
$$= \min \mathbb{E}_q \left[\sum_{t=1}^T \log \frac{q(x_{t-1}|x_t)}{P_{\theta}(x_{t-1}|x_t)} + \sum_{t=1}^T \log \frac{q(x_t)}{q(x_{t-1})} - \log P_{\theta}(x_T) \right]$$

$$= \min \mathbb{E}_q \left[\log \frac{q(x_T|x_0)}{P_{\theta}(x_T)} + \sum_{t=1}^T \log \frac{q(x_{t-1}|x_t, x_0)}{P_{\theta}(x_{t-1}|x_0)} \right]$$

Final objective: given t , $x_t \sim q(x_t | x_0)$, $x_0 \sim p(x_0)$,

$$\min_{\theta} \mathbb{E}_{x_0} \text{KL} [q(x_{t+1} | x_t, x_0) \| p_\theta(x_{t+1} | x_t)]$$

Conditional diffusion model.



The posterior became: $\log P(x_0 | y)$

$$\nabla_{x_0} \log P(x_0 | y) = \nabla_{x_0} \log \frac{P(y | x_0) P(x_0)}{p_y}$$

$$= \underbrace{\nabla_{x_0} \log P(y | x_0)}_{\text{classifier/ forward model}} + \underbrace{\nabla_{x_0} \log P(x_0)}_{\text{unconditioned}}$$

$$\nabla_x \log P(x | y) \rightarrow \text{obj: } \|\mathcal{SD}(x_t, t, y) - \nabla_x \log P(x | y)\|_2^2$$

$$\epsilon_0(x_t, x_i, y) \rightarrow \text{obj: } \mathbb{E}_{t, x_i, x_0} \|\epsilon_0(x_t, x_i, y) - \epsilon\|_2^2$$