

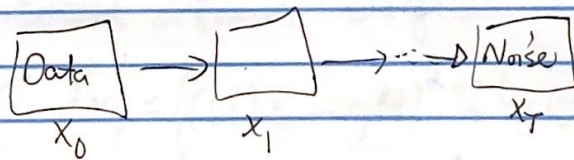
Diffusion Generative Models

Impressive framework that DALL-E2, midjourney etc based on

Better image quality than GANs, but slower.

Original Idea: forward diffusion - noising
gradually add noise to input

Key insight! \rightarrow reverse diffusion - denoising
can gradually remove noise from input!



$$q(x_t | x_{t-1}) = \mathcal{N}(x_t | \mu = \sqrt{1 - \beta_t} x_{t-1}, \sigma = \sqrt{\beta_t})$$

β_t "noise schedule"

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

$\left\{ \begin{array}{l} \text{"Brownian motion"} \\ \text{"Wiener process"} \end{array} \right.$
 $\epsilon_{t-1} \sim \mathcal{N}(0, 1)$

(can prove)

(this is why we need the $\sqrt{1 - \beta_t}$ & $\sqrt{\beta_t}$ factors)

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_0, \sqrt{1 - \alpha_t})$$

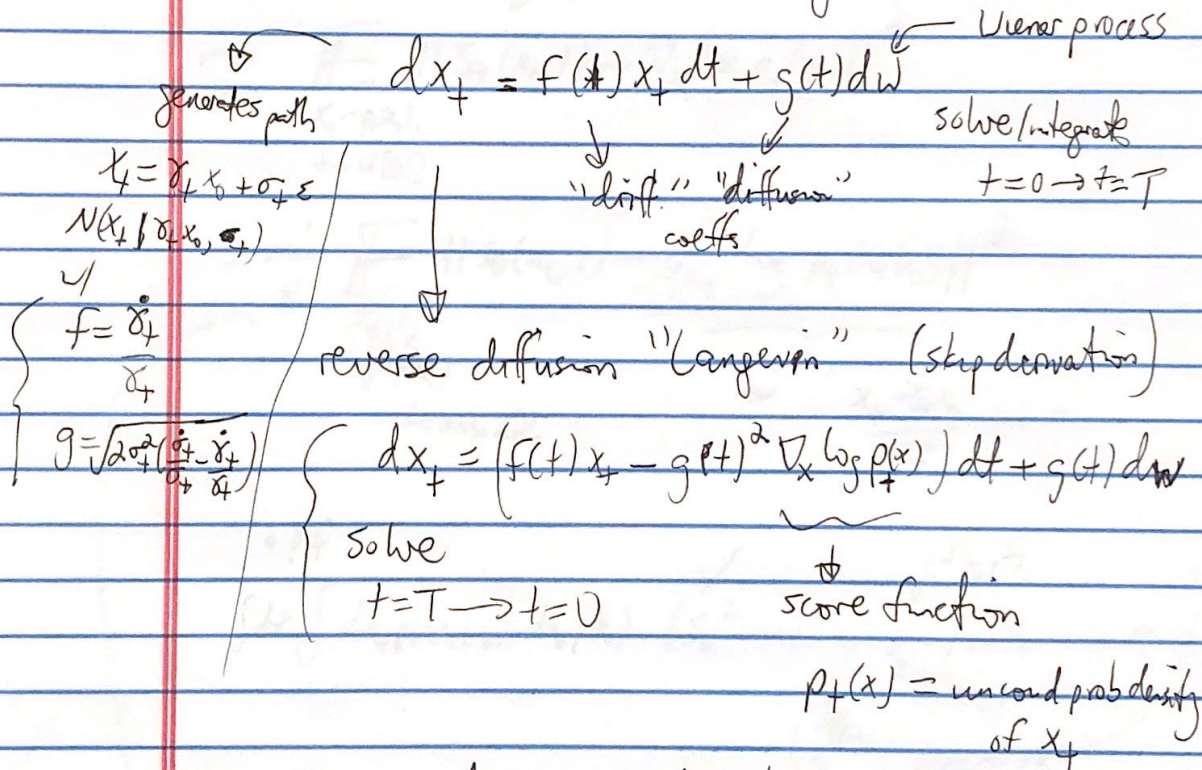
Same Gaussian prob path
as in flow matchy!

$$\alpha_t = \prod_{s=1}^t (1 - \beta_s)$$

β_t s.t. $\alpha_T \rightarrow 0$ and $\mathcal{N}(x_T; 0, 1)$ as $T \rightarrow \infty$
chosen

Original version discrete (e.g. DDPM Ho et al 2020)

↓
continuous time version (Song et al 2020, 2021)
"Score-based generative models"



→ To generate, just need to learn score $\nabla_x \log p_t(x)$!

→ should be much easier than learning $p_t(x)$ itself (don't need samples)

~~not~~

~~MLE~~

Trick to learn $\nabla_x \log P_t(x)$ \rightarrow conditional score noisy

Claim: \checkmark NN

$$\min_{\theta} \mathbb{E}_{\substack{x_t \sim p(x_t) \\ t \sim U[0,1]}} \|s_{\theta}(x_t, t) - \nabla_{x_t} \log P_t(x_t)\|^2$$

$$= \min_{\theta} \mathbb{E}_{\substack{x_t, t \\ x_t \sim p(x_t) \\ \hookrightarrow \text{data dist.}}} \|s_{\theta}(x_t, t) - \underbrace{\nabla_{x_t} \log P_t(x_t | x_0)}_{\sim \frac{x_t - \mu_t | x_0}{\sigma_t^2}}\|^2$$

• pf

$$\int dx_t \int dx_0 p_t(x_t | x_0) \left(s_{\theta}^2 - 2 s_{\theta} \underbrace{\nabla_{x_t} \log P_t(x_t | x_0)}_{\substack{\text{int. by p.t.s} \\ p_t(x_t | x_0)}} + \text{const} \right)$$

$\left(\int dx_t s_{\theta} \right) p_t(x_t | x_0)$ \checkmark

So a very simple objective for learning score!

As in FM, s_{θ} can be any NN, not restricted like in NF