

# Personal Note # 1: Diffusion Models for Hand-Writing generation.

Thứ

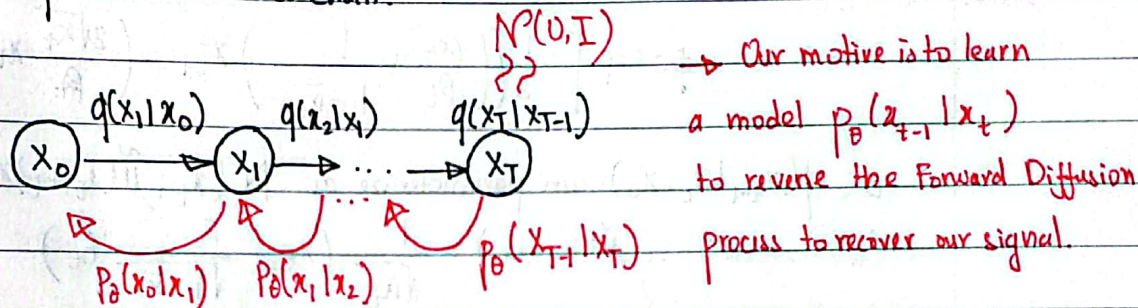
Ngày

No.

## \* Diffusion model re-iteration.

↳ Diffusion model repeatedly distorts an input signal by adding noise after every step.

We can define every step as a conditional data distribution  $q(x_t|x_{t-1})$ . Hence, this can be interpreted as a Markov Chain:



↳ Define  $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t I)$

↳ We can define  $q(x_t|x_0)$  in closed form: Define  $1-\beta_t = \alpha_t$

↳ Using Re-parameterization we have:

$$\begin{aligned} x_t &= \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \epsilon_{t-1} \\ &= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_{t-1}} \epsilon_{t-2}) + \sqrt{1-\alpha_t} \epsilon_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_t \alpha_{t-1}} \epsilon_{t-2} \end{aligned}$$

$$\begin{aligned} &= \dots \\ &= \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1} x_0 + \sqrt{1-\alpha_t \alpha_{t-1} \dots \alpha_1} \epsilon_0 \\ &= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon_0 \quad (\epsilon_0 \sim \mathcal{N}(0, I)) \end{aligned}$$

Where  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

\*  $\Rightarrow q(x_{0:t}|x_0) = \mathcal{N}(x_0; \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I)$  (1)

↳ From (1) we can calculate  $q(x_{t+1}|x_t, x_0)$ . It's hard to calculate  $q(x_{t+1}|x_t)$  directly because of intractability.



$$\Rightarrow q(x_{t+1}|x_t, x_0) = q(x_t|x_{t-1}, x_0) \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

\*  $q(x_t|x_{t-1}, x_0) \propto q(x_t|x_{t-1})$  because  $x_0$  is fixed

\* Not dependent on  $x_t \propto \beta_t$

$$\propto \exp \left[ -\frac{1}{2} \left( \frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t+1}} \right) x_{t+1}^2 - \left( \frac{2\sqrt{\alpha_t}}{\beta_t} x_t + \frac{2\sqrt{\alpha_t}}{1-\bar{\alpha}_t} x_{t-1} + C \right) \right]$$

$\Rightarrow q(x_{t+1}|x_t, x_0)$  can parameterize as  $N(x_{t+1}; \tilde{\mu}(x_t, \epsilon_t), \tilde{\beta}_t I)$  where:

$$\tilde{\mu}(x_t, \epsilon_t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_t \right)$$

$$\tilde{\beta}_t = \beta_t \frac{1-\bar{\alpha}_{t+1}}{1-\bar{\alpha}_t}$$

$\Rightarrow$  Loss function: To optimize  $p_\theta(\cdot)$ , we use the variational lower bound:

\* We know this!

$$E_{\text{vLB}} = \mathbb{E}_q \left[ L_T + \sum_{t>1} D_{\text{KL}}(q(x_{t+1}|x_t, x_0) \| p_\theta(x_{t+1}|x_t)) + L_0 \right]$$

\* Basically constant. \* We are particularly interested in these KL divergence terms.

$\Rightarrow$  The KL divergence terms can be simplified into:

$$L_t^{\text{simple}} = \mathbb{E}_{x_0, \epsilon_t} \left[ \left\| \epsilon_t - \epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \epsilon_t, t) \right\|^2 \right]$$

\* Instead of modelling  $\mu_\theta$  directly we model  $\epsilon_\theta$  that predicts the noise  $\epsilon_t$ .



## \* Diffusion Model for Hand Writing generation.

→ The dataset used in this task will be a bit unique. it comprises of 2 components:

- Hand writing pen-strokes in XML form.
- Text-sequences describing what is written.
- Images specifying handwriting styles.

### \* Pen-strokes.

→ Every datapoint  $x_0 \in \mathbb{R}^2 \times \{0, 1\}$  comprises of  $\{x_1, \dots, x_N\}$  vectors.

Each  $x_t$  contains  $[(x_i, y_i), d_i]$  <sup>①</sup> where  $(x_i, y_i)$  specifies the offset coordinates w.r.t previous stroke. <sup>②</sup>  $d_i$  is a binary value specifying if the pen was down ( $d_i = 0$ ) when writing the stroke.

- We cannot parameterize  $d_i$  as gaussians → We separate  $x_t$  into  $y_t$  and  $d_t$ . Where  $y_t$  specifies the strokes &  $d_t$  specifies if pen is down or nah.
- We specify training loss for both penstrokes and penlifts as:

$$\textcircled{1} L_{\text{stroke}}(\theta) = \| \epsilon_t - \epsilon_\theta(y_t; c, s, \sqrt{\alpha_t}) \|_2^2 \rightarrow L_2 \text{ loss}$$

$$\textcircled{2} L_{\text{draw}}(\theta) = -d_t \log(\hat{d}_t) - (1 - d_t) \log(1 - \hat{d}_t) \rightarrow \text{BCE}$$

→ Where  $\hat{d}_t = d_\theta(y_t; c, s, \sqrt{\alpha_t})$  is the prediction whether pen is down or nah.

### \* Model architecture:

