

Generative Model:

Sample new points from a model which learns the true distribution based on random samples from that true distribution

Variational Autoencoder design that  $p_{\theta}(x|z)$

where latent variable  $z \sim N(0; I)$

why Gaussian?

1. Consistency among transformation

2. it's the limit of many distributions

covariance matrix is  $I$

Evidence Lower Bound (ELBO, Variational Lower Bound)

Objective: Get  $\log P_{\theta}(x)$

$$\log P_{\theta}(x) = \log \int_z p_{\theta}(x|z) dz \quad | \nu(z) \text{ is an arbitrary distribution}$$

$$= \log \int_z p_{\theta}(x|z) \frac{\mu(z)}{\mu(z)} dz$$

$$= \log E_{z \sim \nu(\cdot)} \frac{p_{\theta}(x|z)}{\mu(z)}$$

According to Jensen's Inequality, since  $\log(x)$  is concave,

$$\geq E_{z \sim \nu(\cdot)} \log \frac{p_{\theta}(x|z)}{\mu(z)} \quad \begin{matrix} n \\ \prod_{i=1}^n p_{\theta}(z_i|z_{i-1}) \end{matrix}$$

$$= E_{z \sim \nu(\cdot)} \log \frac{p_{\theta}(x|z) \mu(z)}{\mu(z)} \quad \begin{matrix} n \\ \prod_{i=1}^n \mu(z_i|z_{i-1}) \end{matrix}$$

$$= E_{z \sim \nu(\cdot)} [\log p_{\theta}(x|z) + \log \frac{\mu(z)}{p_{\theta}(z)}]$$

$$= E_{z \sim \nu(\cdot)} \log p_{\theta}(x|z) - E_{z \sim \nu(\cdot)} \log \frac{\mu(z)}{p_{\theta}(z)}$$

$p_{\theta}(z) \sim N(0, I)$

Reconstruction Loss when...

sampling  $N(\text{Decoder}(z), I) \rightarrow \| \hat{x} - x \|_2^2$  MSE  
from

$D_{KL}(\nu(\cdot) || p_{\theta}(\cdot))$

Kind of a regularization KOKUYO

To understand how it learns both mapping of  $x \rightarrow z$  and  $z \rightarrow x$

$$\text{The lower bound } E_{z \sim p(x)} \log \frac{p_\theta(z|x) p_\theta(x)}{p(x)}$$

$$= E_{z \sim p(x)} \log p_\theta(x) + E_{z \sim p(x)} \log \frac{p_\theta(z|x)}{p(z)}$$

not related to  $\theta, x$ , so basically identity here

$$= E_{z \sim p(x)} \log p_\theta(x) - D_{KL}(p(z) || p_\theta(z|x))$$

original object hard to directly compute, so we use

$$\text{For VAE } p(x) = q_\phi(z|x)$$

$$= \begin{cases} \theta e^{-\theta z}, & z \geq 0 \\ 0, & z < 0 \end{cases} = \theta e^{-\theta z} I(z \geq 0)$$

$$\sim N(\hat{\mu}(x), \hat{\sigma}^2(x)) \text{ or } N(x; \hat{\mu}, \hat{\sigma}^2 I)$$

what your network computes

$$z = \hat{\mu}(x) + \hat{\sigma}^2(x) \varepsilon, \quad \varepsilon \sim N(0, I)$$

Reparameterization transfer sampling into non-random  $z$   
so that we can use gradient

\* DKL between two Gaussian Distributions  $D_{KL}(u||v)$

$$u: N(\mu_1, \sigma_1^2) \quad v: N(\mu_2, \sigma_2^2)$$

$$D_{KL}(u||v) = E_{x \sim u(x)} \log \frac{u(x)}{v(x)}$$

$$= E_{x \sim u(x)} \log u(x) - \log v(x)$$

$$= E_{x \sim u(x)} \log \left[ \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \right] - \log \left[ \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \right]$$

$$= E_{x \sim u(x)} \log \left( \frac{\sigma_2}{\sigma_1} \right) + \frac{1}{2\sigma_2^2} (\mu_2 - \mu_1)^2 + \frac{1}{2\sigma_1^2} (x - \mu_1)^2$$

$$\therefore E_{x \sim u(x)} (x - \mu_1)^2 = \sigma_1^2$$

$$E_{x \sim u(x)} (x - \mu_2)^2 = E_{x \sim u(x)} (x - \mu_1 + \mu_1 - \mu_2)^2$$

$$= \sigma_1^2 + (\mu_1 - \mu_2)^2 + 0$$

$$\therefore D_{KL}(u||v) = \log \left( \frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \quad (\text{closed-form solution})$$

Likelihood function:  $L(\theta|x) = p(x|\theta)$

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

- ①  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$
- ②  $D_{KL}(P||Q) \geq 0$ ,  $= 0$  only when  $P=Q$

Diffusion Model - DDPM

With latent variable  $z_1, z_2, \dots, z_n$ ,

we want to generate data with Markov chain.

$$p_\theta(x_1, z_1, \dots, z_n) = p_\theta(z_n) \prod_{i=1}^{n-1} p_\theta(z_i|z_{i+1}) p_\theta(x_i|z_i)$$

similar to sample n times with VAE

$$\text{where } p_\theta(z_{i+1}|z_i) = N(z_{i+1}; \mu_\theta(z_i, t), \Sigma_\theta(z_i, t))$$

forward  $\rightarrow$  Human predefined mapping such configuration makes scaling stable  
 $q(z_i|x, z_1, \dots, z_{i-1}) = q(z_i|z_{i-1}) = N(z_i; \frac{a_i}{a_{i-1}} z_{i-1}, (1 - \frac{a_i}{a_{i-1}})^2 I)$

backward  $\rightarrow$   
Since  $q(z_{i-1}|z_i, x) = \frac{q(z_{i-1}|z_i)}{q(z_i|x)} = \frac{q(z_i|z_{i-1}, x) q(z_{i-1}|x) q(x)}{q(z_i|x) q(x)}$   
 $\propto \exp\left(-\frac{1}{2}\left(\frac{(z_i - \bar{z}_{i-1} - \bar{x})^2}{\beta_i} + \frac{(z_{i-1} - \bar{z}_{i-1} - \bar{x})^2}{1-a_{i-1}} - \frac{(z_i - \bar{z}_{i-1} - \bar{x})^2}{t-a_i}\right)\right)$

$$= \exp\left(-\frac{1}{2}\left(\frac{a_i}{\beta_i} + \frac{1}{1-a_{i-1}}\right) z_{i-1}^2 - \left(\frac{2\bar{z}_i}{\beta_i} z_i + \frac{2\bar{z}_{i-1}}{1-a_{i-1}} x\right) z_{i-1} + C(z_i, x)\right)$$

Variance of  $z_{i-1}$       mean of  $z_{i-1}$       unrelated

which can be deduced in a form  $\frac{1}{2}a(x+\frac{b}{2a})^2 + C$  where  $a = -\frac{b}{2a}$ ,  $\sigma^2 = \frac{1}{a}$ .  
therefore,  $\cong \exp\left(-\frac{1}{2} \frac{(z_{i-1} - \hat{\mu}(z_i, x))^2}{\sigma^2}\right)$   $\hat{\mu}$  is linear

$$\bar{\sigma}^2 = \frac{1}{\alpha_i + \frac{1}{1-\alpha_{i-1}}} = \frac{1-\bar{\alpha}_{i-1}}{1-\bar{\alpha}_i} \cdot \beta_i$$

$$\bar{\mu}_i(\bar{x}_i, x) = \frac{\bar{\alpha}_i(1-\bar{\alpha}_{i-1})}{1-\bar{\alpha}_i} \bar{x}_i + \frac{\sqrt{\bar{\alpha}_{i-1}} \beta_i}{1-\bar{\alpha}_i} x = \frac{1}{\bar{\alpha}_i} (\bar{x}_i - \frac{1-\bar{\alpha}_i}{\sqrt{1-\bar{\alpha}_i}} \varepsilon_i) \quad (2)$$

$\therefore$  with given  $x$ , mean  $\bar{\mu}$  is a function of hyperparam,  $\underline{x}_i, \underline{\varepsilon}_i$  (linear)  
variance  $\bar{\sigma}^2$  is a function of hyperparam

Two tricks

① Reparameterization

To make sampling from Gaussian differentiable,

for  $\underline{x} \sim N(\underline{x}, \underline{\mu}_\theta, \underline{\sigma}_\theta^2 \underline{I})$ , we let

$$\underline{x} = \underline{\mu}_\theta + \underline{\sigma}_\theta \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, \underline{I})$$

randomness is transferred from  $\underline{x}$  to  $\underline{\varepsilon}$

②  $\underline{x}_i$  is a function of  $x$  and  $\beta$  ( $\alpha_i = 1 - \beta_i$ )

$$\underline{x}_i = \sqrt{\bar{\alpha}_i} \underline{x}_{i-1} + \sqrt{1-\bar{\alpha}_i} \underline{\varepsilon}_{i-1}$$

$= \dots$

$$= \sqrt{\bar{\alpha}_i} \underline{x}_0 + \sqrt{1-\bar{\alpha}_i} \underline{\varepsilon} \quad (N(0, \bar{\sigma}_1^2 \underline{I}) + N(0, \bar{\sigma}_2^2 \underline{I}) = N(0, (\bar{\sigma}_1^2 + \bar{\sigma}_2^2) \underline{I}))$$

therefore  $q(\underline{x}_i | x) = N(\underline{x}_i; \sqrt{\bar{\alpha}_i} x, (1-\bar{\alpha}_i) \underline{I})$

$\sqrt{1-\bar{\beta}_i}$  makes sure variance  $\xrightarrow{n \rightarrow \infty} 1$

$$q(\underline{x}_{i-1} | \underline{x}_i, x) = N(\underline{x}_{i-1}, \hat{\mu}_i(\underline{x}_i, x), \bar{\sigma}^2 \underline{I})$$

From ELBO in VAE, we know

$$\begin{aligned}
 \log p_{\theta}(x) &\geq \mathbb{E}_{z \sim p_{\theta}(\cdot)} \log \frac{p_{\theta}(x|z) p_{\theta}(z)}{p_{\theta}(z_i)} \\
 &\quad \checkmark | \rightarrow q \\
 \therefore -\log p_{\theta}(x) &\leq \mathbb{E}_q [-\log p_{\theta}(z_n) + \sum_{i=1}^n \log \frac{q(z_i|x_{i-1})}{p_{\theta}(z_{i-1}|z_i)}] \\
 &= \mathbb{E}[-\log p_{\theta}(z_n) + \sum_{i=2}^n \log (\frac{q(z_i|x_{i-1})}{p_{\theta}(z_{i-1}|z_i)}) + \log \frac{q(z_1|x)}{p_{\theta}(x|z_1)}] \\
 &\quad \text{add } x \text{ to make this computable} \\
 &= \mathbb{E}[-\log p_{\theta}(z_n) + \sum_{i=2}^n \log (\frac{q(z_{i-1}|z_i)x) q(z_i|x)}{q(z_{i-1}|x)p_{\theta}(z_{i-1}, z_i)}) + \log \frac{q(z_1|x)}{p_{\theta}(x|z_1)}] \\
 &= \mathbb{E}[-\log p_{\theta}(z_n) + \sum_{i=2}^n \log (\frac{q(z_{i-1}|z_i, x)}{p_{\theta}(z_{i-1}|z_i)}) + \log \frac{q(z_n|x)}{\cancel{q(z_i|x)}} + \log \frac{q(z_1|x)}{p_{\theta}(x|z_1)}] \\
 &\quad \cancel{\text{cancel out like } \frac{n}{n} \cdot \frac{n}{n-1} \dots \frac{2}{1}} \\
 &= \mathbb{E}[\log \frac{q(z_n|x)}{p_{\theta}(z_n)} + \sum_{i=2}^n \log (\frac{q(z_{i-1}|z_i, x)}{p_{\theta}(z_{i-1}|z_i)}) - \log p_{\theta}(x|z_1)] \\
 &= \mathbb{E}[D_{KL}(q(z_n|x) || p_{\theta}(z_n)) + \sum_{i=2}^n D_{KL}(q(z_{i-1}|z_i, x) || \cancel{p_{\theta}(z_{i-1}|z_i)}) - \log p_{\theta}(x|z_1)] \\
 &\quad \uparrow \quad \nwarrow \text{NLL} \quad \text{pure noise} \quad \text{used for generation} \\
 &= \mathbb{E}[D_{KL}(q(z_n|x) || p_{\theta}(z_n)) + \sum_{i=2}^n D_{KL}(q(z_{i-1}|z_i, x) || \cancel{p_{\theta}(z_{i-1}|z_i)}) - \log p_{\theta}(x|z_1)] \\
 &\quad \uparrow \quad \text{make it learn without} \quad \uparrow \\
 &\quad \text{we don't optimize this, treat as constant} \quad \text{data} \propto \quad \text{simple reconstruct}
 \end{aligned}$$

Recall if  $p, q$  are both Gaussian.

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

By design, the model learn from a bunch of latent variables to match  $q$  and  $p$

Now that 1.  $q(z_{i-1} | z_i, \gamma) = N(z_{i-1}, \hat{\mu}_i(z_i, \gamma), \tilde{\sigma}^2 I)$

$$\text{where } \hat{\mu}_i = \frac{1}{\bar{\alpha}_i} (z_i - \frac{1-\bar{\alpha}_i}{J_{i-1}} \varepsilon_i)$$

$$\tilde{\sigma}^2 = \frac{1-\bar{\alpha}_{i-1}}{1-\bar{\alpha}_i} \beta_i$$

2.  $p_\theta(z_{i-1} | z_i) = N(z_{i-1}; \mu_\theta(z_i, i), \Sigma_\theta(z_i, i))$

$$\therefore \text{Loss} = \text{Eq}[D_{KL}(q(z_{i-1} | z_i, \gamma) || p(z_{i-1} | z_i))]$$

discarded since unrelated

$$\begin{aligned} \text{SNR} &= \frac{\alpha_i^2}{\tilde{\sigma}^2} \in [1-\bar{\alpha}_i^2] \text{ in ddpm} \\ \text{noise rate} &= \text{Eq}[\|\hat{\mu}_i(z_i, \gamma) - \mu_\theta(z_i, i)\|_2^2] \times \frac{1}{2\tilde{\sigma}^2} + C \\ &= \frac{1}{2} E_\varepsilon [(SNR(i-1) - SNR(i)) \|x - \hat{x}_\theta(z_i, i)\|_2^2] \end{aligned}$$

match  $x$  and predicted  $x$

$$\text{we then define } \mu_\theta(z_i, i) = \frac{1}{\bar{\alpha}_i} (z_i - \frac{1-\bar{\alpha}_i}{J_{i-1}} \varepsilon_\theta(z_i, i))$$

*we predict the noise*

$$\text{Objective: } E_{x, \varepsilon} [\|\varepsilon - \varepsilon_\theta(z_i, i)\|^2] \quad \varepsilon \sim N(0, 1)$$

What if infinite  $z_i$  / continuous time?

$$SNR'(t) = \frac{dSNR(t)}{dt}$$

$$\therefore \text{Objective becomes: } -\frac{1}{2} E_\varepsilon \left[ \int_0^1 SNR'(t) \|x - \hat{x}_\theta(z_t, t)\|_2^2 dt \right]$$

SDE perspective:

drift

Weiner motion (random)

$$\text{Forward: } dx = f(x, t)dt + g(t)d\omega \quad \text{steps in vector space, describe}$$

$$\text{Reverse: } dx = [f(x, t) - g^2(t) \nabla_x \log p_\theta(x)] dt + g(t) d\tilde{\omega} \quad \begin{matrix} \text{how distribution} \\ \text{is changing} \end{matrix}$$

Using SDE solver gives us more accurate prediction than discrete version

Score of distribution

(gradient of log prob)

Campus we can view ddpm in a da perspective or score match perspective