Max prob $p(x_0)$

marginalize $\int p(x_0, x_{1:T}) \, dx_{1:T}$

$x_1, x_2 \ldots x_T$ latent variables

Suppose $x_0 \sim \mathcal{D}$ $\qquad \mathcal{D}$ distribution of

$\underbrace{3 \times 64 \times 64}$ images

$d = 3 \times 2^{12} = 12,288$

$$x_T \to x_{T-1} \to \dots \to x_2 \to x_1 \to x_0$$

reverse process to learn
denoise

$$x_T \sim N(0, I^d)$$

$x_0 \sim \mathcal{D}$

given $x_0$ forward process
add noise per noise schedule

$$\beta_1, \beta_2 \dots \beta_T$$

$1e^{-4} \dots$ linear $\dots$ .02

posterior distributions, given $x_0$

$$q(x_t \mid x_{t-1}, x_0) \sim N(x_t; \sqrt{1-\beta_t}\, x_{t-1}, \beta_t I)$$

$$x_1 = \sqrt{1-\beta_1}\, x_0 + \sqrt{\beta_1}\, \varepsilon$$

$x_T \to N(0,1)$
per noise schedule
not exactly

$$var(x_1) = \beta_1$$

Not only $q(x_t | x_{t-1})$ normal, but

$$q(x_t | \boxed{x_0}) \text{ normal}$$

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

and

$$\sim N(x_t, \sqrt{\bar{\alpha}_t} \boxed{x_0} (1-\bar{\alpha}_t) I)$$

only depends on $x_0$

$$\alpha_t = 1 - \beta_t$$

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

Can apply Bayes' Theorem to

$$q(x_t \mid x_{t-1}, x_0) \quad t > 1 \quad \left( \begin{array}{c} q(x_1 \mid x_2, x_0) \\ \text{and onward} \end{array} \right)$$

and it is normal,

$$q(x_{t-1} \mid x_t, x_0) =$$

$$N\left(x_{t-1} \mid c_t x_0 + \gamma_t x_t, \tilde{\beta}_t I\right)$$

$$c_t = \frac{\sqrt{\bar{\alpha}_{t-1}} \, \beta_t}{1 - \bar{\alpha}_t}$$

$$\gamma_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

$$\tilde{\beta}_t = \left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\right)\beta_t$$

Model reverse process

$$p_\theta(x_{t-1} | x_t) \sim N(x_{t-1} | u_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

Simplify as
$$\sim N(x_{t-1} | u_\theta(x_t, t), \beta_t I)$$

assume Markov property

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} p_\theta(x_{t-1} | x_t)$$

likewise for posteriors

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

# Deriving Loss, reducing Variance part 1

want model to maximize $p_\theta(x_0)$

same as minimize NLL $-\ln p_\theta(x_0)$

**marginalize over latent $x_{1:T}$**

$$-\ln p_\theta(x_0) = -\ln \int p_\theta(x_{0:T}) \, dx_{1:T}$$

**bring in posterior**

$$= -\ln \int p_\theta(x_{0:T}) \frac{g(x_{1:T}|x_0)}{g(x_{1:T}|x_0)} \, dx_{1:T}$$

# Deriving Loss, reducing Variance part 2

$$-\ln p_\theta(x_0) = -\ln \int p_\theta(x_{0:T}) \frac{g(x_{1:T}|x_0)}{g(x_{1:T}|x_0)} dx_{1:T}$$

$$= -\ln \underset{x_{1:T} \sim g(\cdot|x_0)}{E} \left( \frac{p_\theta(x_{0:T})}{g(x_{1:T}|x_0)} \right)$$

<span style="color:green">can sample from $g(\cdot|x_0)$</span>

<span style="color:green">apply Jensen's inequality</span>

<span style="color:orange">(sim to deriving ELBO)</span>

$$\leq \underset{x_{1:T} \sim g(\cdot|x_0)}{E} -\ln \frac{p_\theta(x_{0:T})}{g(x_{1:T}|x_0)}$$

# Deriving Loss, reducing Variance part 3

$$loss = \mathop{E}_{x_{1:T} \sim g(\cdot | x_0)} -\ln \frac{P_\theta(x_{0:T})}{g(x_{1:T} | x_0)}$$

$$= \mathop{E}_{x_{1:T} \sim g(\cdot | x_0)} -\ln P_\theta(x_T) - \sum_{t \geq 1} \ln \left( \frac{P_\theta(x_{t-1} | x_t)}{g(x_t | x_{t-1})} \right)$$

$$= \mathop{E}_{x_{1:T} \sim g(\cdot | x_0)} -\ln P_\theta(x_T) - \ln \frac{P_\theta(x_0 | x_1)}{g(x_1 | x_0)} - \sum_{t \geq 2} \ln \left( \frac{P_\theta(x_{t-1} | x_t)}{g(x_t | x_{t-1})} \right)$$

# Deriving Loss, reducing Variance part 4

$$\mathop{E}_{x_{1:T} \sim q(\cdot | x_0)} \left[ -\ln P_\theta(x_T) - \ln \frac{P_\theta(x_0 | x_1)}{q(x_1 | x_0)} - \sum_{t \geq 2} \ln \left( \frac{P_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right) \right]$$

apply Bayes

$$= \mathop{E}_{x_{1:T} \sim q(\cdot | x_0)} \left[ -\ln P_\theta(x_T) - \ln \frac{P_\theta(x_0 | x_1)}{q(x_1 | x_0)} - \sum_{t \geq 2} \ln \frac{P_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t)} \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \right]$$

# Deriving Loss, reducing Variance part 5

$$= \mathop{E}_{x_{1:T} \sim g(\cdot | x_0)} -\ln p_\theta(x_T) - \ln \boxed{\frac{p_\theta(x_0 | x_1)}{g(x_1 | x_0)}} - \sum_{t \geq 2} \ln \frac{p_\theta(x_{t-1} | x_t)}{g(x_{t-1} | x_t)} \boxed{\frac{g(x_{t-1} | x_0)}{g(x_t | x_0)}}$$

these terms form telescoping sum

$t=2$

$+ \ln g(x_1 | x_0) \; -\ln g(x_1 | x_0)$

$+ \ln g(x_2 | x_0)$

$t=3 \; \ldots$

$t=T$

$-\ln g(x_{T-1} | x_0)$

$+\ln g(x_T | x_0)$

only left with this

$$= \mathop{E}_{x_{1:T} \sim g(\cdot | x_0)} -\ln \frac{p_\theta(x_T)}{g(x_T | x_0)} - \ln p_\theta(x_0 | x_1) - \sum_{t \geq 2} \ln \frac{p_\theta(x_{t-1} | x_t)}{g(x_{t-1} | x_t)}$$

# Deriving Loss, reducing Variance part 6

$$= \mathop{E}_{x_{1:T} \sim g(\cdot|x_0)} -\ln \frac{P_\theta(x_T)}{g(x_T|x_0)} - \ln P_\theta(x_0|x_1) - \sum_{t \geq 2} \ln \frac{P_\theta(x_{t-1}|x_t)}{g(x_{t-1}|x_t)}$$

$$= \mathop{E}_{x_1} -\ln P_\theta(x_0|x_1) + \sum_{t \geq 2} KL\left[g(x_{t-1}|x_t) | P_\theta(x_{t-1}|x_t)\right] + KL \, g(x_T$$

$$= \mathop{E}_{x_1} -\ln P_\theta(x_0|x_1) + \sum_{t \gg 2} KL\left[g(x_{t-1}|x_t) | P_\theta(x_{t-1}|x_t)\right] + KL\left[g(x_T|x_0) | p(x_T)\right]$$

$L_0$ loss
term, reconstruction
could do simple
decoder

two normals
analytic formula
weighted MSE on means

constant
can ignore
for training

# Model Noise

$$\text{loss} = \mathop{E}_{x_1} -\ln p_\theta(x_0|x_1) + \sum_{t \geq 2} KL\left[g(x_{t-1}|x_t) \| p_\theta(x_{t-1}|x_t)\right]$$

mean of $g(x_{t-1}|x_t)$

$$L_{t-1} = \mathop{E}_{g} \frac{1}{2\sigma_t^2} \left\| \tilde{u}_t(x_t, x_0) - u_\theta(x_t, t) \right\|^2$$

latent $x_t = c_1 x_0 + c_2 \varepsilon$

where $\varepsilon \sim N(0, I)$

so $x_0 = \tilde{c}_1 x_t + \tilde{c}_2 \varepsilon$

and $\tilde{u}_t = \tilde{\tilde{c}}_1 x_t + \tilde{\tilde{c}}_2 \varepsilon$

so set $u_\theta(x_t, t) = \tilde{\tilde{c}}_1 x_t + \tilde{\tilde{c}}_2 \varepsilon_\theta(x_t, t)$

# Model Noise

$$\text{loss} = \underset{x_1}{E} \underbrace{-\ln p_\theta(x_0|x_1)}_{L_0} + \sum_{t>2} \underset{x_0,\epsilon}{E} \underbrace{\omega_t \left\| \epsilon - \epsilon_\theta\left(\sqrt{\bar{\alpha}_t}\, x_0 + \sqrt{1-\bar{\alpha}_t}\,\epsilon, t\right) \right\|^2}_{L_{t-1}}$$

$$\epsilon \sim N(0, I)$$

$$\omega_t = \frac{B_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)}$$

simple alg 1

ignore $L_0$

$$\omega_t = 1$$