

Let $x \in \mathbb{R}^d$ be an image and D be our image data
 $p_\theta(x)$ is our data distribution or $p_\theta(x_t)$ the probability of
a single image

Training objective is maximize $p_\theta(x_t)$

$$\max_{\theta} \mathbb{E}_{x_t \sim D} [\log p_\theta(x_t)], D = \text{data}$$

$p_\theta(x_t) = \int p_\theta(x_{0:T}) dx_{1:T}$ the marginalization over all
intermediate latents x_1, \dots, x_T
introducing

We observe x_0 an image, all intermediate $x_{1:T}$ are latent.

"The probability of seeing x_0 is the total probability of all
possible reverse trajectories that end at x_0 , weighted by
how likely each trajectory is under the model"

What is $p_\theta(x_{0:T})$?

$$p_\theta(x_{0:T}) = p(x_T) p(x_{T-1}|x_T) p(x_{T-2}|x_{T-1}, x_T) \dots p(x_1|x_{1:T})$$

by chain rule

$$= p(x_T) \prod_{t=1}^T p(x_{t-1}|x_t) \quad \text{by Markov Assumption}$$

$$\Rightarrow p_\theta(x_0) = \int p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) dx_{1:T}$$

The true reverse kernel $p(x_{t-1}|x_t)$ is unknown

"Given a noisy latent x_t , what is the distribution of the
slightly less noisy latent x_{t-1} in the true data generating
process?" - we don't know. So we

introduce known $g(x_{1:T}|x_0)$, the forward process

$$g(x_{1:T}|x_0) = \prod_{t=1}^T g(x_t|x_{t-1}) \quad \text{we don't know this yet}$$

$$g(x_t|x_{t-1}) = N(x_t; f_\alpha(x_0, (I - \bar{\alpha}_t)I))$$

why do we introduce $g(x_{1:T} | x_0)$?

we need a tractable way to sample latent trajectories (x_1, \dots, x_T) given a real data sample x_0 . That's what the forward noising process $g(x_{1:T} | x_0)$ is for

$$g(x_{1:T} | x_0) = \prod_{t=1}^T g(x_t | x_{t-1}).$$
 It is known Gaussian kernel we design ourselves.

Deriving the forward process

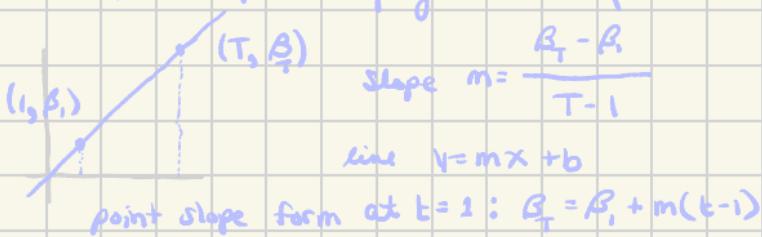
How do we get the parameters of the conditional distribution $g(x_t | x_0)$?

Let's start with $g(x_t | x_{t-1}) = N(a_t x_{t-1}, \beta_t I)$

where Mean = $a_t x_{t-1}$, the a_t is scaling the signal from the previous latent

What about β_t ?

DOPM : w/ the corruption to progress linearly



$$\therefore \beta_t = \beta_1 + \frac{t-1}{T-1} (\beta_T - \beta_1)$$

such $\mu = a_t x_{t-1}$ & $\sigma^2 = \beta_t I$

and using the sampling form of a Gaussian

$$x = \mu + \sigma \epsilon, \epsilon \sim N(0, I) \Rightarrow x_t = a_t x_{t-1} + \sqrt{\beta_t} \epsilon_t, \epsilon_t \sim N(0, I)$$

What is α_t ?

Now, let $\text{Var}(x_{t-1}) = 1$ and we choose to preserve unit variance at each timestep t

That means we want $\text{var}(x_t) = 1$

$$\begin{aligned}\Rightarrow \text{Var}(x_t) &= \text{Var}(\alpha_t x_{t-1} + \beta_t \epsilon_t) \\ &= \alpha_t^2 \text{Var}(x_{t-1}) + \beta_t \text{Var}(\epsilon_t) \\ &= \alpha_t^2 + \beta_t^2 = 1\end{aligned}$$

Now $\alpha_t^2 + \beta_t^2 = 1$

$$\Rightarrow \alpha_t = \sqrt{1 - \beta_t^2} \rightarrow$$

Define: $\alpha_t := 1 - \beta_t$

Then $\alpha_t = \sqrt{\alpha_t}$

$$\therefore g(x_t | x_{t-1}) = N(\sqrt{\alpha_t} x_{t-1}, \beta_t I)$$

Deriving closed form $g(x_t | x_0)$

We've just defined the one-step forward kernel $g(x_t | x_{t-1})$ which tells us: if we know the image at $t-1$ how do we sample the next noisy image x_t

That's useful for simulation, but during training we do something different. We sample α_t and jump straight from x_0 to x_t in one go. Why? So we can avoid all intermediate steps

In the training objective we want draw $x_t | x_0$ and schedule β_0, \dots, β_t then feed that to the network. Thus, mathematically we need $g(x_t | x_0)$

Luckily $g(x_t | x_0)$ comes directly from $g(x_t | x_{t-1})$ by unrolling the recursion and tracking how the scale and noise accumulate.

Let's unroll the recursion

$$\begin{aligned} x_3 &= \overline{\alpha}_3 (\overline{\alpha}_2 x_1 + \overline{\beta}_2 \epsilon_2) + \overline{\beta}_3 \epsilon_3 \\ &= \overline{\alpha}_3 \overline{\alpha}_2 x_1 + \overline{\alpha}_3 \overline{\beta}_2 \epsilon_2 + \overline{\beta}_3 \epsilon_3 \end{aligned}$$

$$\begin{aligned} x_3 &= \overline{\alpha}_3' x_2 + \overline{\beta}_3' \epsilon_3 \\ x_2 &= \overline{\alpha}_2' x_1 + \overline{\beta}_2' \epsilon_2 \\ x_1 &= \overline{\alpha}_1' x_0 + \overline{\beta}_1' \epsilon_1 \end{aligned}$$

$$\begin{aligned} &= \overline{\alpha}_3 \overline{\alpha}_2 (\overline{\alpha}_1 x_0 + \overline{\beta}_1 \epsilon_1) + \overline{\alpha}_3 \overline{\beta}_2 \epsilon_2 + \overline{\beta}_3 \epsilon_3 \\ &= \underbrace{\prod_{s=1}^3 \overline{\alpha}_s}_{(P)} x_0 + \underbrace{\overline{\alpha}_3 \overline{\alpha}_2 \overline{\beta}_1}_{(Q)} \epsilon_1 + \overline{\alpha}_3 \overline{\beta}_2 \epsilon_2 + \overline{\beta}_3 \epsilon_3 \end{aligned}$$

$$= \underbrace{\prod_{s=1}^t \overline{\alpha}_s}_{(P)} x_0 + \sum_{s=1}^t \overline{\beta}_s \left(\prod_{m=s+1}^t \overline{\alpha}_m \right) \epsilon_s$$

(P)

(Q)

$$P, \text{ our mean} \rightarrow \prod_{s=1}^t \overline{\alpha}_s = \sqrt{\prod_{s=1}^t \alpha_s} \equiv \bar{\alpha}_s$$

Property:

Linear combination of independent Gaussians is Gaussian

What about Q?

Since the $\epsilon_s \sim N(0, I)$ (independent Gaussians)

$$\Rightarrow \sum_{s=1}^t c_s \epsilon_s, \text{ with } c_s = \overline{\beta}_s \prod_{m=s+1}^t \overline{\alpha}_m \sim N(0, \sum_{s=1}^t c_s^2)$$

$$\text{Var}(x_t | x_0) = \sum_{s=1}^t c_s^2 = \sum_{s=1}^t \overline{\beta}_s \prod_{m=s+1}^t \overline{\alpha}_m$$

Can we get a better closed form for the variance?

$$\text{SDXL} \rightarrow g(x_t | x_0) \sim N(\bar{\alpha}_t x_0, (1 - \bar{\alpha}_t) I)$$

How do we get $\text{Var} = 1 - \bar{\alpha}_t$?

Recall

$$\alpha_u = 1 - \beta_u \rightarrow \beta_u = 1 - \alpha_u$$

and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$

$$\Rightarrow \sum_{s=1}^t \beta_s \prod \alpha_u = \sum_{s=1}^t (1 - \alpha_s) \prod_{u=s+1}^t \alpha_u$$

$$(1 - \alpha_s) \prod_{u=s+1}^t \alpha_u = \frac{t}{\prod \alpha_u} \alpha_u - \alpha_s \prod_{u=s+1}^t \alpha_u = \frac{t}{\prod \alpha_u} \alpha_u - \frac{t}{\prod \alpha_u}$$

$$\Rightarrow \sum_{s=1}^t \beta_s \prod \alpha_u = \sum_{s=1}^t \left(\frac{t}{\prod \alpha_u} \alpha_u - \frac{t}{\prod \alpha_u} \right)$$

$$\cancel{\frac{t}{\prod \alpha_u}} - \cancel{\frac{t}{\prod \alpha_u}} + \cancel{\frac{t}{\prod \alpha_u}} - \cancel{\frac{t}{\prod \alpha_u}} + \dots + \cancel{\frac{t}{\prod \alpha_u}} - \cancel{\frac{t}{\prod \alpha_u}} = 1 - \bar{\alpha}_t$$

$$t=3 \rightarrow \cancel{2} / \cancel{3} / \cancel{1} / \cancel{0} \\ \cancel{\frac{t}{\prod \alpha_u}}_{u=1} \quad \cancel{1} \rightarrow 1 - \bar{\alpha}_t$$

$$\therefore \text{var}(x_t | x_0) = \sum_{s=1}^t \beta_s \frac{t}{\prod_{u=s+1}^t \alpha_u} = 1 - \bar{\alpha}_t = 1 - \prod_{s=1}^t (1 - \beta_s)$$

$$\therefore g(x_t | x_0) = N(\bar{\alpha}_t x_0, (1 - \bar{\alpha}_t) I)$$

