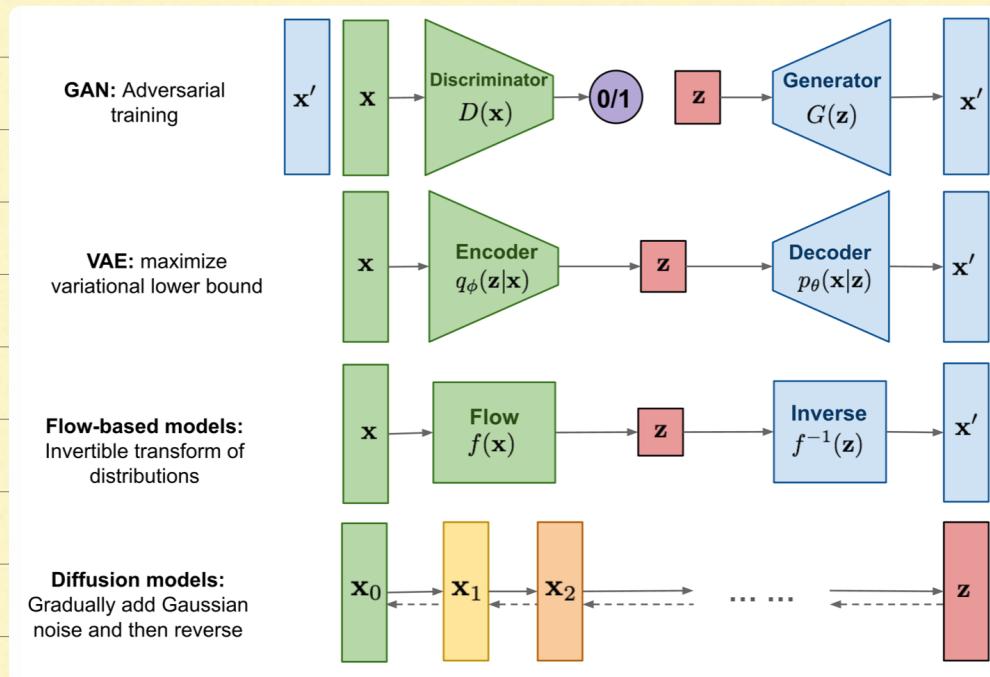


GAN, VAE, Flow-based generative model은 각각 단점이 있다.

GAN은 불안정하고, VAE는 surrogate loss의 미지수.

Flow-based는 reversible 인터프리테이션 architecture의 단점이 있다.

Diffusion model은 markov chain의 algorithm noise를 조정할 수 있는 reverse process를 사용한다



- forward forward process

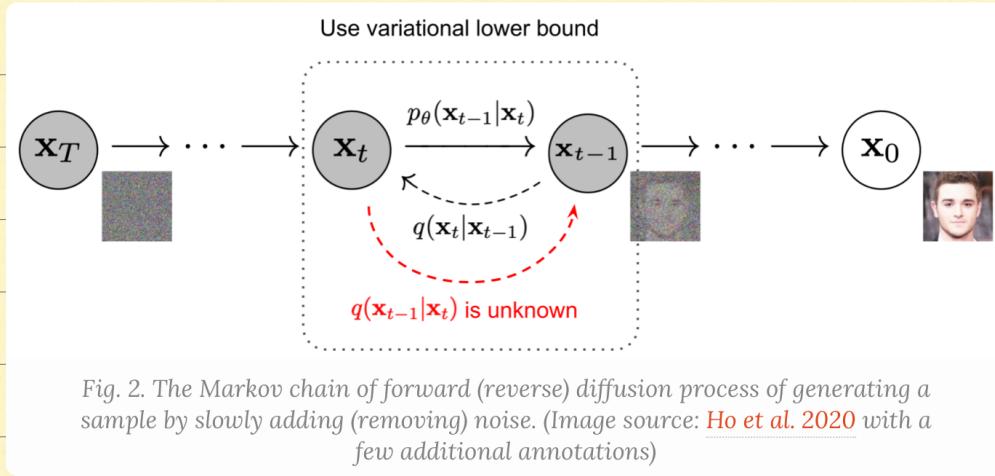
small Gaussian noise를 더하는 걸로 T steps로 수행할 때. x_t 의 noise sample

step size & Variance $\beta_t \in \{0, 1\}$ 를 갖는다.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

forward) \rightarrow 전이학습을 위한 input feature \rightarrow 자연 distinguish한 특성을 사용하기

isotropic gaussian distribution으로 변환.



⇒ reparameterization trick은对抗의, 원래의 단위Normal sample x_t 는

샘플링하는 데 사용. $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ 일 때.

$$x_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} z_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_t \alpha_{t-1}} z_{t-2}$$

⋮

$$= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} z$$

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t)I)$$

이제, noise가 짙어질수록 적용할 수 있는 noise 강도가 작아지며 때로이

$$\beta_1 < \beta_2 < \dots < \beta_T, \quad \bar{\alpha}_1 > \bar{\alpha}_2 > \dots > \bar{\alpha}_T$$

Markov chain update 때 $D_x \log P(x)$ 만을 이용하여

학습을 위한 $p(x)$ 와 sampling을 할 수 있는 방식이

SGLD (stochastic Gradient Langevin Dynamics) 이다.

<https://towardsdatascience.com/langevin-dynamics-29bbb9407b47>

$$x_t = x_{t-1} + \frac{\epsilon}{2} D_x P(x_{t-1}) + \sqrt{\epsilon} z_t, \quad z_t \sim \mathcal{N}(0, I), \quad \epsilon = \text{step size}$$

↳ local minimize collapse 방지를 위해 Gaussian noise를

update parameter of 넣었다.

- Reverse Diffusion process

$q(x_{t-1}|x_t)$ 은 0(정규분포) Gaussian Noise $\mathcal{N}(0, I)$ 로부터 true sample을 얻을 수 있다.

마지막 $q(x_{t-1}|x_t)$ 은 추적하는 모든 target space의 다른 모든 sample을

필로우하기 때문에 어렵다. 대신 $p_\theta(x_{t-1}|x_t)$ 를 이용하여 근사한다.

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

x_0 의 경우 reverse condition probability는 tractable하다.

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

$$\begin{aligned} q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\ &\propto \exp\left(-\frac{1}{2}\left(\frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1})^2}{\beta_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{1-\bar{\alpha}_t}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}^2 - \left(\frac{2\sqrt{\bar{\alpha}_t}}{\beta_t}\mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_t}}{1-\bar{\alpha}_t}\mathbf{x}_0\right)\mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0)\right)\right) \end{aligned}$$

$$\begin{aligned} \tilde{\beta}_t &= 1/\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right) = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \cdot \beta_t \\ \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \left(\frac{\sqrt{\bar{\alpha}_t}}{\beta_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_t}}{1-\bar{\alpha}_t}\mathbf{x}_0\right)/\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right) = \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 \end{aligned}$$

$$\begin{aligned} \tilde{\mu}_t &= \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\mathbf{z}_t) \\ &= \frac{1}{\sqrt{\bar{\alpha}_t}}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\mathbf{z}_t\right) \end{aligned}$$

$$\mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\mathbf{z}_t)$$

- Loss term

VLB (Variational lower bound)

$$\begin{aligned}
 -\log p_\theta(\mathbf{x}_0) &\leq -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \\
 &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})/p_\theta(\mathbf{x}_0)} \right] \\
 &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} + \log p_\theta(\mathbf{x}_0) \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\
 \text{Let } L_{\text{VLB}} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \geq -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0)
 \end{aligned}$$

(CE (Cross-entropy))

$$\begin{aligned}
 L_{\text{CE}} &= -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \right) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \right) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right) \\
 &\leq -\mathbb{E}_{q(\mathbf{x}_{0:T})} \log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \\
 &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = L_{\text{VLB}}
 \end{aligned}$$

Analytically computable <https://arxiv.org/pdf/1503.03585.pdf>

$$\begin{aligned}
 L_{\text{VLB}} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\
 &= \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left(\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \\
 &= \mathbb{E}_q \underbrace{[D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \| p_\theta(\mathbf{x}_T))]}_{L_T} + \underbrace{\sum_{t=2}^T D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0}
 \end{aligned}$$

$$L_{\text{VLB}} = L_T + L_{T-1} + \dots + L_0$$

where $L_T = D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))$

$$L_t = D_{\text{KL}}(q(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_0) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})) \text{ for } 1 \leq t \leq T-1$$

$$L_0 = -\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)$$

L_t 는 두 확률 분포의 비교로, closed form이다.

L_t 는 q 는 learnable parameter이고, \mathbf{x}_t 는 Gaussian Noise

이기 때문이 constant이다.

- Parameterization L_t for training Loss

$$P_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

$$\mu_\theta \Rightarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbf{z}_t \right)$$

\mathbf{x}_t 는 훈련시, input으로 치수가 가능하기 때문에, \mathbf{x}_t 로부터 \mathbf{z}_t 를

어떻게 reparameterize 할 수 있다.

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, t) \right)$$

$$\text{Thus } \mathbf{x}_{t-1} = \mathcal{N}(\mathbf{x}_{t-1}; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, t) \right), \Sigma_\theta(\mathbf{x}_t, t))$$

$\tilde{\mu}$ 의 차이를 minimize하기 위한 loss term L_t

$$\begin{aligned} L_t &= \mathbb{E}_{\mathbf{x}_0, \mathbf{z}} \left[\frac{1}{2\|\Sigma_\theta(\mathbf{x}_t, t)\|_2^2} \|\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{z}} \left[\frac{1}{2\|\Sigma_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbf{z}_t \right) - \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, t) \right) \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{z}} \left[\frac{\beta_t^2}{2\alpha_t(1-\bar{\alpha}_t)\|\Sigma_\theta\|_2^2} \|\mathbf{z}_t - \mathbf{z}_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{z}} \left[\frac{\beta_t^2}{2\alpha_t(1-\bar{\alpha}_t)\|\Sigma_\theta\|_2^2} \|\mathbf{z}_t - \mathbf{z}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\mathbf{z}_t, t)\|^2 \right] \end{aligned}$$

* Simplification

$$L_t^{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, \mathbf{z}_t} \left[\|\mathbf{z}_t - \mathbf{z}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\mathbf{z}_t, t)\|^2 \right]$$

$$L_{\text{simple}} = L_t^{\text{simple}} + C \xrightarrow{\text{constant}}$$

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ 
5:   Take gradient descent step on
       $\nabla_{\theta} \|\epsilon - \mathbf{z}_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t)\|^2$ 
6: until converged

```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \mathbf{z}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

Fig. 4. The training and sampling algorithms in DDPM (Image source: Ho et al. 2020)

- Connection with noise-conditioned score Networks (NCSN)

SCORE-based generative model은 Langevin dynamics에 의한

Sampling 및 score matching을 이용한 방식이다.

각 sample \mathbf{x} 의 pdf에 대한 score는 $D_x \log p(x)$ 로 정의된다.

score 훈련은 $S_{\theta}: \mathbb{R}^D \rightarrow \mathbb{R}^D$ 의 의해 이루어진다.

Sampling은 data가 고사영처럼 보이도록 대부분 저사용에 맞춰진다.

작지만 멀도가 낮은 조건에서는 정확도가 떨어지며 때로 예전의

noise와 함께 잡음을 놓고 score를 추정한다.

- Parameterization of β_t

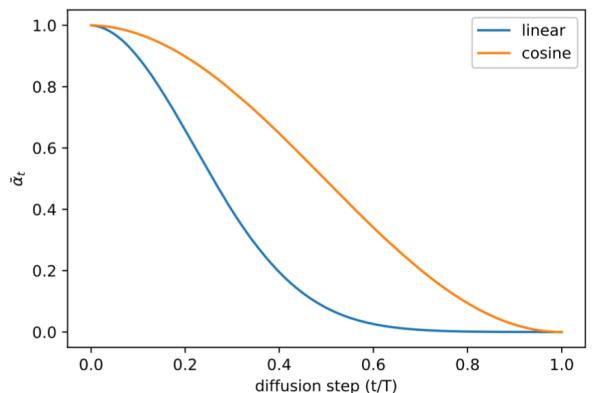
$\beta_t = 10^{-4}$, $\beta_T = 0.02$ 로 forward 시키는 설정 (DDPM)

↳ But, 성능이 조금 떨어짐

cosine-based variance schedule (improved DDPM)

$$\beta_t = \text{clip}\left(1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}}, 0.999\right) \quad \bar{\alpha}_t = \frac{f(t)}{f(0)} \quad \text{where } f(t) = \cos\left(\frac{t/T + s}{1+s} \cdot \frac{\pi}{2}\right)$$

⇒ DDPM보다 선형적이고 일의적이며



- Parameterization of Reverse process Variance Σ_θ

DDPM은 $\Sigma_\theta(x_t, t) = \beta_t^t I$ 이며 β_t 와 β_{t-1} 은 learnable 하지 않지만

$$\beta_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \text{은 설정했다}$$

\Rightarrow learnable한 Σ_θ 가 불안정하기 때문

improved DDPM은 Σ_θ 를 β_t 와 $\tilde{\beta}_t$ 사이의 interpolation으로 설정함

$$\Sigma_\theta(x_t, t) = \exp(v \log \beta_t + (1 - v) \log \tilde{\beta}_t)$$

v mixing vector

하지만 L_{simple} 의 Σ_θ 의 모든 parameter가 서로 충돌

$$L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{ULB}}$$

\hookrightarrow Σ_θ 학습할 때, Σ_θ 의 모든 학습을 멈춤

$L_{\text{gradient}} \text{의 noise} \text{ 떠올리} \text{ optimal} \text{ 를} \text{ 찾음}, \text{ time average smoothing}$

- Speed up diffusion sampling

Markov chain을 이용한 sampling으로 수천 step의 대량 sampling은 시간이 오래 걸린다.

$\Rightarrow 32 \times 32$ 의 5000 sampling이 20시간

1. Sampling은 $[T|s]$ step마다 하는 것. $T > s$

2.

For another approach, let's rewrite $q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ to be parameterized by a desired standard deviation σ_t according to the **nice property**:

$$\begin{aligned}\mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\mathbf{z}_{t-1} \\ &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\mathbf{z}_t + \sigma_t\mathbf{z} \\ &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}} + \sigma_t\mathbf{z} \\ q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2\mathbf{I})\end{aligned}$$

Recall that in $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$, therefore we have:

$$\tilde{\beta}_t = \sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t$$

Let $\sigma_t^2 = \eta \cdot \tilde{\beta}_t$ such that we can adjust $\eta \in \mathbb{R}^+$ as a hyperparameter to control the sampling stochasticity. The special case of $\eta = 0$ makes the sampling process deterministic. Such a model is named the *denoising diffusion implicit model (DDIM; Song et al., 2020)*. DDIM has the same marginal noise distribution but deterministically maps noise back to the original data samples.

During generation, we only sample a subset of S diffusion steps $\{\tau_1, \dots, \tau_S\}$ and the inference process becomes:

$$q_{\sigma, \tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_i}, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{\tau_{i-1}}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{\mathbf{x}_{\tau_i} - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2\mathbf{I})$$

\Rightarrow DDIM ($\eta=0$), DDPM ($\eta=1$)

\hookrightarrow small S when 좋은 결과

Compared to DDPM, DDIM is able to:

1. Generate higher-quality samples using a much fewer number of steps.
2. Have “consistency” property since the generative process is deterministic, meaning that multiple samples conditioned on the same latent variable should have similar high-level features.
3. Because of the consistency, DDIM can do semantically meaningful interpolation in the latent variable.

- Conditioned Generation

ImageNet 같은 복잡한 class의 이미지를 생성합니다.

복잡한 class 정보를 갖는 diffusion model을 확장,

improved DDPM은 noisy image x_t 에서 $f_\phi(y|x_{t,t})$ 를 사용하고,

target class y sampling을 위해 $\nabla_y \log f_\phi(y|x_{t,t})$ 를 사용합니다.

ADM (ablated diffusion model), ADM-G (additional classifier Guidance)

Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $f_\phi(y|x_t)$, and gradient scale s .

```

Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log f_\phi(y|x_t), \Sigma)$ 
end for
return  $x_0$ 
```

Algorithm 2 Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $f_\phi(y|x_t)$, and gradient scale s .

```

Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log f_\phi(y|x_t)$ 
     $x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$ 
end for
return  $x_0$ 
```

Fig. 8. The algorithms use guidance from a classifier to run conditioned generation with DDPM and DDIM. (Image source: [Dhariwal & Nichol, 2021](#))

improved DDPGM의 구조는 DDPGM의 block, depth / width 가 더 크고

더 많은 attention head etc. multi-resolution attention,

BigGAN의 residual, LPIPS의 residual connection rescale, AdaIN

- Summary

- pros

tractability와 flexibility는 원래 상호보수는 성능이다.

tractable은 불투명하지만, flexible(다양한 generative)하지 않지만,

diffusion model은 두 성능 모두 좋다

- cons

sample generate의 markov chain의 의존하는 경향

이 경향은 시간과 계산이 많이 듦다.

IDDPM은 GAN보다 많은 시간과 비용이 듦다

* paper

* Introduction.

energy-based model이 최근 Generative 분야에서 큰 성과

Diffusion probabilistic model을 제시. (이후 DM)

DM은 유한한 시간 후의 data sample로 초기의 sample을

풀기 위해 VE를 사용하여 훈련된 markov chain model이다.

$$VE \text{는 } p(x|z) = \frac{p(z|x)p(x)}{p(x)} \text{인 posterior } z \text{의 확률}$$

더 높은 확률 분포 $q(z)$ 로 근사화 | 유한 네트워크

Chain은 sampling이 반대되는 방향으로, Gaussian noise를 추가하는 것을 data가 처리될 때까지

하는 Markov chain인 DMe의 reverse를 학습한다

diffusion DM (Gaussian noise) 작다면, markov chain sampling을

conditional Gaussian으로 설정해도 되고, MN으로 parameterize할 가능성이 있다

ICA, DMe parameterize multi level modeling with denoising score match et

같다는 걸 알 수 있다.

↓
역할 내용

* Back ground.

DM은 $P_\theta(x_0) = \int P_\theta(x_{0:T}) dx_{1:T}$ 의 형태의 latent variable model이다.

$x_0 \sim q(x_0)$ 이며, $x_t \sim x_{t-1}$ 는 같은 dimension의 latent var.

$P_\theta(x_{0:T})$ 은 $P(x_T) = N(x_T; 0, I)$ 부터 시작하는 Gaussian transition을 흘려가는 markov chain으로, reverse process라고 한다.

$$P_\theta(x_{0:T}) = P(x_T) \prod_{t=1}^T P_\theta(x_{t-1}|x_t), \quad P_\theta(x_{t-1}|x_t) = N(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t I) \quad \text{variance} \quad \dots \quad (1)$$

이것은 VLB(Variational lower bound)을 이용한 NLL을 optim 하는 수법이다.

$$\mathcal{L} : E[-\log P_\theta(x_0)] \leq E_q[-\log \frac{P_\theta(x_{0:T})}{q(x_{1:T}|x_0)}] = E_q[-\log P(x_T) - \sum_{t \geq 1} \log \frac{P_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})}] \quad \dots \quad (2)$$

Variance β_t 는 trainable variable로 학습할 수 있고, constant로 정해 두면 된다.

↳ reverse process는 $P_\theta(\cdot|x_{t-1}, x_t)$ 로 표현할 수 있다.

↳ $\beta_t \rightarrow 0$ 하면 $q = p$

↑↑↑ forward process와 step training의 sample은 closed form으로 표현할 수 있다

$$\alpha_t = 1 - \beta_t, \quad \bar{\alpha}_t = \prod_{s=1}^t \alpha_s \text{ 일때}$$

$$q(x_t|x_0) = N(x_t; \overbrace{\bar{\alpha}_t x_0}^{\text{new}}, \overbrace{(1-\bar{\alpha}_t)I}^{\text{var}}) \quad \dots \quad (4)$$

∴ \mathcal{L} 는 SGD 학습으로, training이 이뤄지면, \mathcal{L} 는 다시 표현된다.

$$\mathcal{L} : E_q[\underbrace{D_{KL}(q(x_T|x_0) || p(x_T))}_{L_T} + \underbrace{\sum_{t \geq 1} D_{KL}(q(x_{t-1}|x_t, x_0) || P_\theta(x_{t-1}|x_t))}_{L_{t-1}} - \underbrace{\log P_\theta(x_0)}_{L_0}] \quad \dots \quad (5)$$

x_0 이 tractable한 $q(x_{t-1}|x_t, x_0)$ 과 P_θ 의 분포를 간단히 찾음

$$q(x_{t-1}|x_t, x_0) = N(x_{t-1}; \tilde{M}_t(x_t, x_0), \tilde{\beta} I), \quad \leftarrow \text{sampling} \quad \dots \quad (6)$$

$$\therefore \tilde{M}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_t} \beta_t}{1 - \bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t, \quad \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad \dots \quad (7)$$

↳ Rao-Blackwell theorem은 closed form으로 계산

- Rao - Blackwell theorem

모든 조정시, 충분한 data에 M conditional expectation을 취함으로써,

estimator의 효율성을 향상시킨다.

예를 들어, $\Theta = (\theta_1, \dots, \theta_q)$ 가 unknown vector인 때, $X = (X_1, \dots, X_n)$ 는

학률분포 $P(X, \Theta)$ 에서의 random sample이라고 할 때,

Θ 의 대수는 estimator $t(X) = (t_1(X), \dots, t_q(X))$, $C(t) = (C_{ij})$, $i, j = 1 \sim q$

$$\hookrightarrow C_{ij} = E[(t_i(X) - \theta_i)(t_j(X) - \theta_j)]$$

$E(t|S) = T(S)$ 가 Θ 의 독립적이도록, S 를 vector valued statistic이라 칭한다.

이때, Rao - Blackwell theorem이 따르면

$C(t) - C(T)$ 은 PD matrix이다.

S 의 표본이 충분히 크다면, $T(S)$ 가 Θ 의 독립적이라는 것은 충족된다.

따라서, $E(t_r - \theta_r)^2 \geq E(T_r - \theta_r)^2$ 이고, t_r 를 T_r 로 대체하면 MSE가 증명에서

효율성이 향상된다.

또한 Φ 가 convex function이라면

$$E[\Phi(t_r - \theta_r)] \geq E[\Phi(T_r - \theta_r)], r = 1 \sim q \text{ 가 성립한다.}$$

X· estimator가 충분통제량이 아닐 경우, 충분통제량의 선수인 다른 estimator가 있고

ol estimator는 MSE의 관점에서 더 효율적이다.

* Diffusion model and denoising autoencoder

DM 구현에는 β_t 의 reverse architecture 및 Gaussian distribution의 parameterization이 필요하다.

이를 위해 DM의 denoising score matching을 명시적으로 연결한다.

- Forward process and L_T

β_t 는 learnable 시점 reparameterization 할 수 있지만, paper는 fixed constant로 써준다.

$\therefore L_T$ 은 constant이다. (noise가 일정한 분포로 따른다)

- Reverse process and $L_{1:T-1}$

$$P_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad 1 < t \leq T \text{ 인 경우}$$

1. $\Sigma_\theta(x_t, t) = \sigma_t^2 I$ 를 훈련되지 않은 time dependent constant로 설정한다.

실증적으로, $\sigma_t^2 = \beta_t$ 일 때와 $\sigma_t^2 = \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$ 일 때의 결과는 유사했다.

$x_0 \sim N(\alpha I)$ 은 optimal, x_t 는 sample of optimal

이는 coordinate wise unit variance data의 case reverse process entropy와 upper, lower bound이며
좌표별로 독립적으로 계산(책임집)

극단적인 선택이다.

2. $\mu_\theta(x_t, t)$ 은 L_t 로부터, 특정 parameterization을 통해 표현된다.

$$P_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I) \text{이며}$$

$$L_{t-1} = E_q \left[\frac{1}{2\sigma_t^2} \| \tilde{\mu}_\theta(x_t, x_0) - \mu_\theta(x_t, t) \|_F^2 \right] + C, \quad C \in \text{independent constant} \quad \cdots (8)$$

로 LLT를 수 있다.

따라서 μ_θ 를 가장 잘 parameterization하는 방법은 $\tilde{\mu}_\theta$ 를 예측하는 것이다.

Eq. 8은 eq. 4인 $x_t(x_0, \epsilon) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim N(0, I)$ 의 reparameterization이다

eq. 7을 적용함으로서 확장할 수 있다.

$$L_{t-1} - C = E_{x_0, \epsilon} \left[\frac{1}{2\delta_t} \left\| \tilde{M}_t \left(x_t(x_0, \epsilon), \frac{1}{\sqrt{\alpha_t}} (x_t(x_0, \epsilon) - \sqrt{1-\bar{\alpha}_t} \epsilon) \right) - M_\theta(x_t(x_0, \epsilon), t) \right\|^2 \right] \dots (9)$$

$$= E_{x_0, \epsilon} \left[\frac{1}{2\delta_t} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon \right) - M_\theta(x_t(x_0, \epsilon), t) \right\|^2 \right] \dots (10)$$

$\hookrightarrow M_\theta$ 는 $\frac{1}{\sqrt{\alpha_t}} (x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon)$ 을 예측해야 한다.

x_t 는 model의 input이 될 수 있다. (대문자)

$$M_\theta(x_t, t) = \tilde{M}_t \left(x_t, \frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_\theta(x_t)) \right) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right),$$

ϵ_θ 는 x_t 로부터의 t 를 예측한 function

즉 parameterization을 할 수 있다.

Algorithm 1 Training

- 1: **repeat**
- 2: $x_0 \sim q(x_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(0, I)$
- 5: Take gradient descent step on
 $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon, t)\|^2$
- 6: **until** converged

Algorithm 2 Sampling

- 1: $x_T \sim \mathcal{N}(0, I)$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $z \sim \mathcal{N}(0, I)$ if $t > 1$, else $z = 0$
- 4: $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z$
- 5: **end for**
- 6: **return** x_0

$x_{t-1} \sim p_\theta(x_{t-1} | x_t)$ 및 sampling은 $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$ 로 계산할 수 있다.

sampling process는 Algorithm 2는 데이터의 빛도의 간접된 gradient로 ϵ_θ 를 사용하는

Langevin dynamics와 유사하다.

Eq.11을 이용하여 Eq.10을 Simplify 할 수 있다.

$$\Rightarrow E_{x_0, \epsilon} \left[\frac{\beta_t}{2\delta_t \alpha_t (1-\bar{\alpha}_t)} \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon, t) \right\|^2 \right] \dots (11)$$

\hookrightarrow multi noise level tnm denoising score matching CL 대처이다.

Eq.12는 Eq.11의 variational bound인 유사하다.

\hookrightarrow denoising score matching과 유사한 확률적 optimization은

Algorithm 2

Langevin dynamics와 비슷한 finite-time sample chain의 marginal을 맞추기 위해

Variational Inference를 사용하는 것 같다.

즉, μ_θ 가 \tilde{m}_θ 를 근사하도록 reverse process를 훈련시킬 수 있고,

parameterization을 변경해서 σ 를 근사하도록 훈련시킬 수도 있다.

ϵ -parameterization은 효율적이지만 아직 $P_\theta(x_t | x_0)$ 의 또 다른 parameterization이기 때문에

Section 4.6.1M ϵ 와 \tilde{m}_θ 를 비교한다.

- Data scaling, reverse process decoder, and Lo

image의 range를 $-1 \sim 1$ 로 scaling step 사용한다

→ pixel별 이라는 뜻인듯...?
discrete log likelihood를 위해 $N(x_0; \mu_\theta(x_1, 1), \sigma^2, I)$ 로부터

$$P_\theta(x_0 | x_1) = \prod_{i=1}^D \frac{\delta_+(x_{0i})}{\delta_-(x_{0i})} N(x_0; \mu_\theta(x_1, 1), \sigma^2, I) dx$$

D는 data dimension

$$\left\{ \begin{array}{ll} \delta_+(x) = & \infty, \text{ if } x = 1 \\ & x + \frac{1}{255}, \text{ if } x < 1 \\ \delta_-(x) = & -\infty, \text{ if } x = -1 \\ & x - \frac{1}{255}, \text{ if } x > -1 \end{array} \right.$$

→ Sampling의 결과 $\mu_\theta(x_1, 1)$ 을 display

- Simplified training objective

위에서 정의된 reverse process의 decoder로 eq 12, 13으로 구성된 variational bound를

명확히 배운다. 훈련할 수 있다.

하지만 아래와 같이 하는게 더 간단함.

$$L_{\text{simple}}(\theta) = E_{t, x_0, \epsilon} \left[\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \|^2 \right], \quad t = 1 \sim T \text{ 사이}$$

$t=1$ 일 때는 eq 13과 비슷하고, $t \sim T-1$ 은 eq 12와 같다. T 는 fixed인 경우가 고려됨

정확한 동작은 Algorithm 1과 같다.

Eq.14는 (2)의 weight coefficient를 각각의 대응의 가중치로 VI의 비중을 정함

weighted VI라 할 수 있다.

이전 낮은 티미 대비 가중치가 낮아지는가 티가 높아질수록 어려운 task이기 때문

어려운 task이 더 깊증할 수 있게 하여, 더 좋은 성능을 이끌어낸다.

* Experiments:

$T = 1000$, $\beta_t = 10^{-4} \sim 0.02$ until linear decay 풀기

↳ forward or reverse가 거의 동일하도록 작게 설정

reverse process only, unmasked PixelCNN++ or UNet backbone w/ Group normalize etc 한다.

network에 transformer의 position encoding 같이 사용에 대한 고려

↳ Appendix B

- Sample quality

code length?

Table 1: CIFAR10 results. NLL measured in bits/dim.

Model	IS	FID	NLL Test (Train)
Conditional			
EBM [11]	8.30	37.9	
JEM [17]	8.76	38.4	
BigGAN [3]	9.22	14.73	
StyleGAN2 + ADA (v1) [29]	10.06	2.67	
Unconditional			
Diffusion (original) [53]			≤ 5.40
Gated PixelCNN [59]	4.60	65.93	3.03 (2.90)
Sparse Transformer [7]			2.80
PixelIQN [43]	5.29	49.46	
EBM [11]	6.78	38.2	
NCSNv2 [56]			31.75
NCSN [55]	8.87 ± 0.12	25.32	
SNGAN [39]	8.22 ± 0.05	21.7	
SNGAN-DDLS [4]	9.09 ± 0.10	15.42	
StyleGAN2 + ADA (v1) [29]	9.74 ± 0.05	3.26	
Ours (L , fixed isotropic Σ)	7.67 ± 0.13	13.51	≤ 3.70 (3.69)
Ours (L_{simple})	9.46 ± 0.11	3.17	≤ 3.75 (3.72)

Table 2: Unconditional CIFAR10 reverse process parameterization and training objective ablation. Blank entries were unstable to train and generated poor samples with out-of-range scores.

Objective	IS	FID
$\tilde{\mu}$ prediction (baseline)		
L , learned diagonal Σ	7.28 ± 0.10	23.69
L , fixed isotropic Σ	8.06 ± 0.09	13.22
$\ \tilde{\mu} - \tilde{\mu}_\theta\ ^2$	-	-
ϵ prediction (ours)		
L , learned diagonal Σ	-	-
L , fixed isotropic Σ	7.67 ± 0.13	13.51
$\ \tilde{\epsilon} - \epsilon_\theta\ ^2$ (L_{simple})	9.46 ± 0.11	3.17

- reverse process parameterization and training objective ablation

learned variance는 불안정하다.

이 예측하면, 그리고 비슷한 성능을 내지만, simplified objective로 훈련을 때, 성능이 좋다.

- Progressive coding

고객님 예상한 Appendix D.

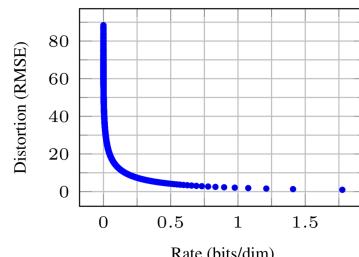
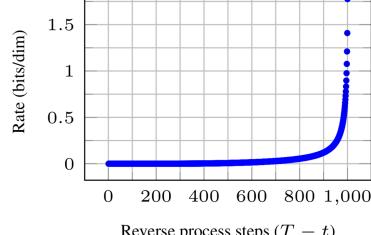
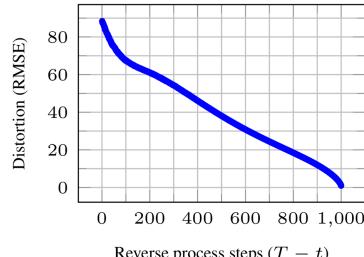
- Progressive lossy compression

Algorithm 3 Sending \mathbf{x}_0

- 1: Send $\mathbf{x}_T \sim q(\mathbf{x}_T | \mathbf{x}_0)$ using $p(\mathbf{x}_T)$
- 2: **for** $t = T - 1, \dots, 2, 1$ **do**
- 3: Send $\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_0)$ using $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$
- 4: **end for**
- 5: Send \mathbf{x}_0 using $p_\theta(\mathbf{x}_0 | \mathbf{x}_1)$

Algorithm 4 Receiving

- 1: Receive \mathbf{x}_T using $p(\mathbf{x}_T)$
- 2: **for** $t = T - 1, \dots, 1, 0$ **do**
- 3: Receive \mathbf{x}_t using $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$
- 4: **end for**
- 5: **return** \mathbf{x}_0



- Progressive generation

random bit에 progressive decomposition을 의해,

progressive unconditional generation process를 실행한다.

Fig 6, 10은 Algorithm 2를 이용한 progressive generation을 나타내고.

Fig 7은 다양한 t 에 대한 prediction $\hat{x}_0 \sim p_\theta(x_0 | x_t)$ 을 의미한다.

t 가 크면 전처리된 feature만 유지되지만, 작으면 detail한 부분도 보존된다.

- Connection to Autoregressive decoding

Variational bound eq.5는 다음과 같이 표현할 수 있다

$$L = D_{KL}(q(x_t) \| p(x_t)) + E_q \left[\sum_{t \geq 1} D_{KL}(q(x_{t-1} | x_t) \| p(x_{t-1} | x_t)) \right] + H_0(x_0) \dots \quad (16)$$

DM은 data 차트를 재생성하는 것으로는 표현할 수 없는 일반화된 bit 순서를

autoregressive 차트로 해석된다

- Interpolation



Figure 8: Interpolations of CelebA-HQ 256x256 images with 500 timesteps of diffusion.

⇒ 두 이미지의 중간으로 생기는 artifact는 x_t 의 interpolation 흡수로 해결할 수 있다

* Conclusion

학산모델은 고품질 이미지 생성을 제안했고, markov chain 훈련, denoising score matching,

annealed Langevin dynamics, autoregressive - progressive lossy compression을 제안한

VIEt DM과의 연결성을 찾았다

* Appendix A.

$$L = E_q \left[D_{KL} (q(x_t | x_0) || p(x_t)) + \sum_{t \geq 1} D_{KL} (q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) - \log p_\theta(x_0 | x_0) \right]$$

$$L = E_q \left[-\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right] = E_q \left[-\log p(x_T) - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$$

$$= E_q \left[-\log p(x_T) - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} - \log \frac{p_\theta(x_0 | x_0)}{q(x_0 | x_0)} \right]$$

$$= E_q \left[-\log p(x_T) - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_{t-1}, x_0)} \cdot \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} - \log \frac{p_\theta(x_0 | x_0)}{q(x_0 | x_0)} \right]$$

$$= E_q \left[-\log \frac{p(x_T)}{q(x_T | x_0)} - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_{t-1}, x_0)} - \log p_\theta(x_0 | x_0) \right]$$

$$= E_q \left[D_{KL} (q(x_T | x_0) || p(x_T)) + \sum_{t \geq 1} D_{KL} (q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) - \log p_\theta(x_0 | x_0) \right]$$

\Rightarrow progressive

$$L = E_q \left[-\log p(x_T) - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$$

$$= E_q \left[-\log p(x_T) - \sum_{t \geq 1} \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \cdot \frac{q(x_{t-1})}{q(x_t)} \right]$$

$$= E_q \left[-\log \frac{p(x_T)}{q(x_T)} - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t)} - \log q(x_0) \right]$$

$$= D_{KL} (q(x_T) || p(x_T)) + E_q \left[\sum_{t \geq 1} D_{KL} (q(x_{t-1} | x_t) || p_\theta(x_{t-1} | x_t)) \right] + H(x_0)$$

*Appendix B

Backbone 은 Wide resnet 2 를 기반으로 한 UNet 의 PixelCNN++ 가방법이다

Weight Normalization은 group normalization과 비슷

32x32는 32x32 ~ 4x4 resolution, 256는 64x64 resolution

resolution별 2x2 residual conv block, 16x16은 conv-block 3x3 self-attention

→ resblock의 token transformer의 sinusoidal PE 사용

EMA의 decay 0.9999로 사용

± Denoising score matching

score matching은 널널한(연이 끊긴) 널널한(연이 끊긴) 가 단위 힘은 비정규화 확률 분포 모델의
최대한 maximum log-likelihood 이 되는!

⇒ Fisher Divergence는 score matching, Kullback-Leibler Divergence와 ML을 조합함

Score matching의 주는 parameter는 훈련 데이터의 small-noise perturbation의 robustness.

Score matching은 generalized score matching과 minimum probability flow의

보다 general framework이다. 이산 널포에 대한 score matching의 generalization을 허용한다.

minimum probability flow는 다른 alternative 방식으로 재구성할 수 있다.

↳ data 널포와 예측 상태 D_{KL} 을 최소화 시킴으로써

이전 연구에서는 극소 variance의 Gaussian noise의 case는 optimal denoising의 SM을 연결하고.

SM으로 Gaussian binary RBM을 훈련하는 것이 추가 regularization을 사용하여, 일반적인 autoencoder를

훈련하는 것과 동일하고, SM은 Contrastive Divergence 간의 연결을 연구했다

논문은 regularized score matching의 형식으로 DAE의 훈련 재구성

denoising은 continuous value로는 GMM 대비 CE나 MSE로 훈련

- Explicit Score matching (ESM)

Score matching은 intractable한 partition function $Z(\theta)$ 의 probability density model $p(x; \theta)$ 의

θ 찾음을 위해 계산

→ energy function

$$p(x; \theta) = \frac{1}{Z(\theta)} \exp(-E(x; \theta))$$

$$\text{Score: } \psi(x; \theta) = \frac{\partial \log p(x; \theta)}{\partial x}$$

\Rightarrow score \in param of Gradient \in elai.

$$J_{ESMq}(\theta) = E_{q(x)} \left[\frac{1}{2} \left\| \psi(x; \theta) - \frac{\partial \log q(x)}{\partial x} \right\|^2 \right]$$

G_q

but q 를 모르면 푸는게 어렵다

- Implicit score matching

$$\underbrace{E_{q(x)} \left[\frac{1}{2} \left\| \psi(x; \theta) - \frac{\partial \log q(x)}{\partial x} \right\|^2 \right]}_{J_{ESMq}(\theta)} = \underbrace{E_{q(x)} \left[\frac{1}{2} \left\| \psi(x; \theta) \right\|^2 + \sum_{i=1}^d \frac{\partial \psi_i(x; \theta)}{\partial x_i} \right]}_{J_{ISMq}(\theta)} + C_1$$

$$\psi_i(x; \theta) = \psi(x; \theta)_i = \frac{\partial \log p(x; \theta)}{\partial x_i}$$

- Finite Sample Version of Implicit score matching

구로 넉터 D_n 대신 sample 만 있으면 (마진 데이터 data)

이제 대신 optimize 를 하면

$$J_{ISM_{q_0}}(\theta) = E_{q_0(x)} \left[\frac{1}{2} \left\| \psi(x; \theta) \right\|^2 + \sum_{i=1}^d \frac{\partial \psi_i(x; \theta)}{\partial x_i} \right]$$

$$= \frac{1}{n} \sum_{t=1}^n \left(\frac{1}{2} \left\| \psi(x_t; \theta) \right\|^2 + \sum_{i=1}^d \frac{\partial \psi_i(x_t; \theta)}{\partial x_i} \right)$$

$$\bar{J}_{ESMq} = \bar{J}_{ISMq} = \lim_{n \rightarrow \infty} J_{ISM_{q_0}}$$

Criterion of 안정성 향상을 위한 regularizer 추가.

$$J_{ISM_{reg}}(\theta) = J_{ISM_{q_0}}(\theta) + \lambda \sum_{i=1}^d \left(\frac{\partial \psi_i(x^{(i)}; \theta)}{\partial x_i} \right)^2$$

- Linking Score Matching to the Denoising Autoencoder Objective

- Matching the score of a Non-Parametric Estimator

ESM with Parzen window density estimator $q_\sigma(\tilde{x})$

$$J_{ESMq_\sigma}(\theta) = \mathbb{E}_{q_\sigma(\tilde{x})} \left[\frac{1}{2} \left\| \psi(\tilde{x}; \theta) - \frac{\partial \log q_\sigma(\tilde{x})}{\partial \tilde{x}} \right\|^2 \right]$$

6>0 이고, 미분 가능성이 있다. $\mathbb{E}_{q_\sigma} \left[\left\| \frac{\partial \log q_\sigma(\tilde{x})}{\partial \tilde{x}} \right\|^2 \right]$ 는 유한이다.

$$\Rightarrow J_{ESMq_\sigma} \sim J_{ISMq_\sigma}$$

$\lim_{\delta \rightarrow 0}$ 일정한 확률로, 미분 가능한 경우에 대해서만 J_{ISM} 과 J_{ESM} 가 일치한다.

- Denoising score matching

$$(clean, corrupted) = (x, \tilde{x}) \quad , \quad q_\sigma(\tilde{x}|x) = q_\sigma(\tilde{x}|x)q_0(x)$$

$$J_{DSMq_\sigma}(\theta) = \mathbb{E}_{q_\sigma(x, \tilde{x})} \left[\frac{1}{2} \left\| \psi(\tilde{x}; \theta) - \frac{\partial \log q_\sigma(\tilde{x}|x)}{\partial \tilde{x}} \right\|^2 \right]$$

$$\tilde{x} \rightarrow x \text{ 일 때 } \frac{\partial \log q_\sigma(\tilde{x}|x)}{\partial \tilde{x}} = \frac{1}{\sigma^2}(x - \tilde{x}).$$

$$J_{DSMq_\sigma} \sim J_{DAE\sigma}$$

$$\begin{aligned} J_{DSMq_\sigma}(\theta) &= \mathbb{E}_{q_\sigma(x, \tilde{x})} \left[\frac{1}{2} \left\| \psi(\tilde{x}; \theta) - \frac{\partial \log q_\sigma(\tilde{x}|x)}{\partial \tilde{x}} \right\|^2 \right] \\ &= \mathbb{E}_{q_\sigma(x, \tilde{x})} \left[\frac{1}{2} \left\| \frac{1}{\sigma^2} (\mathbf{W}^T \text{sigmoid}(\mathbf{W}\tilde{x} + \mathbf{b}) + \mathbf{c} - \tilde{x}) - \frac{1}{\sigma^2}(x - \tilde{x}) \right\|^2 \right] \\ &= \frac{1}{2\sigma^4} \mathbb{E}_{q_\sigma(x, \tilde{x})} \left[\left\| \mathbf{W}^T \text{sigmoid}(\mathbf{W}\tilde{x} + \mathbf{b}) + \mathbf{c} - x \right\|^2 \right] \\ &= \frac{1}{2\sigma^4} J_{DAE\sigma}(\theta). \end{aligned}$$

We have thus shown that

$$J_{DSMq_\sigma} \sim J_{DAE\sigma} \tag{14}$$

Proof that $J_{ESMq_\sigma} \succsim J_{DSMq_\sigma}$ (11)

The explicit score matching criterion using the Parzen density estimator is defined in Eq. 7 as

$$J_{ESMq_\sigma}(\theta) = \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \psi(\tilde{\mathbf{x}}; \theta) - \frac{\partial \log q_\sigma(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \right\|^2 \right]$$

which we can develop as

$$J_{ESMq_\sigma}(\theta) = \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\psi(\tilde{\mathbf{x}}; \theta)\|^2 \right] - S(\theta) + C_2 \quad (16)$$

where $C_2 = \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \frac{\partial \log q_\sigma(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \right\|^2 \right]$ is a constant that does not depend on θ , and

$$\begin{aligned} S(\theta) &= \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \right\rangle \right] \\ &= \int_{\tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}}) \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}}) \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\frac{\partial}{\partial \tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}})}{q_\sigma(\tilde{\mathbf{x}})} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial}{\partial \tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}}) \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial}{\partial \tilde{\mathbf{x}}} \int_{\mathbf{x}} q_0(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \left\langle \psi(\tilde{\mathbf{x}}; \theta), \int_{\mathbf{x}} q_0(\mathbf{x}) \frac{\partial q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} d\mathbf{x} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \left\langle \psi(\tilde{\mathbf{x}}; \theta), \int_{\mathbf{x}} q_0(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \frac{\partial \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} d\mathbf{x} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \int_{\mathbf{x}} q_0(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\rangle d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \int_{\mathbf{x}} q_\sigma(\tilde{\mathbf{x}}, \mathbf{x}) \left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\rangle d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}}, \mathbf{x})} \left[\left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\rangle \right]. \end{aligned}$$

Substituting this expression for $S(\theta)$ in Eq. 16 yields

$$\begin{aligned} J_{ESMq_\sigma}(\theta) &= \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\psi(\tilde{\mathbf{x}}; \theta)\|^2 \right] \\ &\quad - \mathbb{E}_{q_\sigma(\mathbf{x}, \tilde{\mathbf{x}})} \left[\left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\rangle \right] + C_2. \end{aligned} \quad (17)$$

We also have defined in Eq. 9,

$$J_{DSMq_\sigma}(\theta) = \mathbb{E}_{q_\sigma(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \psi(\tilde{\mathbf{x}}; \theta) - \frac{\partial \log q_\sigma(\tilde{\mathbf{x}} | \mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\|^2 \right],$$

which we can develop as

$$\begin{aligned} J_{DSMq_\sigma}(\theta) &= \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\psi(\tilde{\mathbf{x}}; \theta)\|^2 \right] \\ &\quad - \mathbb{E}_{q_\sigma(\mathbf{x}, \tilde{\mathbf{x}})} \left[\left\langle \psi(\tilde{\mathbf{x}}; \theta), \frac{\partial \log q_\sigma(\tilde{\mathbf{x}} | \mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\rangle \right] + C_3 \end{aligned} \quad (18)$$

where $C_3 = \mathbb{E}_{q_\sigma(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \left\| \frac{\partial \log q_\sigma(\tilde{\mathbf{x}} | \mathbf{x})}{\partial \tilde{\mathbf{x}}} \right\|^2 \right]$ is a constant that does not depend on θ .

Looking at equations 17 and 18 we see that $J_{ESMq_\sigma}(\theta) = J_{DSMq_\sigma}(\theta) + C_2 - C_3$. We have thus shown that the two optimization objectives are equivalent.

Reference paper: Nonequilibrium

2. Algorithm

목적은 forward process의 복잡한 distribution을 간단하고 tractable한 분포로 바꾸기

Generative model로 reverse를 계산하는 것

2.1 forward trajectory

$\hat{q}(x_t | x_{t-1}) \sim q(x_t | x_{t-1})$

후련도(연산) $\sim \pi(y) \rightarrow \sum_t$ 까지 목표

↳ Markov chain kernel $T_\pi(y|y';\beta)$ 을 빌려서 계산

$$\pi(y) = \int dy' T_\pi(y|y';\beta) \pi(y')$$

$$q(x_t | x_{t-1}) = T_\pi(x_t | x_{t-1}; \beta_t)$$

$\downarrow t=0 \sim T$

$$q(x_{t=0:T}) = q(x_0) \prod_{t=1}^T q(x_t | x_{t-1})$$

2.2 reverse trajectory

$$p(x_r) = \pi(x_r)$$

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

$$\beta_1, \beta_2, \dots, \beta_l \sim \dots \beta_l$$

$$\sqrt{\alpha_t} (\sqrt{\alpha_3} (\sqrt{\alpha_1} x_0 + \beta_1 z) + \beta_2 z) + \beta_3 z$$

$$\sqrt{\alpha_1} x_0 + \sqrt{\alpha_2} \beta_1 z + \beta_2 z$$

$$\sqrt{\alpha_3} x_0 + \sqrt{\alpha_3} (\sqrt{\alpha_2} \beta_1 z + \sqrt{\alpha_3} \beta_2 z + \beta_3 z)$$

$$\sqrt{\alpha_t} x_0 + \sqrt{\alpha_t} \sqrt{\alpha_3} \sqrt{\alpha_2} \beta_1 z + \sqrt{\alpha_t} \sqrt{\alpha_3} \beta_2 z + \sqrt{\alpha_t} \beta_3 z + \beta_4 z$$

$$z (\sqrt{\alpha_t} \sqrt{\alpha_3} \sqrt{\alpha_2} \beta_1 + \sqrt{\alpha_t} \sqrt{\alpha_3} \beta_2 + \sqrt{\alpha_t} \beta_3 + \beta_4)$$

$$\alpha_4 \alpha_3 \alpha_2 (1 - \alpha_1) + \alpha_4 \alpha_3 (1 - \alpha_2) + \alpha_4 (1 - \alpha_3) + (1 - \alpha_4)$$

$$1 - \alpha_1 \alpha_2 \alpha_3 \dots \alpha_r$$

$$\alpha_4 \alpha_3 \alpha_2 = \alpha_1 \alpha_2 \alpha_3 \alpha_4 - \alpha_1 \alpha_2 \alpha_3 \alpha_4 - \alpha_1 \alpha_2 \alpha_3 + \alpha_1 \alpha_2 + 1 - \alpha_1$$

$$t-1 \rightarrow t : N(\sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

$$0 \rightarrow t : N(\sqrt{\alpha_t} x_t, (1-\alpha_t) I) \rightarrow \text{one row at a time}$$

$$t-1 \rightarrow t : N(M_t(x_{t-1}), \Sigma_t(x_{t-1}))$$

$$L = f_t = \frac{1 - \alpha_t}{1 - \alpha_t} \beta_t$$

$$Q(x_{t-1} | x_t, \Sigma_t) = \frac{\sqrt{\alpha_t} \beta_t x_0 + \frac{\sqrt{\alpha_t} (1 - \alpha_{t-1})}{1 - \alpha_t} x_t}{1 - \alpha_t}$$

$$Q(x_t | x_0) = \sqrt{1 - \beta_t} x_0 + \beta_t I$$

$$Q(x_t | x_0) = \sqrt{\alpha_t} x_0 + (1 - \alpha_t) I$$

$$\sqrt{1 - \beta_t} (\sqrt{\alpha_{t-1}} x_0 + (1 - \alpha_{t-1})) + \beta_t = x_t$$

$$x_t = \sqrt{d_t} x_0 + (1-d_t) = \sqrt{1-d_t} x_{t-1} + \hat{x}_t$$

$$x_{t-1} = \frac{\sqrt{d_{t-1}} \beta_t}{1-d_t} x_0 + \frac{\sqrt{d_t}(1-d_{t-1})}{1-d_t} x_t + \frac{1-d_{t-1}}{1-d_t} \beta_t$$

$$\begin{aligned} & \underbrace{\downarrow}_{\frac{\sqrt{d_t}(1-d_{t-1})}{1-d_t} (\sqrt{d_{t-1}} x_0 + (1-d_t))} \\ &= \frac{d_t(1-d_{t-1})\sqrt{d_{t-1}}}{1-d_t} x_0 + \sqrt{d_t}(1-d_{t-1}) \\ &\Rightarrow \cancel{\frac{(d_t - d_{t-1} + \beta_t)\sqrt{d_{t-1}}}{1-d_t} x_0 + \sqrt{d_t}(1-d_{t-1}) + \frac{1-d_{t-1}}{1-d_t} \beta_t}^{\cancel{(1-d_t)}} \\ &= \sqrt{d_{t-1}} x_0 + \sqrt{d_t}(1-d_{t-1}) + \frac{1-d_{t-1}-d_t+\beta_t}{1-d_t} \end{aligned}$$

$$= \sqrt{d_{t-1}} x_0 + \frac{\sqrt{d_t}(1-d_{t-1})(1-d_t) + 1-d_{t-1}-d_t+\beta_t}{1-d_t}$$

$$\frac{\sqrt{d_t}(1-d_t-d_{t-1}-\beta_t)}{1-d_t}$$

$$\frac{\sqrt{d_t}(1-d_{t-1})(1-d_t) + (1+d_t-d_{t-1}-d_t)}{1-d_t} \Rightarrow 1-d_t$$

$$\sqrt{d_t}(1-d_{t-1}) + \frac{-1+d_t+1-d_{t-1}-d_t}{1-d_t}$$

$$= \sqrt{d_t}(1-d_{t-1}) - 1 + \frac{1-d_t+1-d_{t-1}}{1-d_t}$$

220V

30A

