

Understanding Diffusion Models: A Unified Perspective



NTTコミュニケーションズ 藤原大悟

本日のモチベーション

- 個人的に後回しにしていた **Diffusion Model** について **真面目に**理解する
 - ざっくりノイズを足すとか拡散過程を使ってることは知っているけど正確な定式化とか知らない
 - 数式が多いですがお気持ちを説明するようにします
 - なんで上手くいくのか直感的に理解できてない
 - 色んな生成モデルを復習して関係性を理解したい
- ちょっと業務逼迫で急ピッチで読んだ&作ったので粗いところがあったらごめんなさい

生成モデルとは

- 生成モデルの大まかな目的
 - データ $D = (x_1, x_2, \dots, x_N)^T$ が与えられた時に、データ x_i を生み出す分布関数 $p(x)$ をモデリングし、そこからのサンプリングを行うこと
- 生成モデルの系譜 (知ってる人向けの説明)
 - Generative Adversarial Networks (GANs)
 - みなさんご存知Diffusion登場以前の最強モデル。Discriminator(識別器)とGenerator(生成器)を戦わせて精度アップ。
 - 背後には密度比推定とかの理論が
 - "Likelihood-Based" (原文ママ)
 - パラメトリックな確率モデル p_θ を考えて、データに対する尤度を最大化してフィッティング
 - Autoregressive Models
 - (入出力を繰り返し適用すればデータ生成できるからそういう意味で生成モデルとして挙げられている?)
 - Normalizing Flows
 - 正規ノイズにNNなどによる非線形変換を何回もかましてデータ分布に近づけようみたいなやつ
 - カップリングとかで逆変換可能な非線形変換を作って、尤度Lossを計算可能にする工夫。変分推論を利用。
 - Variational Auto Encoders (VAEs)
 - Auto Encoder構造を基礎に、中間層(低次元)は正規ノイズから生成されるものとして尤度ベースのLossで学習。
 - Encoder $p_\theta(z|x)$ とDecoder $p_\phi(x|z)$ を同時に学習。変分推論の利用。

生成モデル(エネルギーモデル,スコアベースモデル)

- 生成モデルの系譜 (知ってる人向けの説明)
 - "Likelihood-Based" (原文ママ)
 - パラメトリックな確率モデルを考えて、データに対する尤度を最大化してフィッティング
 - Energy Based Models
 - よく統計力学とかでありそうな分布形、 $p(x) = \frac{1}{Z_\theta} \exp(-E_\theta(x))$ を仮定して尤度Lossで $E_\theta(x)$ を学習
 - 統計力学の知見を流用しつつ、尤度に登場する期待値(規格化定数 Z_θ)計算を近似とサンプリングで上手く回避
 - Z_θ を求めず分布 $p(x)$ からサンプリングする方法 (MCMC+ランジュバンダイナミクス)
 - $x^{t+1} = x^t + \eta \nabla_x \log p_\theta(x) + \omega$ (ω : ガウスノイズ, η : ハイパーパラメタ)
 - ※スコア $s_\theta(x) := \nabla_x \log p_\theta(x) = \nabla_x \log(\frac{1}{Z_\theta} \exp(-E_\theta(x))) = \nabla_x \{\log(\exp(-E_\theta(x))) - \log(Z_\theta)\} = -\nabla_x E_\theta(x)$

Score Based Models

- 発想: $E_\theta(x)$ を介さずに直接 $s_\theta(x)$ だけ学習したらいいのでは??
 - スコアマッチング:

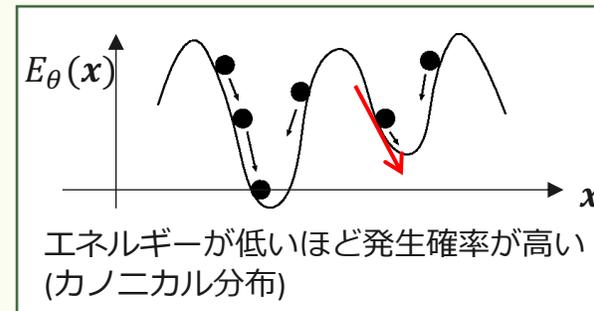
Fisher Divergence
の最小化問題

- データからスコアをフィッティングするLossを計算するための工夫(変形)

$$E_{x \sim p(x)} \left[\frac{1}{2} \|\nabla_x \log p(x) - s_\theta(x)\|_2^2 \right] = E_{x \sim p(x)} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right]$$

ここがわからなくてもOK

筆者曰く、Diffusion Modelは尤度ベース、スコアベースどちらとしても解釈可能とのこと (楽しみ)



変分推論ってなんだっけ？

前提

データ x に対して、その根本構造を有する(低次元)潜在変数 z の存在を考えるとうまく行くことが多い。
(MNIST画像 x に対して数字や手癖 z の存在を考えるのは自然)

ということでその関係性 $p(x, z)$ に興味があり、モデル化したい。
(※実際は多くの場合 $p(z|x)$ をモデル化する)

そこで、我々は真の分布 $p(z|x)$ に対して 適当な提案分布のクラス $q(z|x)(= q_\phi(z|x))$ で近似を図る

で、天下りのだが右が
成り立つ

KL divergence 最小化を
ELBOを介して行うこと
ができる

ここをでかくすればOK

=ELBO

ゼロにしたら近似成功

$$\log p(\mathbf{x}) = \log p(\mathbf{x}) \int q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \quad (\text{Multiply by } 1 = \int q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z}) \quad (9)$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) (\log p(\mathbf{x})) d\mathbf{z} \quad (\text{Bring evidence into integral}) \quad (10)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})] \quad (\text{Definition of Expectation}) \quad (11)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Apply Equation 2}) \quad (12)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}) q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Multiply by } 1 = \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})}) \quad (13)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Split the Expectation}) \quad (14)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \quad (\text{Definition of KL Divergence}) \quad (15)$$

$$\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \quad (\text{KL Divergence always } \geq 0) \quad (16)$$

変分推論ってなんだっけ？

ちなみに

VAEなんかでよくやる手は、 $p(\mathbf{x}) = p_{\theta}(\mathbf{x})$ (や、 $p(\mathbf{x}|\mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})$) とおいて、 θ に対する尤度の最大化もいっぺんにやってしまう (※EMアルゴリズムっぽい)

とにかく ELBO をでかくすれば KL も小さくなって尤度 ($\log p_{\theta}(\mathbf{x})$) もでかくなるだろうと

※EMアルゴリズムはKL divergence 最小化と尤度最大化を交互に行う

$\log p_{\theta}(\mathbf{x})$

$$\log p(\mathbf{x}) = \log p(\mathbf{x}) \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \quad (\text{Multiply by } 1 = \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}) \quad (9)$$

定数
(for ϕ)

$$= \int q_{\phi}(\mathbf{z}|\mathbf{x}) (\log p(\mathbf{x})) d\mathbf{z} \quad (\text{Bring evidence into integral}) \quad (10)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})] \quad (\text{Definition of Expectation}) \quad (11)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Apply Equation 2}) \quad (12)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Multiply by } 1 = \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})}) \quad (13)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right] \quad (\text{Split the Expectation}) \quad (14)$$

ここをでかくすれば OK
= ELBO

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \quad (\text{Definition of KL Divergence}) \quad (15)$$

ゼロにしたなら近似成功

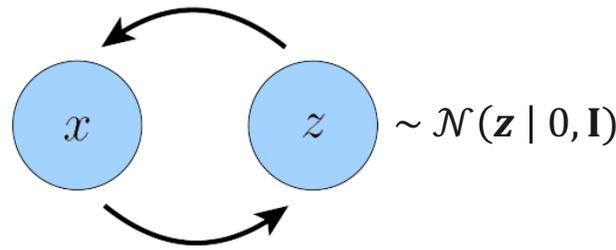
$$\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (\text{KL Divergence always } \geq 0) \quad (16)$$

VAE → Hierarchical VAE

例えば白黒画像ならベルヌーイ分布

$$p(x|z) = \text{Bern}(x | \mu_\theta(z)) \text{とか}$$

VAE



$$q(z|x) = \mathcal{N}(z | \mu_\phi(x), \sigma_\phi^2(x)\mathbf{I})$$

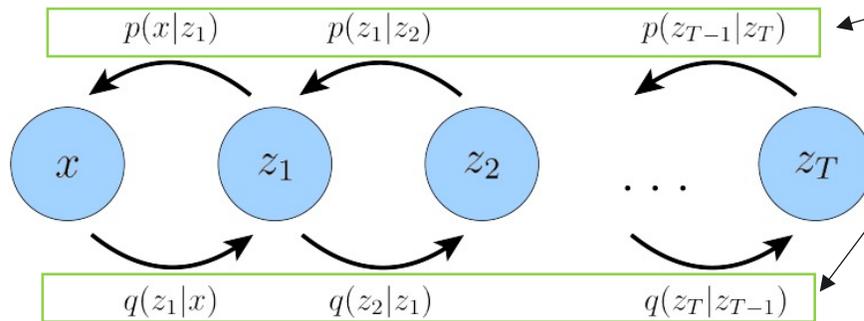
Loss (ELBO)

$$\begin{aligned} \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(\mathbf{x}, z)}{q_\phi(z|x)} \right] &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(\mathbf{x}|z)p(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(\mathbf{x}|z)] + \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(z)}{q_\phi(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(\mathbf{x}|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \| p(z))}_{\text{prior matching term}} \end{aligned}$$

Figure 1: A Variational Autoencoder graphically represented. Here, encoder $q(z|x)$ defines a distribution over latent variables z for observations x , and $p(x|z)$ decodes latent variables into observations.

階層化しただけ

HVAE



$$\begin{aligned} p(\mathbf{x}, z_{1:T}) &= p(z_T)p_\theta(\mathbf{x}|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t) \\ q_\phi(z_{1:T}|\mathbf{x}) &= q_\phi(z_1|\mathbf{x}) \prod_{t=2}^T q_\phi(z_t|z_{t-1}) \end{aligned}$$

Loss (ELBO)

$$\mathbb{E}_{q_\phi(z_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z_{1:T})}{q_\phi(z_{1:T}|\mathbf{x})} \right] = \mathbb{E}_{q_\phi(z_{1:T}|\mathbf{x})} \left[\log \frac{p(z_T)p_\theta(\mathbf{x}|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)}{q_\phi(z_1|\mathbf{x}) \prod_{t=2}^T q_\phi(z_t|z_{t-1})} \right]$$

Figure 2: A Markovian Hierarchical Variational Autoencoder with T hierarchical latents. The generative process is modeled as a Markov chain, where each latent z_t is generated only from the previous latent z_{t+1} .

Hierarchical VAE → Diffusion Models

Variational Diffusion Models

この文脈でのDiffusion Modelsは以下の制約を持つHVAEと解釈できる

- 潜在次元はデータ次元と等しい。
- 各タイムステップにおけるエンコーダー $q(\mathbf{z}|\mathbf{x})$ は学習されない。
前のタイムステップの出力を中心とする線形ガウス分布モデルとして(ハイパラが与えられ)定義されている。
- 潜在変数は、最終タイムステップ T における分布が標準ガウス分布となるように、各ステップで時間発展する
 - (T が十分大きければ最終ステップの分布は勝手に標準ガウス分布になる)

- 前後で分散が同程度になるよう調整されている
- ハイパラ α_t は t によって異なる(学習してもいい)

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

データ空間サイド

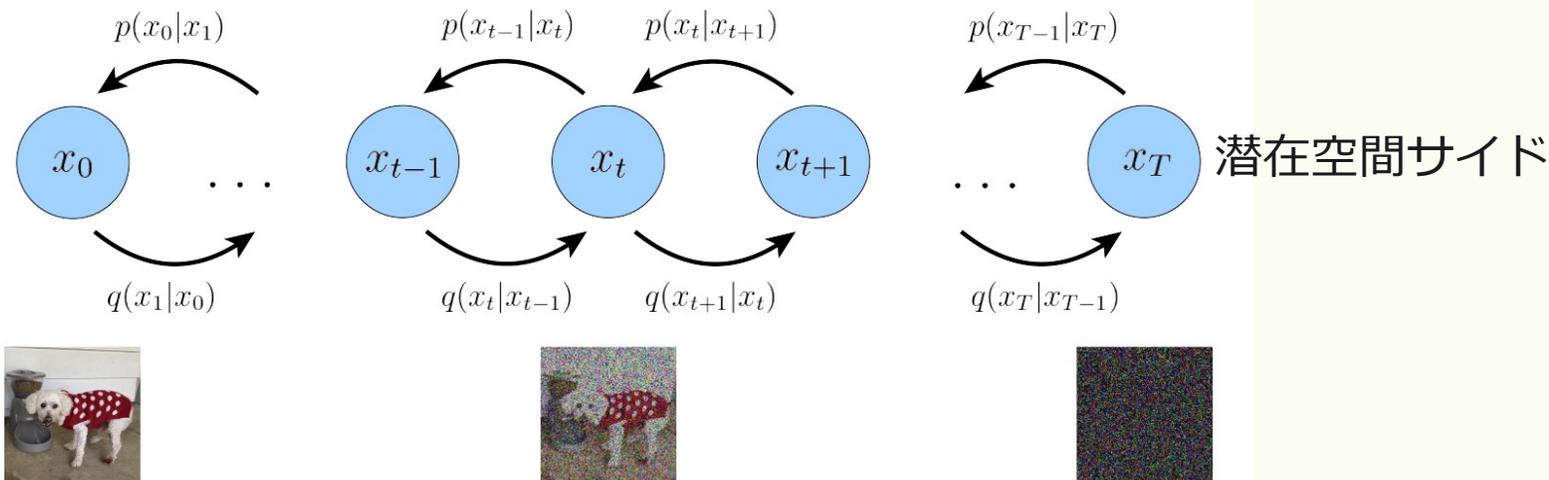


Figure 3: A visual representation of a Variational Diffusion Model; \mathbf{x}_0 represents true data observations such as natural images, \mathbf{x}_T represents pure Gaussian noise, and \mathbf{x}_t is an intermediate noisy version of \mathbf{x}_0 . Each $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ is modeled as a Gaussian distribution that uses the output of the previous state as its mean.

Diffusion Models の新しいLoss

Variational Diffusion Models

(H)VAEと同様にELBOの最大化に落とせる

まあ変形は置いといて最後の方を見ると以下の3つに分解できる。

- 再構成誤差
 - データとの直接のDecoder Loss
- 事前分布のKL divergence
 - 学習可能なパラメータがないため実質定数
 - Tが十分大きければ0に
- EncoderとDecoder過程の一致を要求する項
 - Decoderは θ が学習可能なので、Encoder (=与えられたガウス分布拡散)に合わせなければいけない

$$\begin{aligned}
 \log p(\mathbf{x}) &= \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\
 &= \log \int \frac{p(\mathbf{x}_{0:T})q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
 &= \log \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad \text{ここをでかくすればOK} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \text{=ELBO} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=1}^{T-1} p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))]}_{\text{prior matching term}} \\
 &\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}
 \end{aligned}$$

参考：VAEのELBO

$$\begin{aligned}
 \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right] &= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}|z)p(z)}{q_{\phi}(z|\mathbf{x})} \right] \\
 &= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|z)] + \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(z)}{q_{\phi}(z|\mathbf{x})} \right] \\
 &= \underbrace{\mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \parallel p(z))}_{\text{prior matching term}}
 \end{aligned}$$

Diffusion Models の新しいLoss

Variational Diffusion Models

補足

- EncoderとDecoder過程の一致を要求する項

- Decoderは θ が学習可能なので、Encoder (=与えられたガウス分布拡散)に合わせなければいけない
- x_t を中心に両側から挟む感じ (ピンクと緑の矢印)

- 悪い点

- 実際のロスの計算は期待値のモンテカルロ推定が必要だが、第3項については2変数(x_{t-1}, x_{t+1})で期待値を取るなので、**推定値の誤差(分散)が大きくなってしま**

→修正

$$-\sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(x_{t-1}, x_{t+1}|x_0)} [D_{\text{KL}}(q(x_t|x_{t-1}) \parallel p_{\theta}(x_t|x_{t+1}))]}_{\text{consistency term}}$$

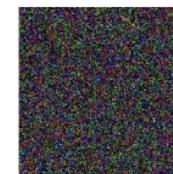
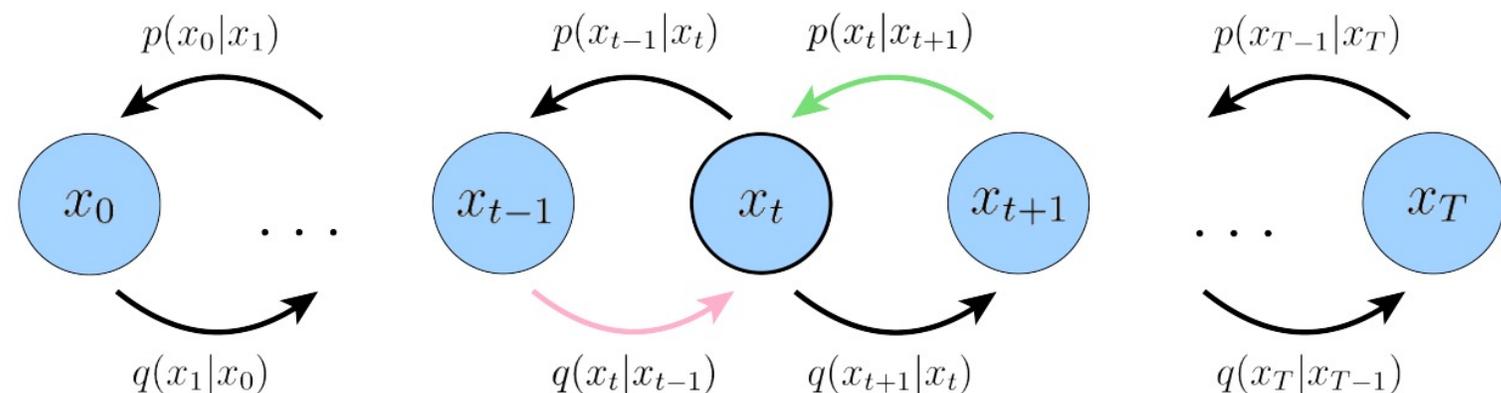


Figure 4: Under our first derivation, a VDM can be optimized by ensuring that for every intermediate x_t , the posterior from the latent above it $p_{\theta}(x_t|x_{t+1})$ matches the Gaussian corruption of the latent before it $q(x_t|x_{t-1})$. In this figure, for each intermediate x_t , we minimize the difference between the distributions represented by the pink and green arrows.

Diffusion Models の修正したLoss

計算が
めっちゃめっちゃ長いのでざっくり説明 (次ページ)

$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] && \text{同じELBOからスタート} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)} \right]
 \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} + \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}
 \end{aligned}$$

同じ

ほぼ同じ

新しいLoss

Diffusion Models の修正したLoss

Variational Diffusion Models

$$\underbrace{- \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \leftarrow \underbrace{- \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}$$

新しいLoss (時間が揃って嬉しい)

元のLoss (時間がずれててうざい)

- EncoderとDecoder過程の一致を要求する項を修正

基本発想:

- ベイズの定理を使えば $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ を逆転して $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ みたいな形に持っていけるんじゃないか?



- ピンクと緑の矢印の場所が揃った!!

良い点

- 期待値推定が \mathbf{x}_t についてのみに なり精度向上

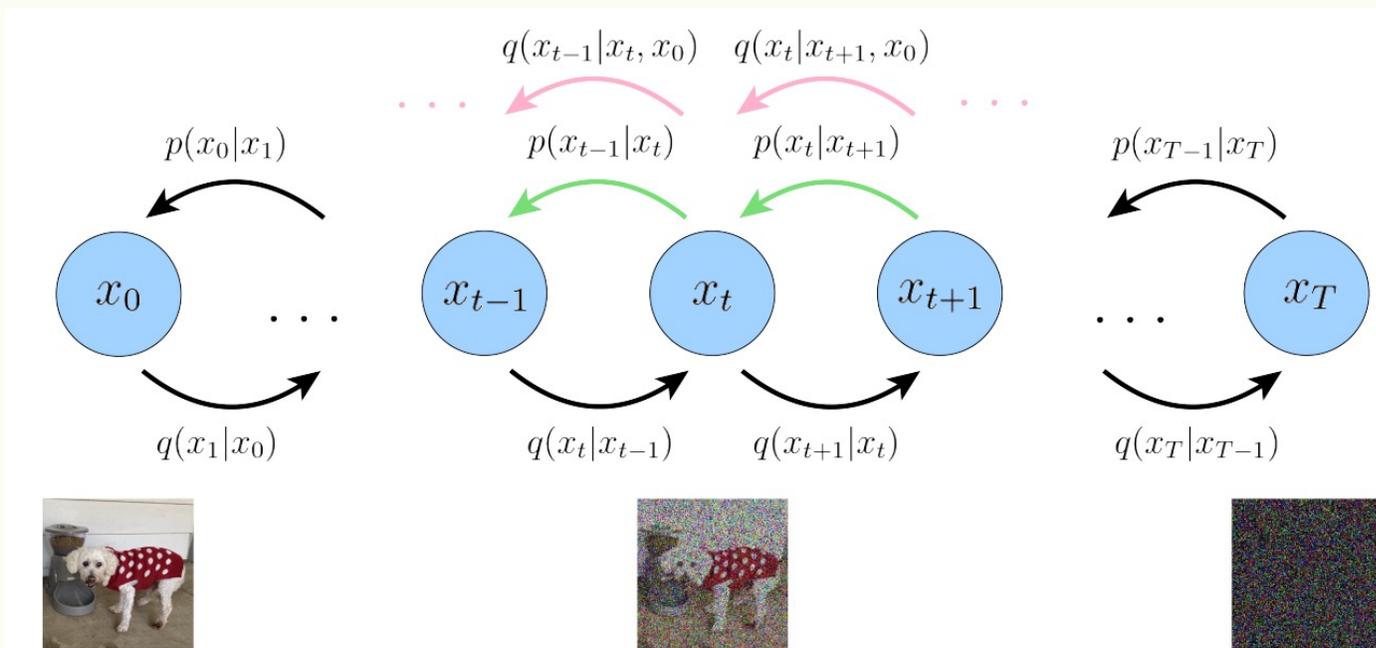


Figure 5: Depicted is an alternate, lower-variance method to optimize a VDM; we compute the form of ground-truth denoising step $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ using Bayes rule, and minimize its KL Divergence with our approximate denoising step $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$. This is once again denoted visually by matching the distributions represented by the green arrows with those of the pink arrows. Artistic liberty is at play here; in the full picture, each pink arrow must also stem from \mathbf{x}_0 , as it is also a conditioning term.

Diffusion Models の修正したLoss

Variational Diffusion Models

$$\underbrace{- \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \leftarrow \underbrace{- \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}$$

新しいLoss (時間が揃ってて嬉しい)

元のLoss (時間がずれててうざい)

- EncoderとDecoder過程の一致を要求する項を修正

基本発想:

- ベイズの定理を使えば $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ を逆転して $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ みたいな形に持っていけるんじゃないか?



- ピンクと緑の矢印の場所が揃った!!

良い点

- 期待値推定が \mathbf{x}_t についてのみに なり精度向上

ポイント!!!

- $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ が解析的に計算できるのは今回だけ!
 - なぜなら単純な線形ガウス変換(拡散過程)を考えていたから
 - 一般のHVAEでは計算困難

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$$

$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ 拡散過程(既知)

拡散過程

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

↓等式で書き下すところ

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon \quad \text{with } \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

- $q(\mathbf{x}_t|\mathbf{x}_0)$ について、
 - \mathbf{x}_t に再起的に繰り返し代入すると、

$$\mathbf{x}_t = \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0 = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$$

$$\sim q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

ということで計算できる(代入するだけ)

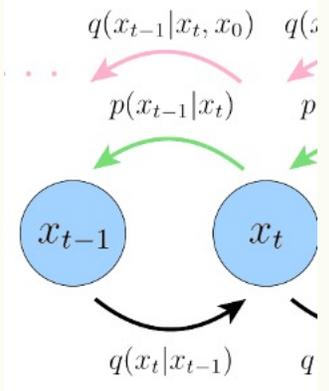
Diffusion Models による画像生成(元画像推定)

Variational Diffusion Models

先の結果を $q(x_{t-1}|x_t, x_0)$ に代入して具体的に計算すると、

$$-\sum_{t=2}^T \underbrace{\mathbb{E}_{q(x_t|x_0)} [D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]}_{\text{denoising matching term}}$$

新しいLoss (時間が揃って嬉しい)



$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})$$

のようにガウス分布になる。(ガウス分布×ガウス分布÷ガウス分布なので自明)

さて我々は $p_{\theta}(x_{t-1}|x_t)$ にどんな分布を持ってきてても良いわけだが、形が分かっているので、**ガウス分布を仮定してみよう。** しかも分散は既知(α_t から計算可)なので、**正解の分散を教えてあげることにしよう。**

真の平均がどんなだったかという、**こんな**感じ。
 では、ここも形が分かっているので真似して、**こういう形**を仮定しよう。
 Loss, KL最小化は二乗誤差 $\|\hat{x}_{\theta}(x_t, t) - x_0\|^2$ の最小化に帰着する。

$\hat{x}_{\theta}(x_t, t)$ は元画像 x_0 の推定器になる!!!

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}$$

$$\mu_{\theta}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_{\theta}(x_t, t)}{1 - \bar{\alpha}_t}$$

NNで与えてあげる

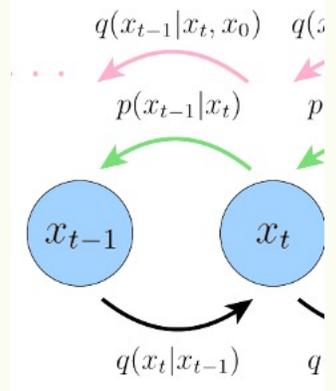
Diffusion Models による画像生成(元画像推定)

Variational Diffusion Models

先の結果を $q(x_{t-1}|x_t, x_0)$ に代入して具体的に計算すると、

$$-\sum_{t=2}^T \underbrace{\mathbb{E}_{q(x_t|x_0)} [D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]}_{\text{denoising matching term}}$$

新しいLoss (時間が揃って嬉しい)



$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})$$

のようにガウス分布になる。(ガウス分布×ガウス分布÷ガウス分布なので自明)

さて我々は $p_{\theta}(x_{t-1}|x_t)$ にどんな分布を持ってきてても良いわけだが、形が分かっているので、**ガウス分布を仮定してみよう**。
しかも分散は既知(α_t から計算可)なので、**正解の分散を教えてあげることにしよう**。

真の平均がどんなだったかという、**こんな**感じ。
では、ここも形が分かっているので真似して、**こういう形**を仮定しよう。
Loss, KL最小化は二乗誤差 $\|\hat{x}_{\theta}(x_t, t) - x_0\|^2$ の最小化に帰着する。

$\hat{x}_{\theta}(x_t, t)$ は元画像 x_0 の推定器になる!!!

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}$$

$$\mu_{\theta}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_{\theta}(x_t, t)}{1 - \bar{\alpha}_t}$$

NNで与えてあげる

素朴な疑問：この結果がどこまで見えていて階層的な拡散過程考えたんだろう、、、

- 解析的に逆拡散過程 $q(x_{t-1}|x_t, x_0)$ が計算できるため、好きなパラメトライズを選択できる $p_\theta(x_{t-1}|x_t)$ に分かっている情報を全て組み込むと、必然的に二つの分布の間で足りない情報(= x_0)を推定する形になるのかもしれないと思った
- (ノイズ込みで)分布をモデリングして、そこからサンプリングするといった生成方法ではなく、あくまで「元画像の推定」のテイなので、めちゃくちゃ解像度の高い画像を出力できるのかなと思った
- ただ原理的に「元画像を推定する」点で、Training Dataを再現する傾向が強うそうなので汎化性能に疑問が残るので、AI絵師、著作権関連で言われる「既存画像のパクリを出力してるだけ」説は結構正しそうな気がしてしまった

ハイパラ(α_t)の学習

ここでは割愛

- α_t ではなく $\bar{\alpha}_t$ をNNでモデリングすること
 - 結局計算には $\bar{\alpha}_t$ だけ使えば十分だし、全ての時間 t で掛け合わせないといけないので非効率
- 時間と共に単調減少する関数としてモデリングする必要があること
 - 最終的にノイズに終着してもらう必要があるため。

がポイント

付加ノイズ推定としての解釈

ここでは割愛。ざっくり言うと、

- 元画像そのものを推定する問題として定式化したが、等価な方法として、 t ステップ時点で元画像に加わった合計ノイズ、を推定する方法もある。(簡単な式変形で求まる)
- なぜかノイズを学習した方が経験的に精度が良いという話もあるらしい
 - 等価なのになぜ。。。？

スコア推定としての解釈

- Tweedieの公式というものを使う

指数関数型分布の真の平均は、その分布から得られた標本が与えられたとき、標本による最尤推定値（別名：経験平均）に推定値のスコアを含む何らかの補正項を加えたもので推定できることを述べている。

ガウス分布の標本
が一つある場合：

Mathematically, for a Gaussian variable $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$, Tweedie's Formula states that:

$$\mathbb{E}[\boldsymbol{\mu}_z | \mathbf{z}] = \mathbf{z} + \boldsymbol{\Sigma}_z \nabla_{\mathbf{z}} \log p(\mathbf{z})$$

さっきまでの結果に流用

真の分布

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Tweedieの公式を適用

(\mathbf{x}_t から $\sqrt{\bar{\alpha}_t} \mathbf{x}_0$ を推定するテイ)

$$\mathbb{E}[\boldsymbol{\mu}_{\mathbf{x}_t} | \mathbf{x}_t] = \mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

$$\sqrt{\bar{\alpha}_t} \mathbf{x}_0 = \mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)$$

$$\therefore \mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}$$

スコア推定としての解釈

前ページ
Tweedieの公式より
以下に代入

$$\therefore \mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}$$

$$-\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

Diffusion ModelのLoss (時間が揃ってて嬉しい)

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ の平均

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ の平均

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$



良い感じに式変形



$$= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t)$$

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \mathbf{s}_{\theta}(\mathbf{x}_t, t)$$

$\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)$ と定数などを
適当に置き換え

先述と同様にLoss, KL最小化は二乗誤差 $\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|^2$ の最小化に帰着する。

よって、
Diffusion Modelの学習は**スコア推定の問題**としても解釈できる (※スコア関数も、NNでモデリングする)

(再掲)生成モデル(エネルギーモデル,スコアベースモデル)

- 生成モデルの系譜 (知ってる人向けの説明)
 - "Likelihood-Based" (原文ママ)
 - パラメトリックな確率モデルを考えて、データに対する尤度を最大化してフィッティング
 - Energy Based Models
 - よく統計力学とかでありそうな分布形、 $p(x) = \frac{1}{Z_\theta} \exp(-E_\theta(x))$ を仮定して尤度Lossで $E_\theta(x)$ を学習
 - 統計力学の知見を流用しつつ、尤度に登場する期待値(規格化定数 Z_θ)計算を近似とサンプリングで上手く回避
 - Z_θ を求めず分布 $p(x)$ からサンプリングする方法 (MCMC+ランジュバンダイナミクス)
 - $x^{t+1} = x^t + \eta \nabla_x \log p_\theta(x) + \omega$ (ω : ガウスノイズ, η : ハイパーパラメタ)
 - ※スコア $s_\theta(x) := \nabla_x \log p_\theta(x) = \nabla_x \log(\frac{1}{Z_\theta} \exp(-E_\theta(x))) = \nabla_x \{\log(\exp(-E_\theta(x))) - \log(Z_\theta)\} = -\nabla_x E_\theta(x)$

Score Based Models

- 発想: $E_\theta(x)$ を介さずに直接 $s_\theta(x)$ だけ学習したらいいのでは??
 - スコアマッチング:

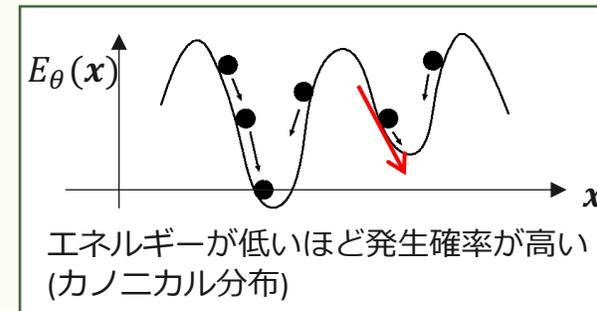
Fisher Divergence
の最小化問題

- データからスコアをフィッティングするLossを計算するための工夫(変形)

$$E_{x \sim p(x)} \left[\frac{1}{2} \|\nabla_x \log p(x) - s_\theta(x)\|_2^2 \right] = E_{x \sim p(x)} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right]$$

ここがわからなくてもOK

筆者曰く、Diffusion Modelは尤度ベース、スコアベースどちらとしても解釈可能とのこと (楽しみ)



(再掲)生成モデル(エネルギーモデル,スコアベースモデル)

- 生成モデルの系譜 (知ってる人向けの説明)
 - "Likelihood-Based" (原文ママ)
 - パラメトリックな確率モデルを考えて、データに対する尤度を最大化してフィッティング
 - Energy Based Models
 - よく統計力学とかでありそうな分布形、 $p(x) = \frac{1}{Z_\theta} \exp(-E_\theta(x))$ を仮定して尤度Lossで $E_\theta(x)$ を学習
 - 統計力学の知見を流用しつつ、尤度に登場する期待値(規格化定数 Z_θ)計算を近似とサンプリングで上手く回避
 - Z_θ を求めず分布 $p(x)$ からサンプリングする方法 (MCMC+ランジュバンダイナミクス)
 - $x^{t+1} = x^t + \eta \nabla_x \log p_\theta(x) + \omega$ (ω : ガウスノイズ, η : ハイパーパラメタ)
 - ※スコア $s_\theta(x) := \nabla_x \log p_\theta(x) = \nabla_x \log(\frac{1}{Z_\theta} \exp(-E_\theta(x))) = \nabla_x \{\log(\exp(-E_\theta(x))) - \log(Z_\theta)\} = -\nabla_x E_\theta(x)$

Score Based Models

- 発想: $E_\theta(x)$ を介さずに直接 $s_\theta(x)$ だけ学習したらいいのでは??
 - スコアマッチング:

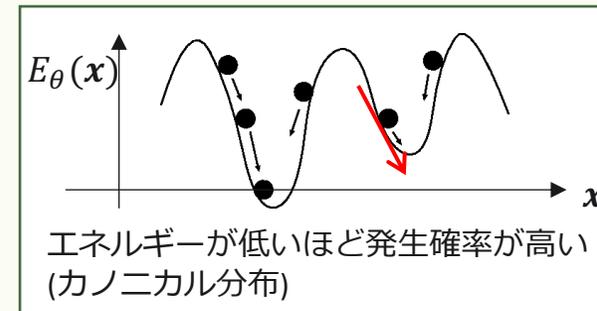
Fisher Divergence
の最小化問題

- データからスコアをフィッティングするLossを計算するための工夫(変形)

$$E_{x \sim p(x)} \left[\frac{1}{2} \|\nabla_x \log p(x) - s_\theta(x)\|_2^2 \right] = E_{x \sim p(x)} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right]$$

ここがわからなくてもOK

筆者曰く、Diffusion Modelは尤度ベース、スコアベースどちらとしても解釈可能とのこと (楽しみ)



(再掲)生成モデル(ディフュージョン)

ランジュバンダイナミクスが逆拡散過程と対応。

また以下のような式も導出できる

$$\mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}} = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

$$\therefore (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t) = -\sqrt{1 - \bar{\alpha}_t} \epsilon_0$$

$$\nabla \log p(\mathbf{x}_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_0$$

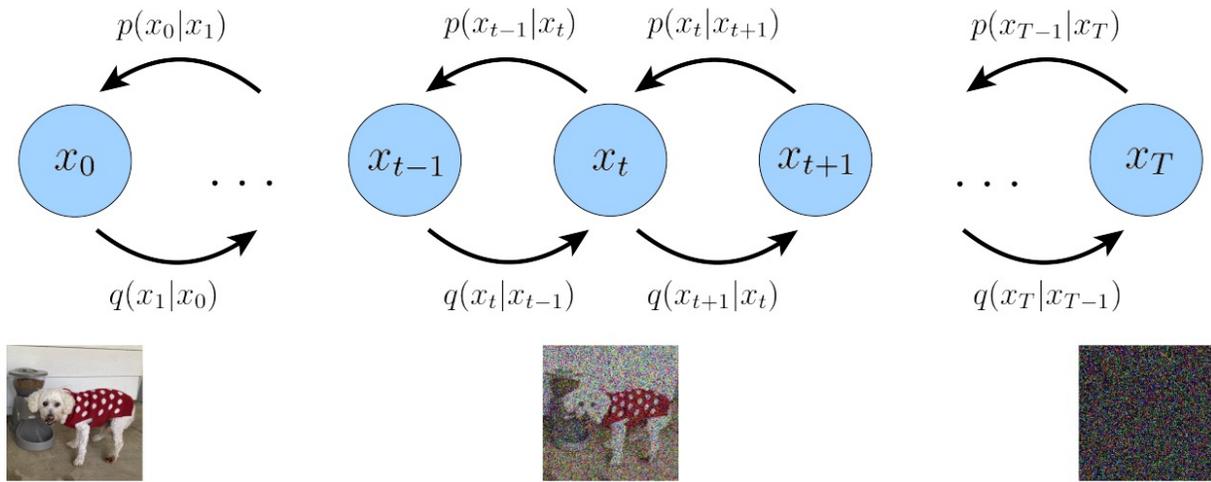


Figure 3: A visual representation of a Variational Diffusion Model; \mathbf{x}_0 represents true data observations such as natural images, \mathbf{x}_T represents pure Gaussian noise, and \mathbf{x}_t is an intermediate noisy version of \mathbf{x}_0 . Each $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ is modeled as a Gaussian distribution that uses the output of the previous state as its mean.

- Z_θ を求めず分布 $p(x)$ からサンプリングする方法 (MCMC+ランジュバンダイナミクス)
 - $\mathbf{x}^{t+1} = \mathbf{x}^t + \eta \nabla_x \log p_\theta(\mathbf{x}) + \omega$ (ω : ガウスノイズ, η : ハイパーパラメタ)
- ※スコア $s_\theta(\mathbf{x}) := \nabla_x \log p_\theta(\mathbf{x}) = \nabla_x \log(\frac{1}{Z_\theta} \exp(-E_\theta(\mathbf{x}))) = \nabla_x \{\log(\exp(-E_\theta(\mathbf{x})) - \log(Z_\theta))\} = -\nabla_x E_\theta(\mathbf{x})$

Score Based Models

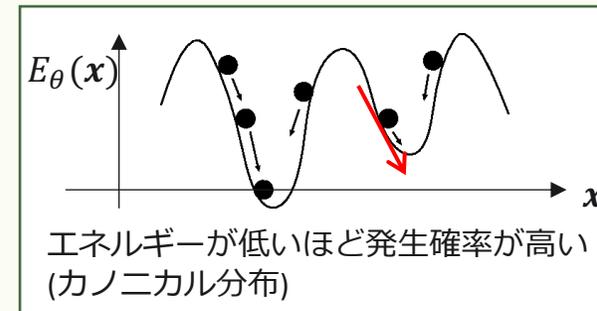
- 発想: $E_\theta(\mathbf{x})$ を介さずに直接 $s_\theta(\mathbf{x})$ だけ学習したらいいのでは??
 - スコアマッチング:

Fisher Divergence
の最小化問題

- データからスコアをフィッティングするLossを計算するための工夫(変形)

$$E_{\mathbf{x} \sim p(\mathbf{x})} \left[\frac{1}{2} \|\nabla_x \log p(\mathbf{x}) - s_\theta(\mathbf{x})\|_2^2 \right] = E_{\mathbf{x} \sim p(\mathbf{x})} \left[\frac{1}{2} \|s_\theta(\mathbf{x})\|_2^2 + \text{tr}(\nabla_x s_\theta(\mathbf{x})) \right]$$

ここがわからなくてもOK



筆者曰く、Diffusion Modelは尤度ベース、スコアベースどちらとしても解釈可能とのこと (楽しみ)

まとめ

- 先ほどのスコア推定としての解釈結果 $\|s_{\theta}(x_t, t) - \nabla \log p(x_t)\|^2$ を利用して、Diffusion Model を Score Based Models として再解釈可能である
- Score Based Models のランジュバンダイナミクスと逆拡散過程(デコーダ)が対応している
- Score Based Models と同様、ランジュバンダイナミクスを使ったサンプリングも可能である
 - つまり言い換えると、デコーダ $p_{\theta}(x_{t-1}|x_t)$ を繰り返し適用したサンプリングが可能である
- 連続時間のダイナミクスとかも考えられるらしい、Neural ODE との組み合わせとか
- あとは解説しなかったですが条件付きサンプリングとか(DALL-E とかで使われてるやつですね)

素朴な疑問

- じゃあここで言うエネルギーってなんやる??
 - データの尤度でしょうか...??

参考) ランジュバンダイナミクス

