

確率数理要論(12)

機械学習応用

- 勾配ランジュバン動力学
- 拡散モデル

• 勾配ランジュバン動力学

Stochastic Gradient Langevin Dynamics (SGLD)

非凸最適化：

$$\min_{x \in \mathbb{R}^d} L(x) = \min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell_i(x) \quad \text{(非凸でも良い)}$$

サンプリング：

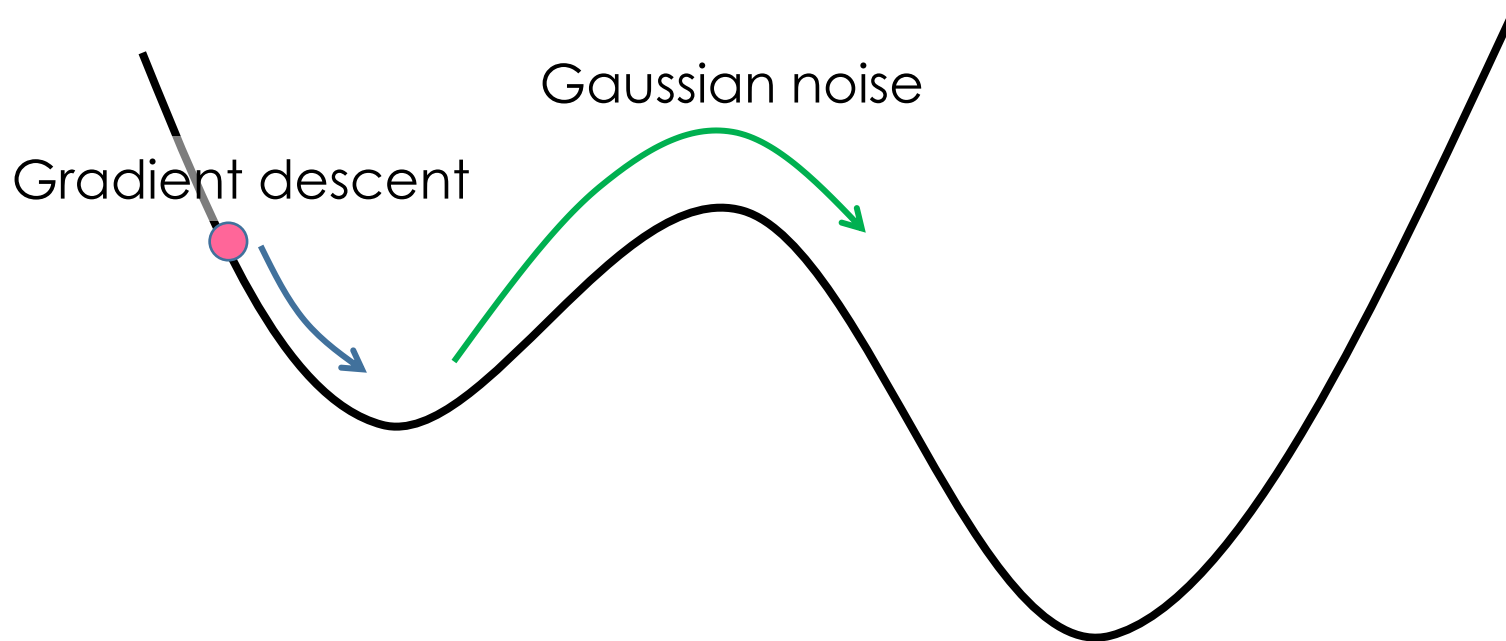
$$X \sim \exp(-\beta L(X)) / Z_{\beta, L}$$

β : 逆温度パラメータ

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t \quad \text{(勾配ランジュバン動力学)}$$

$$\text{定常分布： } \pi \propto \exp(-\beta L(X))$$

[Gelfand and Mitter (1991); Borkar and Mitter (1999); Welling and Teh (2011)]



GLDのFokker-Planck方程式

$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t$$

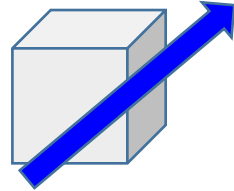
μ_t : X_t の分布 (密度関数であると考えてもらって問題ない)

Fokker-Planck方程式

$$\nabla \cdot [\mu_t \nabla L] = \sum_{j=1}^d \partial_i [\mu_t \partial_i L]$$

$$\partial_t \mu_t = \frac{1}{\beta} \Delta_x \mu_t + \nabla \cdot [\mu_t \nabla L]$$

Mass: $\mu_t(x)$



Vector field: v_t

次のように解釈できる:

$$\partial_t \mu_t = \nabla \cdot \left[\underbrace{\left(\frac{1}{\beta} \nabla \log(\mu_t) + \nabla L \right)}_{-v_t \text{とおく}} \mu_t \right] = -\nabla \cdot [v_t \mu_t]$$

[連続の方程式]

連続の方程式

「連続の方程式」
$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (v_t \rho_t)$$

この方程式の意味

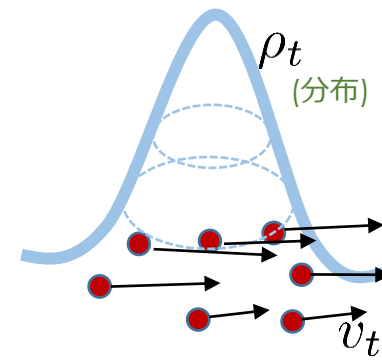
$$\frac{d}{dt} \int f(w) d\rho_t(w) = \int (\nabla f(w))^\top v_t(w) d\rho_t(w)$$

($= - \int f(w) d[\nabla \cdot (v_t \rho_t)]$) ($\forall f$: コンパクトサポート, C^∞ -級)

- ベクトル場 v_t で生成される写像を T_t とする: $\frac{dT_t}{dt}(w) = v_t(T_t(w))$.
- ρ_t は写像 $T_t: R^d \rightarrow R^d$ による ρ_0 の押し出し: $\rho_t = T_{t\#}\rho_0$.
つまり, $w \sim \rho_0$ に対する $T_t(w)$ の分布が ρ_t .

($t = 0$ で導出: $T_0 = I$ (恒等写像))

$$\begin{aligned} \frac{d}{dt} \int f(w) d\rho_t(w) &= \frac{d}{dt} \int f(T_t(w)) d\rho_0(w) \\ &= \int \nabla f(T_0(w))^\top \frac{dT_t(w)}{dt} d\rho_0(w) \\ &= \int \nabla f(w)^\top v_0(w) d\rho_0(w) \quad \text{[連続の方程式]} \end{aligned}$$



連続の方程式

「連続の方程式」
$$\frac{\partial \rho_t}{\partial t} = -\nabla \cdot (v_t \rho_t)$$

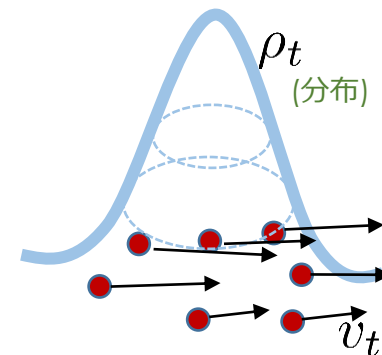
この方程式の意味

$$\frac{d}{dt} \int f(w) d\rho_t(w) = \int (\nabla f(w))^\top v_t(w) d\rho_t(w)$$

($= - \int f(w) d[\nabla \cdot (v_t \rho_t)]$) ($\forall f$: コンパクトサポート, C^∞ -級)

- ベクトル場 v_t で生成される写像を T_t とする: $\frac{dT_t}{dt}(w) = v_t(T_t(w))$.
- ρ_t は写像 $T_t: R^d \rightarrow R^d$ による ρ_0 の押し出し: $\rho_t = T_{t\#}\rho_0$.
つまり, $w \sim \rho_0$ に対する $T_t(w)$ の分布が ρ_t .

(一般の t)
$$\begin{aligned} \frac{d}{dt} \int f(w) d\rho_t(w) &= \frac{d}{dt} \int f(T_t(w)) d\rho_0(w) \\ &= \int \nabla f(T_t(w))^\top \frac{dT_t(w)}{dt} d\rho_0(w) \\ &= \int \nabla f(T_t(w))^\top v_t(T_t(w)) d\rho_0(w) \\ &= \int \nabla f(w)^\top v_t(w) d\rho_t(w). \quad (\text{連続の方程式}) \end{aligned}$$



$$\partial_t \mu_t = \nabla \cdot \left[\left(\frac{1}{\beta} \nabla \log(\mu_t) + \nabla L \right) \mu_t \right] = -\nabla \cdot [v_t \mu_t]$$

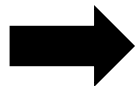
定常分布: $\partial_t \mu_t = 0 \Rightarrow v_t = 0$ (分布がこれ以上動かない)

$$\frac{1}{\beta} \nabla \log(\mu^*) + \nabla L = 0 \Rightarrow \mu^*(x) \propto \exp(-L(x))$$

実は、これは以下の目的関数を最小化するWasserstein勾配流である:

$$\mu^* = \arg \min_{\mu \in \mathcal{P}} \int L(x) d\mu(x) + \frac{1}{\beta} \text{Ent}(\mu) =: \mathcal{L}(\mu)$$

$$(\text{Ent}(\mu) = \int \log(\mu) d\mu)$$



確かにこの最適解は定常分布と等しい:

$$\mu^*(x) \propto \exp(-\beta L(x))$$

$$\begin{aligned}\beta\mathcal{L}(\mu) &= \int \beta L(x) d\mu(x) + \text{Ent}(\mu) && \boxed{\mu^*(x) \propto \exp(-\beta L(x))} \\ &= \int -\log(\mu^*) d\mu + \int \log(\mu) d\mu + (\text{const.}) \\ & && \text{以下, 無視} \\ &= \int \log\left(\frac{\mu}{\mu^*}\right) d\mu = \text{KL}(\mu \parallel \mu^*)\end{aligned}$$

連続の方程式 $\mu_t = -\nabla \cdot [v_t \mu_t]$ に従っているなら

$$\begin{aligned}\frac{d}{dt} \text{KL}(\mu_t \parallel \mu^*) &= \frac{d}{dt} \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \mu_t(x) dx \\ &= \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \partial_t \mu_t(x) dx + \int \frac{\partial_t \mu_t(x)}{\mu_t(x)} \mu_t(x) dx \\ &= \int \log\left(\frac{\mu_t(x)}{\mu^*(x)}\right) \nabla \cdot (-v_t \mu_t(x)) dx \\ &= - \int \langle v_t, \nabla \log(\mu^*) - \nabla \log(\mu_t) \rangle d\mu_t\end{aligned}$$

$$\frac{d}{dt} \text{KL}(\mu_t || \mu^*) = - \int \langle v_t, \nabla \log(\mu^*) - \nabla \log(\mu_t) \rangle d\mu_t$$

特に

$$v_t = - \left(\frac{1}{\beta} \nabla \log(\mu_t) + \nabla L \right) = \frac{1}{\beta} (\nabla \log(\mu^*) - \nabla \log(\mu_t)) \quad (\text{GLD})$$

は最急降下方向で、以下が成り立つ。

$$\partial_t \mu_t = \nabla \cdot \left[\underbrace{\left(\frac{1}{\beta} \nabla \log(\mu_t) + \nabla L \right)}_{=:-v_t} \mu_t \right]$$

$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t || \mu^*) &= -\frac{1}{\beta} \int \|\nabla \log(\mu^*) - \nabla \log(\mu_t)\|^2 d\mu_t \\ &= -\frac{1}{\beta} I(\mu_t || \mu^*) \end{aligned}$$

定常分布 μ^* からのKL-divを最小化するWasserstein勾配流

Fisher divergence:

$$I(\mu || \nu) = \int \|\nabla \log(\nu) - \nabla \log(\mu)\|^2 d\mu$$

μ, ν : 距離空間 (\mathcal{X}, c) 上の確率測度 (通常 \mathcal{X} はPoland空間)

$$W_p(\mu, \nu) = \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} c(x, y)^p d\pi(x, y) \right)^{1/p}$$

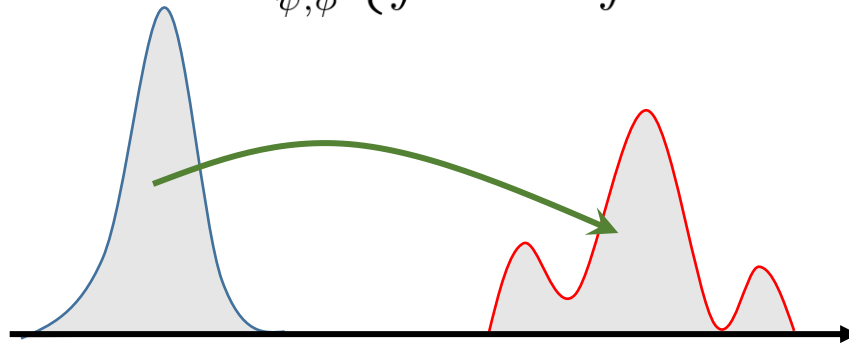
$\Pi(\mu, \nu)$: 周辺分布が μ, ν である $\mathcal{X} \times \mathcal{X}$ 上の同時分布の集合
周辺分布を固定した同時分布の中で最小化

$$(\mathcal{X} = \mathbb{R}^d: c(x, y) = \|x - y\|)$$

- 分布のサポートがずれていてもwell-defined
- 底空間の距離が反映されている
※KL-divergenceは距離が反映されない。

(双対表現: Kantorovich双対)

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y)^p d\pi(x, y) = \sup_{\psi, \phi} \left\{ \int \psi d\mu + \int \phi d\nu \mid \psi(x) + \phi(y) \leq c(x, y)^p \right\}$$



「輸送距離」とも言われる

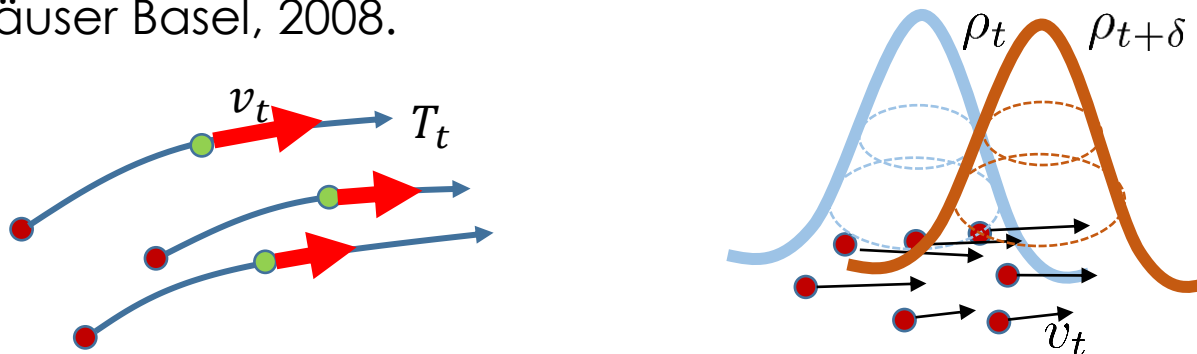
定理

- $\rho_t = T_t \# \rho_0$
- $\frac{dT_t}{dt}(w) = v_t(T_t(w))$
- ある ϕ_t を用いて $v_t = \nabla \phi_t$ と書けるとする。
この時、以下が成り立つ:

$$\lim_{\delta \rightarrow 0} \frac{W_2(\rho_{t+\delta}, (\text{id} + \delta v_t) \# \rho_t)}{\delta} = 0$$

詳細は以下を参照:

Ambrosio, Gigli, and Savaré. *Gradient Flows in Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel, 2008.



Brenierの定理

ρ_0, ρ_1 が確率密度関数を持つ時, 以下が成り立つ:

$$W_2^2(\rho_0, \rho_1) = \inf_{T: T\#\rho_0=\rho_1} \mathbb{E}_{X \sim \rho_0} [\|X - T(X)\|^2]$$

- Infを達成する写像 T^* が存在する.
- しかも, ある凸関数 ψ が存在して $T^*(x) \in \partial\psi(x)$ と書ける.
- この T^* を最適輸送写像という.

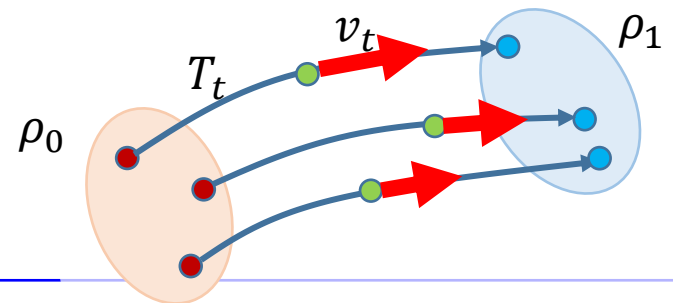
Benamou-Brenier formula (連続の方程式と W_2 距離の関係):

同条件のもと

$$W_2^2(\rho_0, \rho_1) = \inf_{\{v_t\}_t} \int_0^1 \|v_t\|_{L_2(\rho_t)}^2 dt$$

ただし, infは ρ_0 から ρ_1 へ連続の方程式で“繋ぐ”
全ての速度ベクトル場 v_t に関して取る.

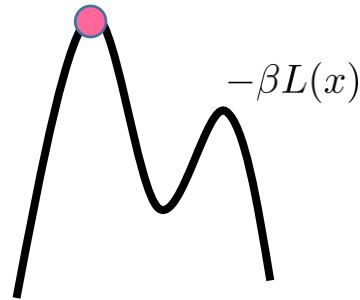
- $\rho_t = T_t\#\rho_0$
- $\frac{dT_t}{dt}(w) = v_t(T_t(w))$



$$dX_t = -\nabla L(X_t)dt + \sqrt{2\beta^{-1}}dB_t$$

適当な条件のもとGLDの定常分布は以下で与えられる:

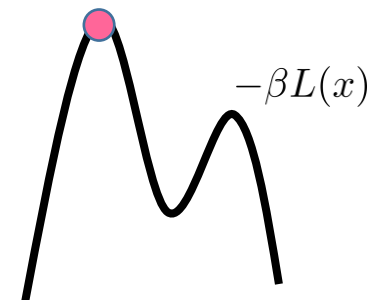
$$\pi_\infty(dx) \propto \exp(-\beta L(x))dx$$



- β が大きければ, 定常分布は最適解周りに集中する.
- 定常分布からのサンプリングにも使える (というかこっちが本来の目的).

定常分布

$$\pi_\infty(dx) \propto \exp(-\beta L(x)) dx$$



対数ソボレフ不等式 (π_∞ の性質):

例:

- 二次関数+有界関数
- Weak Morse型関数

任意の(π_∞ に対して絶対連続な)確率分布 $d\nu = f d\pi_\infty$ に対し,

$$\int f \log(f) d\pi_\infty \leq \frac{c_{LS}}{2} \int \frac{\|\nabla f\|^2}{f} d\pi_\infty$$

$$(D(\nu||\pi_\infty) \leq \frac{c_{LS}}{2} I(\nu||\pi_\infty))$$

$$D(\mu||\nu) = \int \log\left(\frac{d\mu}{d\nu}\right) d\mu, \quad I(\mu||\nu) = \int \left\| \nabla \log \frac{d\mu}{d\nu} \right\|^2 d\mu$$

KL-div Fisher-div

➔ 幾何的エルゴード性

ρ_t : X_t の周辺分布

$$D(\rho_t||\pi_\infty) \leq \exp(-2t/c_{LS}) D(\rho_0||\pi_\infty)$$

定常分布へKL-divergenceの意味で指数オーダーの収束

強凸な場合: $L(x)$ が μ -強凸 $\Rightarrow c_{\text{LS}} \leq 1/(\mu\beta)$

[Bakry and Émery, 1985]

$$L(x) = h(x) + \lambda_1 \|x\|^2$$

Bounded perturbation lemma (Holley-Stroock):

$$|h(x)| \leq B \ (\forall x) \quad \longrightarrow \quad c_{\text{LS}} \leq \frac{1}{2\lambda_1\beta} \exp(4\beta B)$$

[R. Holley and D. Stroock. Logarithmic sobolev inequalities and stochastic Ising models. Journal of statistical physics, 46(5-6):1159-1194, 1987.]

• **散逸的 (dissipative):**

$$\exists m > 0, b \geq 0, \langle x, \nabla L(x) \rangle \geq m\|x\|^2 - b$$

• **平滑性:**

$$\exists M, \|\nabla L(x) - \nabla L(y)\| \leq M\|x - y\|$$



$$\longrightarrow \quad c_{\text{LS}} \leq \frac{2m^2 + 8M^2}{m^2 M \beta} + \left(\frac{6M(d + \beta) + 2}{m} \right) e^{O(\beta+d)}$$

[Raginsky, Rakhlin and Telgarsky, 2017]

過程: L は M -平滑: $\exists M, \|\nabla L(x) - \nabla L(y)\| \leq M\|x - y\|$

定理

[Vempala and Wibisono, 2019]

ν_k : Marginal distribution of X_k (discrete time dynamics)

$$D(\nu_k || \pi_\infty) \lesssim \exp(-k\eta/c_{LS})D(\nu_0 || \pi_\infty) + 8c_{LS}dM^2\eta$$

定理 (informal)

散逸性と平滑性の条件のもと (and other technical condition),

$$E[L(X_k)] - L(X^*) \lesssim \exp(-ck\eta/c_{LS}) + c_{c_{LS},\beta,d}\eta + \frac{d \log(\beta + 1)}{\beta}$$

幾何的エルゴード性

時間離散化の誤差

$E_{\pi_\infty}[L(X)] - L(X^*)$

where $c, c_{c_{LS},\beta,d} > 0$ are constants.

定常分布が最適解まわりにどれだけ集中しているか

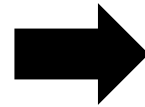
[Raginsky, Rakhlin and Telgarsky, 2017; Xu, Chen, Zou, and Gu, 2018; Erdogdu, Mackey and Shamir, 2018]

- 逆温度パラメータ β が十分大きければ、目的関数が非凸でも最適解の近くに到達できる。
- ただし、一般には対数ソボレフ不等式は β に指数的に依存することに注意。
(そうでない場合もある: 強凸目的関数, Weak Morse関数)

拡散モデル

文章による説明から画像を生成するモデル

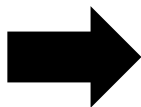
「An astronaut riding a horse in a photorealistic style」



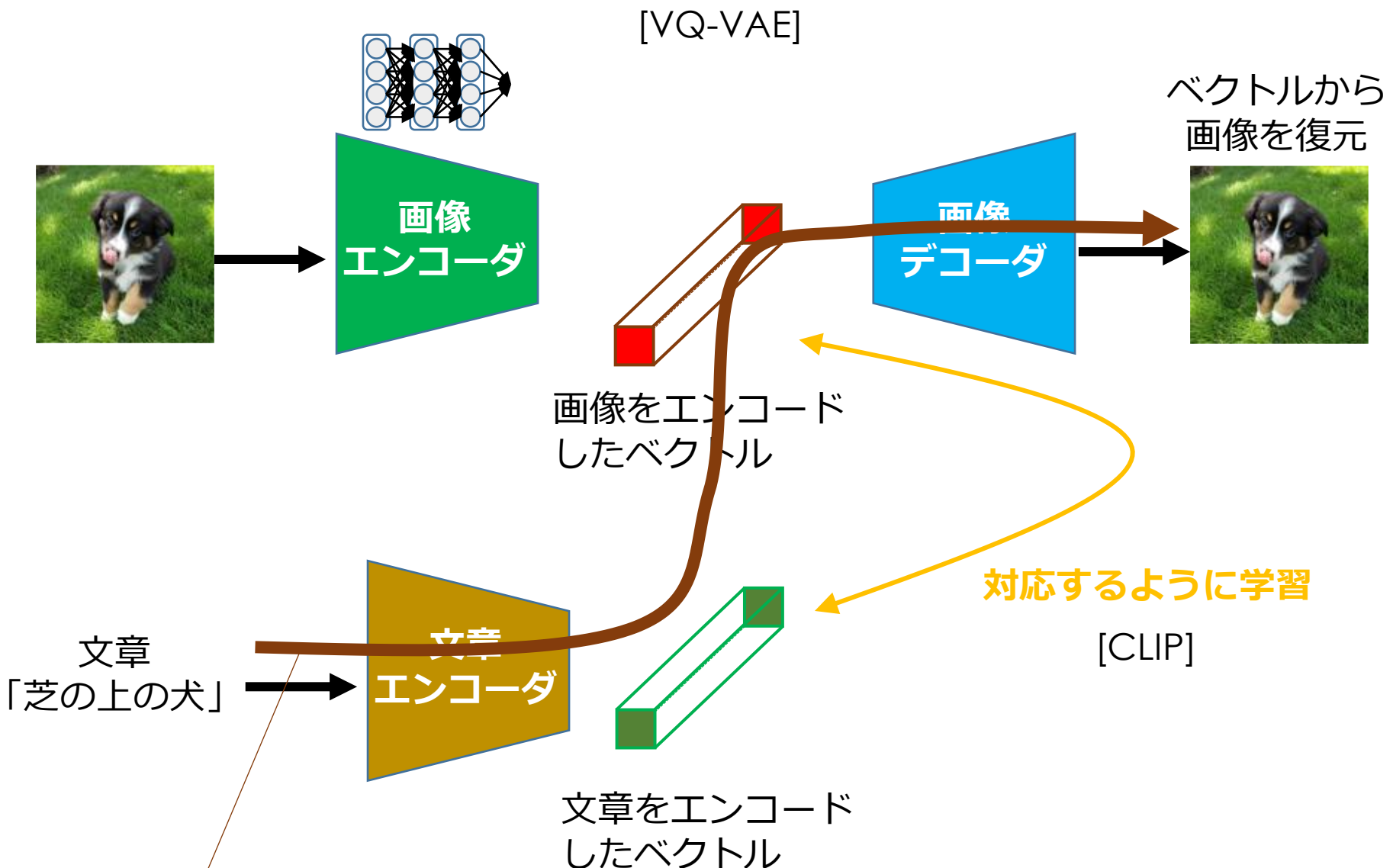
DALL·E: [Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, Ilya Sutskever: Zero-Shot Text-to-Image Generation. ICML2021.]

DALL·E2:[Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, Mark Chen: Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125]

「Teddy bears shopping for groceries in the style of ukiyo-e」



特徴ベクトルの役割: DALL-Eで説明



ベクトルから
画像を復元

画像をエンコード
したベクトル

対応するように学習
[CLIP]

文章をエンコード
したベクトル

画像生成時

同様の考え方が翻訳など、ほとんどの
深層学習モデルで使われている。



Jason Allen "Théâtre D'opéra Spatial" generated by **Midjourney**. Colorado State Fair's fine art competition, 1st prize in digital art category

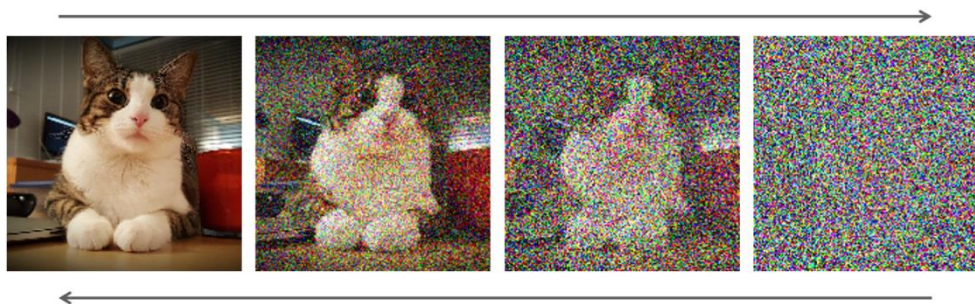


Generated by NovelAI

デコーダー：拡散モデル (GLIDE)

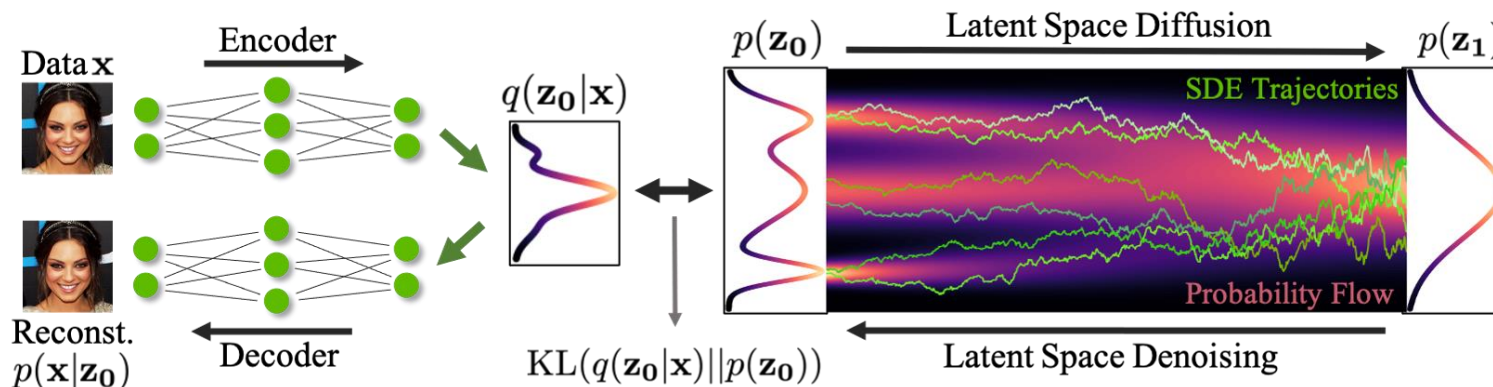
順モデル：所望の分布からノイズに変換

$$dX_t = f(x, t)dt + dB_t$$



$$dX_t = [f(x, t) - \nabla \log(p_t(X_t))]dt + dB_t$$

逆モデル：ノイズから所望の分布へ変換



- 順過程：ノイズの分布が標準正規分布の場合

標準正規分布へ向かう勾配ランジュバン動力学

$$\mu_\infty \propto \exp\left(-\frac{\|x\|^2}{2}\right) = \exp(-L(x))$$

$$\Rightarrow L(x) = \frac{\|x\|^2}{2}$$

$$\Rightarrow dX_t = -X_t dt + \sqrt{2}dB_t$$

(OU-過程)

Fokker-Planck方程式： $\partial_t \mu_t = \Delta \mu_t + \nabla \cdot (x \mu_t)$

$$\partial_t \mu_t = \frac{1}{\beta} \Delta_x \mu_t + \nabla \cdot [\mu_t \nabla L]$$

- 逆過程 (reverse SDE) :

\tilde{t} を逆向き時間として ($\infty \rightarrow 0$ へ向かう)

$$dX_{\tilde{t}} = (X_{\tilde{t}} + 2\nabla \log(\mu_{\tilde{t}}(x)))d\tilde{t} + \sqrt{2}dB_{\tilde{t}}$$

$b(x) = -x = -\nabla L(x)$ として次のページで導出.

逆向きSDEのFK-方程式

$$\partial_s p(s, x|t, y) = \partial_s \left[p(t, y|s, x) \frac{p(s, x)}{p(t, y)} \right] \quad (s < t \text{を想定})$$

$$= \partial_s p(t, y|s, x) \frac{p(s, x)}{p(t, y)} + p(t, y|s, x) \frac{\partial_s p(s, x)}{p(t, y)}$$

$$= \left[\underline{-b^\top \nabla_x p(t, y|s, x)} - \underline{\Delta_x p(t, y|s, x)} \right] \frac{p(s, x)}{p(t, y)}$$

$$+ p(t, y|s, x) \frac{\underline{-\nabla_x \cdot (bp(s, x))} + \cancel{\Delta_x p(s, x)}}{p(t, y)}$$

$$\text{--- } b^\top \nabla_x p(t, y|s, x) \frac{p(s, x)}{p(t, y)} = b^\top \nabla_x p(s, x|t, y) - p(t, y|s, x) \frac{\cancel{b^\top \nabla_x p(s, x)}}{p(t, y)}$$

$$\text{--- } \Delta_x p(t, y|s, x) \frac{p(s, x)}{p(t, y)} = \Delta_x p(s, x|t, y) - 2 \nabla_x^\top p(t, y|s, x) \frac{\nabla_x p(s, x)}{p(t, y)} - p(t, y|s, x) \frac{\Delta_x p(s, x)}{p(t, y)}$$

$$= \Delta_x p(s, x|t, y) - 2 \nabla_x^\top p(s, x|t, y) \nabla_x \log(p(s, x))$$

$$+ 2p(s, x|t, y) \|\nabla_x p(s, x)\|^2 - p(t, y|s, x) \frac{\Delta_x p(s, x)}{p(t, y)}$$

$$= \Delta_x p(s, x|t, y) - \underline{2 \nabla_x^\top p(s, x|t, y) \nabla_x \log(p(s, x))}$$

$$\underline{-2p(s, x|t, y) \Delta_x \log(p(s, x))} + \cancel{p(t, y|s, x) \frac{\Delta_x p(s, x)}{p(t, y)}}$$

$$\text{--- } p(t, y|s, x) \frac{-\nabla_x \cdot (bp(s, x))}{p(t, y)} = -(\underline{\nabla_x \cdot b} + \cancel{b^\top \nabla_x \log(p(s, x))}) p(s, x|t, y)$$

まとめると,

$$\partial_s p(s, x|t, y) = -\nabla_x \cdot [(b - 2\nabla_x \log(p(s, x)))p(s, x|t, y)] - \Delta_x p(s, x|t, y)$$

時間を反転させて, $d\tilde{s} \leftarrow -ds$ とすると,

$$\partial_{\tilde{s}} p(\tilde{s}, x|t, y) = \nabla_x \cdot [(b - 2\nabla_x \log(p(\tilde{s}, x)))p(\tilde{s}, x|t, y)] + \Delta_x p(\tilde{s}, x|t, y)$$

これはドリフト項が

$$-(b - 2\nabla_x \log(p(s, x))) = x + 2\nabla_x \log(\mu_s(x))$$

の拡散過程の前向き方程式に他ならない.

$\tilde{s} \rightarrow 0$ とすることで, 時刻0における分布を得ることができる.

つまり, ドリフト項をデータから推定し, 逆過程を走らせることでデータの分布からのサンプリングができるようになる.