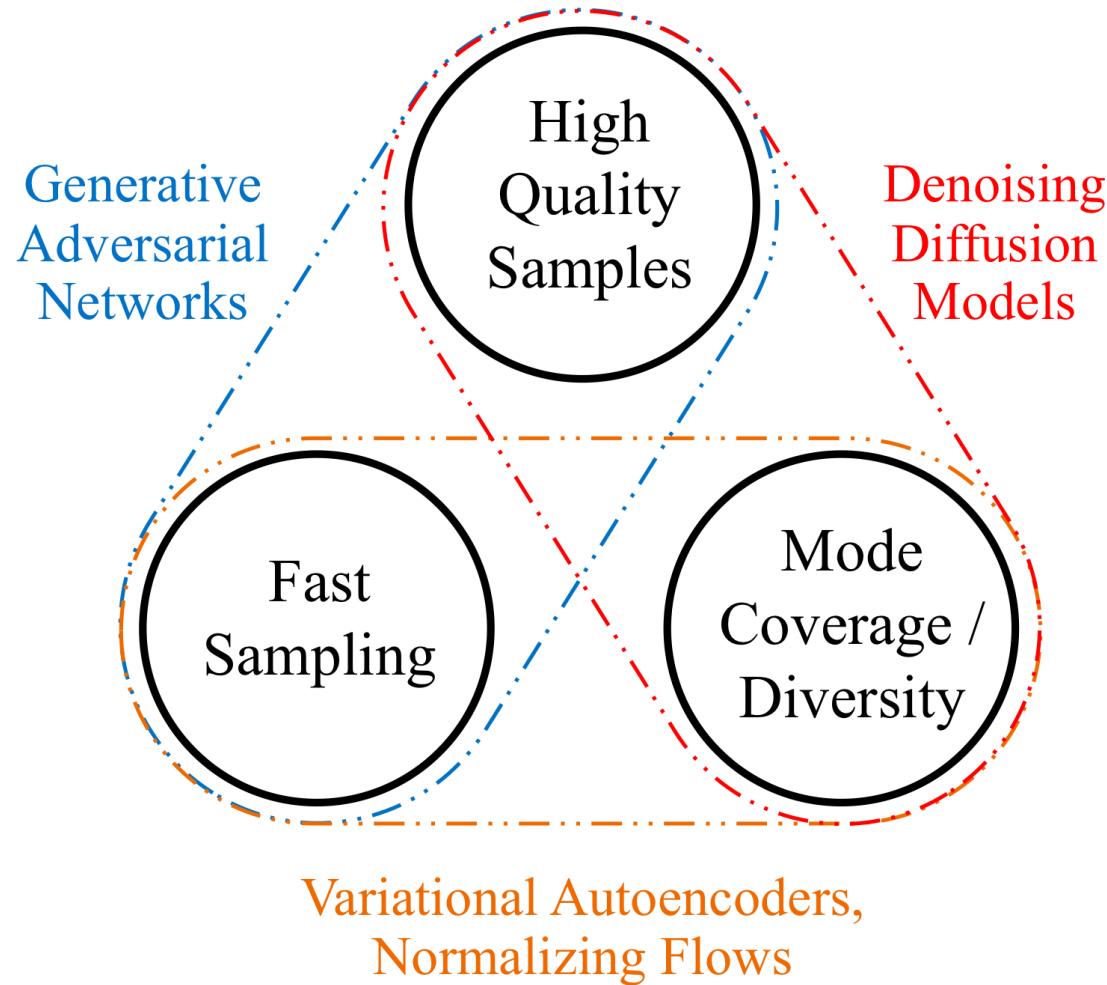


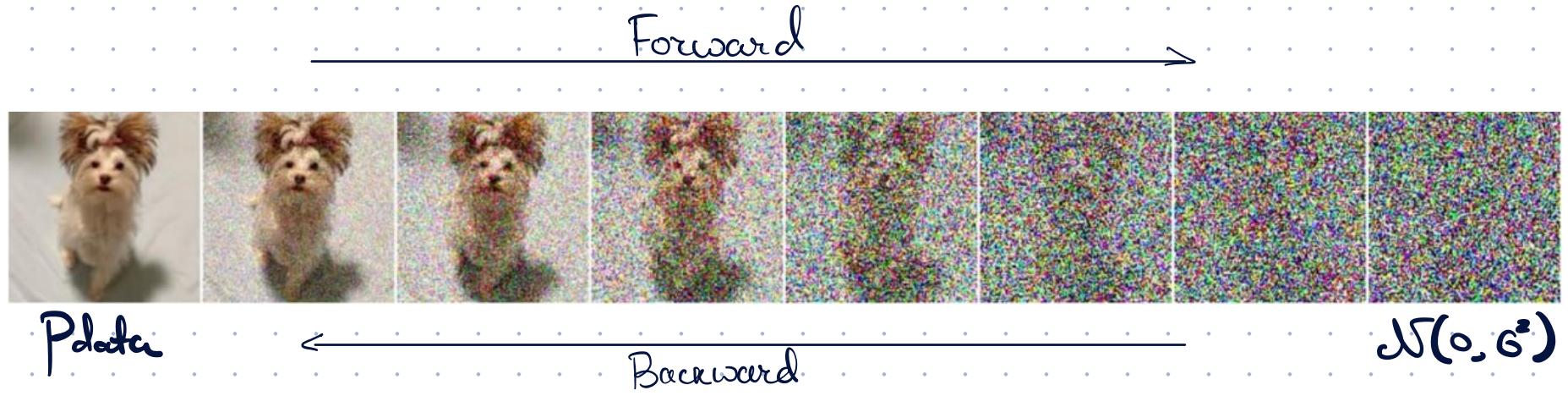
Diffusion models

Denis Rakitin

Generative trilemma



Diffusion



Idea: define forward noising process and try to reverse it.

Discrete time

- Variance preserving (VP) process: $t=1, \dots, T$

$$X_{t+1} = \sqrt{1-\beta_t} X_t + \sqrt{\beta_t} \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, I)$$

$$(I \nparallel \text{Var } X_0 : I \Rightarrow \text{Var } X_t = I)$$

$$P_{t|0}(x_t | X_0) = \mathcal{N}(x_t | \alpha_t \cdot x_0, \sigma_t^2 I) \quad \begin{array}{l} \alpha_t \rightarrow 0 \\ \sigma_t^2 \rightarrow 1 \end{array}$$
$$P_T(x_T) \approx \mathcal{N}(x_T | 0, I)$$

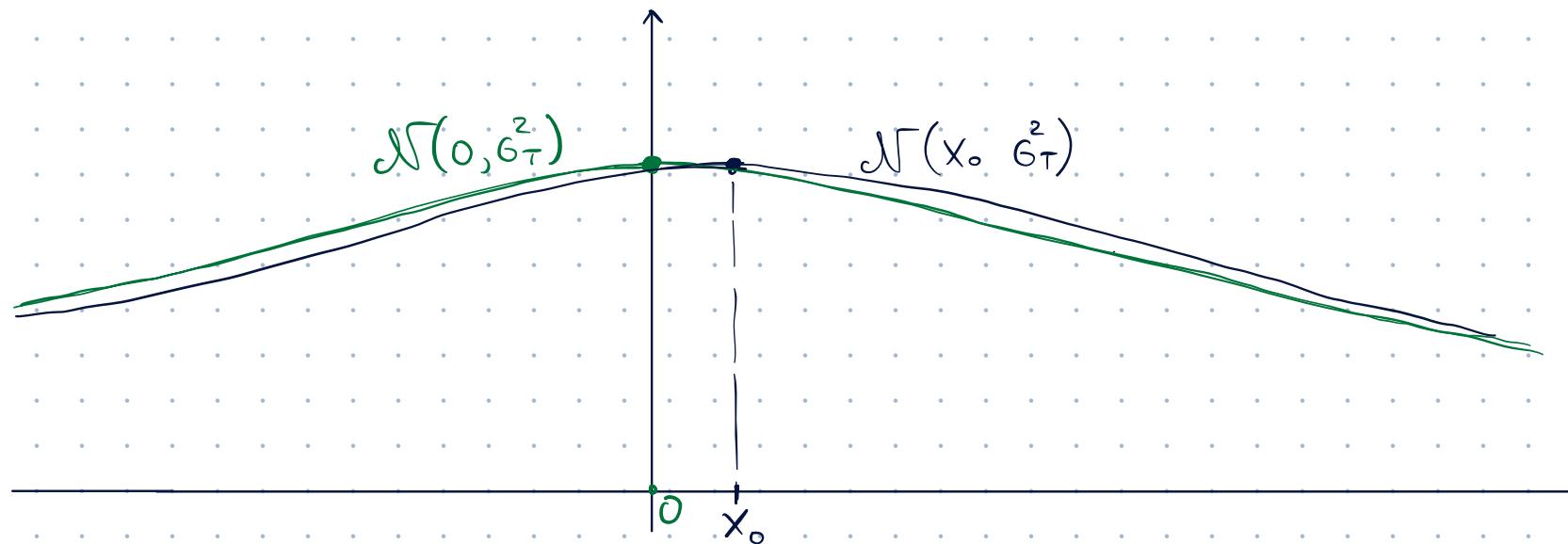
Discrete time

- Variance exploding (VE) process: $t=1, \dots, T$

$$x_{t+1} = x_t + g(t) \cdot \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, I)$$

$$P_{t|0}(x_t|x_0) = \mathcal{N}(x_t|x_0, G_t^2 I), \quad G_t^2 = \sum_{s \leq t} g^2(s)$$

If G_T^2 is big $\Rightarrow P_T(x_T) \approx \mathcal{N}(0, G_T^2 I)$



Basic idea (VE)

$$P_{t+1|t}(X_{t+1} | X_t) = \mathcal{N}(X_{t+1} | X_t, g^2(t) I)$$

$$P_T(X_T) \approx \mathcal{N}(X_T | 0, G_T^2), \quad G_T^2 = \sum_{t=1}^T g^2(t).$$

Ideally:

- Sample $X_T \sim \mathcal{N}(0, G_T^2)$ intractable ↓
- Iterate $X_{t-1} \sim P_{t-1|t}(X_{t-1} | X_t) = \frac{P_{t|t-1}(X_t | X_{t-1}) P_{t-1}(X_{t-1})}{P_t(X_t)}$ $P_t(X_t)$
↑
intractable

Basic idea (VE)

$$P_{t+1|t}(X_{t+1}|X_t) = \mathcal{N}(X_{t+1}|X_t, g^2(t)I)$$

$$P_T(X_T) \approx \mathcal{N}(X_T|0, G_T^2), \quad G_T^2 = \sum_{t=1}^T g^2(t).$$

Practical replacement:

- $P_{t-1|t,0}(X_{t-1}|X_t, X_0) = \frac{P_{t|t-1}(X_t|X_{t-1}) P_{t-1|0}(X_{t-1}|X_0)}{P_{t|0}(X_t|X_0)}$

gaussian
gaussian

 \nwarrow
 \nearrow

- Make (maybe not ideal) prediction $\hat{X}_0 = D(X_t, t)$
- Sample $X_{t-1} \sim P_{t-1|t,0}(X_{t-1}|X_t, \hat{X}_0)$

Train NN: $\sum_{t=1}^T \omega_t \mathbb{E} \|D_\theta(X_t, t) - X_0\|^2 \rightarrow \min_\theta$

Basic idea (summary)

Training: $\sum_{t=1}^T \omega_t \mathbb{E} \|D_\theta(x_t, t) - x_0\|^2 \rightarrow \min_\theta$

Sampling:

- Claim: $P_{t-1|t,0}(x_{t-1} | x_t, \hat{x}_0) = \mathcal{N}(x_{t-1} | \mu_t, \Sigma_t)$
- $$\mu_t = x_t \cdot \left(1 - \frac{g^2(t)}{6_t}\right) + \hat{x}_0 \cdot \frac{g^2(t)}{6_t}, \quad \Sigma_t = \boxed{\frac{6_{t-1}}{6_t^2}} g^2(t) \cdot I$$

- Sample $x_t \sim \mathcal{N}(0, 6_t^2)$

- Iterate

$$x_{t-1} = x_t \left(1 - \frac{g^2(t)}{6_t^2}\right) + D_\theta(x_t, t) \cdot \frac{g^2(t)}{6_t^2} + g(t) \varepsilon_t$$

Links

<https://arxiv.org/abs/1907.05600> (VE process)

<https://arxiv.org/abs/2006.11239> (VP process, method derivation)

ODEs, SDEs, PDEs

ODEs

$$dX_t = f(X_t, t) dt \quad / \quad \dot{X} = f(X, t)$$

Euler scheme: $X_{t+h} \approx X_t + h \cdot f(X_t, t)$

- Part of VP process: $X_{t+1} = \sqrt{1 - \beta_t^2} X_t + \sqrt{\beta_t^2} \varepsilon_t$

- Cont. modification:

$$X_{t+h} = \sqrt{1 - \beta_t^2 h} X_t = \left(1 - h \frac{\beta_t}{2} + \bar{o}(h) \right) X_t$$

$$\frac{X_{t+h} - X_t}{h} = \frac{-h \frac{\beta_t}{2} + \bar{o}(h)}{h} \cdot X_t$$



$$\dot{X}_t$$



$$- \frac{\beta_t}{2} X_t$$

$$\text{Solution: } X_t = X_0 \cdot \exp \left(-\frac{1}{2} \int_0^t \beta_s ds \right)$$

Wiener process

VE process with $g(t) = 1$:

$$X_{t+1} = X_t + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, I) \text{ ind. with } X_t.$$

- Let $X_0 = 0$
- Independent increments:
 $t_1 < \dots < t_m \Rightarrow X_{t_1}, X_{t_2} - X_{t_1}, \dots, X_{t_m} - X_{t_{m-1}}$ are indep.
- $X_t - X_s = \sum_{i=s+1}^t \varepsilon_i \sim \mathcal{N}(0, (t-s)I)$.

Def Wiener process is a continuous-time process

- $\omega_0 = 0$
- It has indep. increments
- $\omega_t - \omega_s \sim \mathcal{N}(0, (t-s)I)$
- It has continuous trajectories

SDEs

- VE process: $X_{t+1} = X_t + g(t) \varepsilon_t$
- VP process: $X_{t+1} = \sqrt{1 - \beta_t^2} X_t + \sqrt{\beta_t^2} \varepsilon_t$

Continuous limits?

Def $dX_t = f(X_t, t) dt + G(X_t, t) d\omega_t$, if
 $X_{t+h} \approx X_t + f(X_t, t) \cdot h + G(X_t, t) (\omega_{t+h} - \omega_t)$
indep. with X_t

Euler scheme:

$X_{t+h} \approx X_t + f(X_t, t) \cdot h + G(X_t, t) \cdot \sqrt{h} \cdot \varepsilon_t$,
 $\varepsilon_t \sim \mathcal{N}(0, I)$ - indep. with X_t

SDEs

Discrete VE: $X_{t+1} = X_t + g(t) \varepsilon_t$

$$\Leftrightarrow X_t + g(t)(\omega_{t+1} - \omega_t)$$

Continuous VE: $X_{t+h} = X_t + g(t)(\omega_{t+h} - \omega_t)$

$$VE-SDE: dX_t = g(t) d\omega_t$$

$$X_t = X_0 + \int_0^t g(s) d\omega_s$$

$$(P_{X_t|X_0}(x_t | x_0) = \text{[Redacted]})$$

SDEs

$$\text{Discrete VP: } X_{t+1} = \sqrt{1-\beta_t} X_t + \sqrt{\beta_t} \varepsilon_t$$

$$\text{Continuous VP: } X_{t+h} = \sqrt{1-\beta_t h} X_t + \sqrt{\beta_t h} \varepsilon_t$$

$$\Leftrightarrow \left(1 - h \frac{\beta_t}{2} + O(h) \right) X_t + \sqrt{\beta_t} (\omega_{t+h} - \omega_t)$$

$$\text{VP-SDE: } dX_t = -\frac{\beta_t}{2} X_t dt + \sqrt{\beta_t} d\omega_t$$

$$X_t = d_t X_0 + d_t \int_0^t d_s^{-\frac{1}{2}} \sqrt{\beta_s} d\omega_s, \quad d_t = \exp \left(- \int_0^t \beta_s ds \right)$$

$$(P_{X_t|X_0}(x_t|x_0) = \mathcal{N}(x_t|x_0, d_t, (1-d_t)I))$$

Backward process

VP/VE (simple)

Forwarded: $\left\{ \begin{array}{l} dX_t = f(X_t, t)dt + g(t)dW_t \\ X_0 \sim P_{\text{data}} \end{array} \right. \leftarrow \text{difficult}$

s.t. $P_1 := P_{X_1}$ is simple and known (s.t. $\mathcal{N}(0, \sigma^2 I)$).

Goal: define $X_t^{(b)}$: $\left\{ \begin{array}{l} dX_t^{(b)} = f^{(b)}(X_t^{(b)}, t)dt + g^{(b)}(t)dW_t \\ X_0^{(b)} \sim P_1 \end{array} \right. \leftarrow \text{difficult}$

s.t. $q_t := P_{X_t^{(b)}}$ equals $P_{t-t} = P_{X_{t-t}}$.

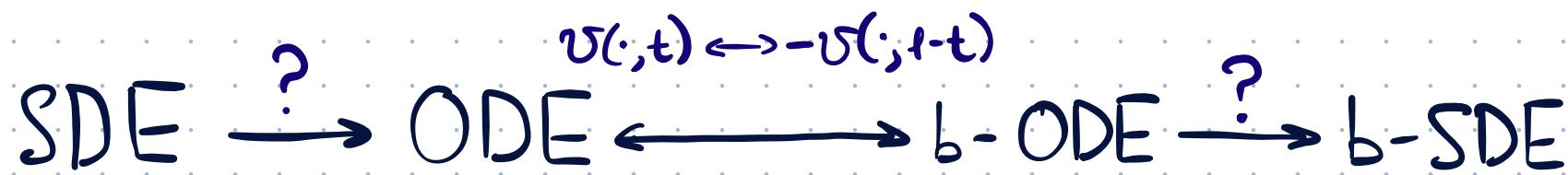
Backward ODE is easy

Let y_t satisfy $\begin{cases} dy_t = v(y_t, t)dt \\ y_0 = y_0 \end{cases}$

Define $y_t^{(b)}$ as $\begin{cases} dy_t^{(b)} = -v(y_t^{(b)}, t-t)dt \\ y_0^{(b)} = y_1 \end{cases}$

Then, $y_t^{(b)} = y_{t-t}$ $\forall t$.

Backward SDE scheme



Goal : $\begin{cases} dX_t = f(X_t, t)dt + g(t)dW_t \\ X_0 \sim P_0 \end{cases}$

↑ ↓ save marginals: $P_{X_t} = P_{Y_t}$

$$\begin{cases} dY_t = \sigma(Y_t, t)dt \\ Y_0 \sim P_0 \end{cases}$$

How? Study evolution of density.

Evolution of density

- $dX_t = f(X_t, t) dt$

$$\frac{\partial}{\partial t} P_t(x) = - \frac{\partial}{\partial x} (P_t(x) f(x, t))$$

Continuity
equation

- $dX_t = f(X_t, t) dt + g(t) dW_t$

$$\frac{\partial}{\partial t} P_t(x) = - \frac{\partial}{\partial x} (P_t(x) f(x, t)) + \frac{g^2(t)}{2} \frac{\partial^2}{\partial x^2} P_t(x)$$

Fokker-Planck (-Kolmogorov) equation

Interpretation

Conservation laws

thanks to S.V. Shaposhnikov
and T.I. Krasovitsky
for their PDE course

Let $u(x, t)$ be density of a substance at time t .

Amount of substance in $[x_0, x_0 + \Delta x]$ is

$$\int_{x_0}^{x_0 + \Delta x} u(x, t) dx.$$

Change from time t_0 to $t_0 + \Delta t$ is

$$\int_{x_0}^{x_0 + \Delta x} u(x, t_0 + \Delta t) dx - \int_{x_0}^{x_0 + \Delta x} u(x, t_0) dx.$$

Interpretation

Let $F(x, t)$ be flow of the substance, i.e. the amount of flowed substance in the point x_0 from t_0 to $t_0 + \Delta t$ is

$$\int_{t_0}^{t_0 + \Delta t} F(x, t) dt.$$

Then,

$$\int_{x_0}^{x_0 + \Delta x} u(x, t_0 + \Delta t) - u(x, t_0) dx = \int_{t_0}^{t_0 + \Delta t} F(x_0, t) - F(x_0 + \Delta x, t) dt$$

$$\int_{x_0}^{x_0 + \Delta x} \int_{t_0}^{t_0 + \Delta t} \frac{\partial}{\partial t} u(x, t) dx dt = - \int_{x_0}^{x_0 + \Delta x} \int_{t_0}^{t_0 + \Delta t} \frac{\partial}{\partial x} F(x, t) dx dt$$

Conservation laws

$$\frac{\partial}{\partial t} u(x, t) = - \frac{\partial}{\partial x} F(x, t)$$

General FPE:

$$\frac{\partial}{\partial t} P_t(x) = - \frac{\partial}{\partial x} (P_t(x) f(x, t)) + \frac{g^2(t)}{2} \frac{\partial^2}{\partial x^2} P_t(x)$$

- $dX_t = f(X_t, t)dt$ $F(x, t) = f(x, t) \cdot P_t(x)$
- $dX_t = g(t)dW_t$ $F(x, t) = - \frac{g^2(t)}{2} \frac{\partial}{\partial x} P_t(x)$
- $dX_t = f(X_t, t)dt + g(t)dW_t$ $F(x, t) = f(x, t) \cdot P_t(x) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} P_t(x)$

Uniqueness

Thm Under regularity conditions, FPE

$$\begin{cases} \frac{\partial}{\partial t} P_t(x) = - \frac{\partial}{\partial x} (P_t(x) f(x, t)) + \frac{g^2(t)}{2} \frac{\partial^2}{\partial x^2} P_t(x) \\ P_0(x) = \hat{P}(x) \end{cases}$$

has unique solution, i.e. \Leftrightarrow corresponds to evolution of density of SDE

$$\begin{cases} dX_t = f(X_t, t)dt + g(t)dW_t \\ X_0 \sim \hat{P} \end{cases}$$

SDE to ODE

$$dX_t = f(X_t, t) dt + g(t) dW_t$$

$$\frac{\partial}{\partial t} P_t(x) \downarrow = - \frac{\partial}{\partial x} (P_t(x) f(x, t)) + \frac{g^2(t)}{2} \frac{\partial^2}{\partial x^2} P_t(x)$$

$$= - \frac{\partial}{\partial x} \left(P_t(x) f(x, t) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} P_t(x) \right)$$

$$\uparrow = - \frac{\partial}{\partial x} \left(P_t(x) [f(x, t) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(x)] \right)$$

$$dY_t = \left(f(Y_t, t) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(Y_t) \right) dt$$

- corresponding forward ODE.

ODE 2 b-ODE

$$dY_t = \nu(Y_t, t) dt, \quad \nu(x, t) = f(x, t) - \\ - \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(x)$$



$$dY_t^{(b)} = -\nu(Y_t^{(b)}, t-t) dt = \\ = \left(-f(Y_t^{(b)}, t-t) + \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_{t-t}(Y_t^{(b)}) \right) dt$$

- backward ODE

b-ODE 2 b-SDE

$$dY_t^{(b)} = -\underline{\nu(Y_t^{(b)}, t-t)} dt \quad (q_t := P_{Y_t^{(b)}} = P_{Y_{t-t}} = P_{t-t})$$

↓

$$\frac{\partial}{\partial t} q_t(x) = - \frac{\partial}{\partial x} (q_t(x) \underline{(-\nu(x, t-t))})$$

$$= - \frac{\partial}{\partial x} \left(q_t(x) \left(-\nu(x, t-t) + \frac{g^2(1-t)}{2} \frac{\partial}{\partial x} \log q_t(x) \right) \right) + \\ + \frac{g^2(1-t)}{2} \frac{\partial^2}{\partial x^2} q_t(x)$$

↑

$$dX_t^{(b)} = \left(-\nu(X_t^{(b)}, t-t) + \frac{g^2(1-t)}{2} \frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)}) \right) dt \\ + g(1-t) dW_t$$

b-ODE 2 b-SDE

$$dX_t^{(b)} = \left(-v(X_t^{(b)}, t-t) + \frac{g^2(t-t)}{2} \frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)}) \right) dt + g(t-t) dW_t$$

Recall $v(x, t) = f(x, t) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(x)$

$$dX_t^{(b)} = \left(-f(X_t^{(b)}, t-t) + g^2(t-t) \frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)}) \right) dt + g(t-t) dW_t$$

Summary

$$dX_t = f(X_t, t) dt + g(t) d\omega_t$$



$$dY_t = \left(f(Y_t, t) - \frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(Y_t) \right) dt$$



$$dY_t^{(b)} = \left(-f(Y_t^{(b)}, t-t) + \frac{g^2(1-t)}{2} \frac{\partial}{\partial x} \log P_{t-t}(Y_t^{(b)}) \right) dt$$



$$dX_t^{(b)} = \left(-f(X_t^{(b)}, t-t) + g^2(1-t) \frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)}) \right) dt + g(1-t) d\omega_t$$

More precise result

$$dX_t = f(X_t, t)dt + g(t)dW_t$$



$$\begin{aligned} dX_t^{(b)} = & \left(-f(X_t^{(b)}, t-t) + g^2(t-t) \frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)}) \right) dt \\ & + g(t-t) dW_t \end{aligned}$$

Thm (Anderson)

Under reg. conditions, if $X_0^{(b)} \sim P_{X_0^{(b)}}$, then
for all measurable $B \in \mathcal{G}_l(C[0,1])$
 $P((X_{t-t}, t \in [0,1]) \in B) = P((X_t^{(b)}, t \in [0,1]) \in B).$

VE-SDE

$$dX_t = g(t) dW_t$$

$$dY_t = -\frac{g^2(t)}{2} \frac{\partial}{\partial x} \log P_t(Y_t) dt$$

$$dY_t^{(b)} = \frac{g^2(1-t)}{2} \frac{\partial}{\partial x} \log P_{t-t}(Y_t^{(b)}) dt$$

$$dX_t^{(b)} = g^2(1-t) \frac{\partial}{\partial x} \log P_{t-t}(Y_t^{(b)}) dt + g(1-t) dW_t$$

Sampling

Both b-ODE $Y_t^{(b)}$ and b-SDE $X_t^{(b)}$ generate the reversed dynamics P_{t-t} , given $\begin{cases} X_0 \sim P_1 \\ Y_0 \sim P_1 \end{cases}$.

$$dY_t^{(b)} = \left(-f(Y_t^{(b)}, t-t) + \frac{g^2(t)}{2} \underbrace{\frac{\partial}{\partial x} \log P_{t-t}(Y_t^{(b)})}_{\text{green underline}} \right) dt$$

$$\begin{aligned} dX_t^{(b)} = & \left(-f(X_t^{(b)}, t-t) + g^2(t) \underbrace{\frac{\partial}{\partial x} \log P_{t-t}(X_t^{(b)})}_{\text{green underline}} \right) dt \\ & + g(t-t) dW_t \end{aligned}$$

Need to know $\frac{\partial}{\partial x} \log P_t(x) \leftarrow$ learn by neural net.

Summary

- # • Fokker-Planck equation

$$\frac{\partial}{\partial t} P_t(x) = - \frac{\partial}{\partial x} (P_t(x) f(x, t)) + \frac{g^2(t)}{2} \frac{\partial^2}{\partial x^2} P_t(x)$$

- $SDE \rightarrow ODE \rightarrow ODE \rightarrow SDE$
 \uparrow
 $V \xrightarrow{\uparrow} -V$
 $FPE \rightarrow CE$
 $CE \rightarrow FPE$

$$\bullet dX_t^{(b)} = \left(-f(X_t^{(b)}, t) + g^2(t) \frac{\partial}{\partial x} \log p_{t-t}(X_t^{(b)}) \right) dt + g(t) dW_t$$

Links

[http://wiki.cs.hse.ru/Уравнения_с_частными_производными_\(2022-2023\)](http://wiki.cs.hse.ru/Уравнения_с_частными_производными_(2022-2023))

<https://arxiv.org/abs/2011.13456> (*SDE diffusion models*)

Continuous-time diffusion models

Score function

$$P_t(x) = \int P_{t|0}(x|x_0) P_0(x_0) dx_0$$

easy \nearrow difficult ($= P_{\text{data}}$)

$$(\text{VE-SDE} : P_{t|0}(x|x_0) = \mathcal{N}(x|x_0, \sigma_t^2 I))$$

$$\frac{\partial}{\partial x} \log P_t(x) = \frac{\partial}{\partial x} \frac{P_t(x)}{P_t(x)} = \frac{\partial}{\partial x} \frac{\int P_{t|0}(x|x_0) P_0(x_0) dx_0}{P_t(x)} =$$

$$= \int \left(\frac{\partial}{\partial x} P_{t|0}(x|x_0) \right) \frac{P_0(x_0)}{P_t(x)} dx_0 =$$

$\nwarrow P_{0|t}(x_0|x)$

$$= \int \left(\frac{\partial}{\partial x} \log P_{t|0}(x|x_0) \right) \boxed{\frac{P_0(x_0) P_t(x|x_0)}{P_t(x)}} dx_0$$

Score function

$$\frac{\partial}{\partial x} \log P_t(x) = E \left[\frac{\partial}{\partial x_t} \log P_{t|0}(x_t | X_0) \mid X_t = x \right]$$

easy ↓

Interpretation:

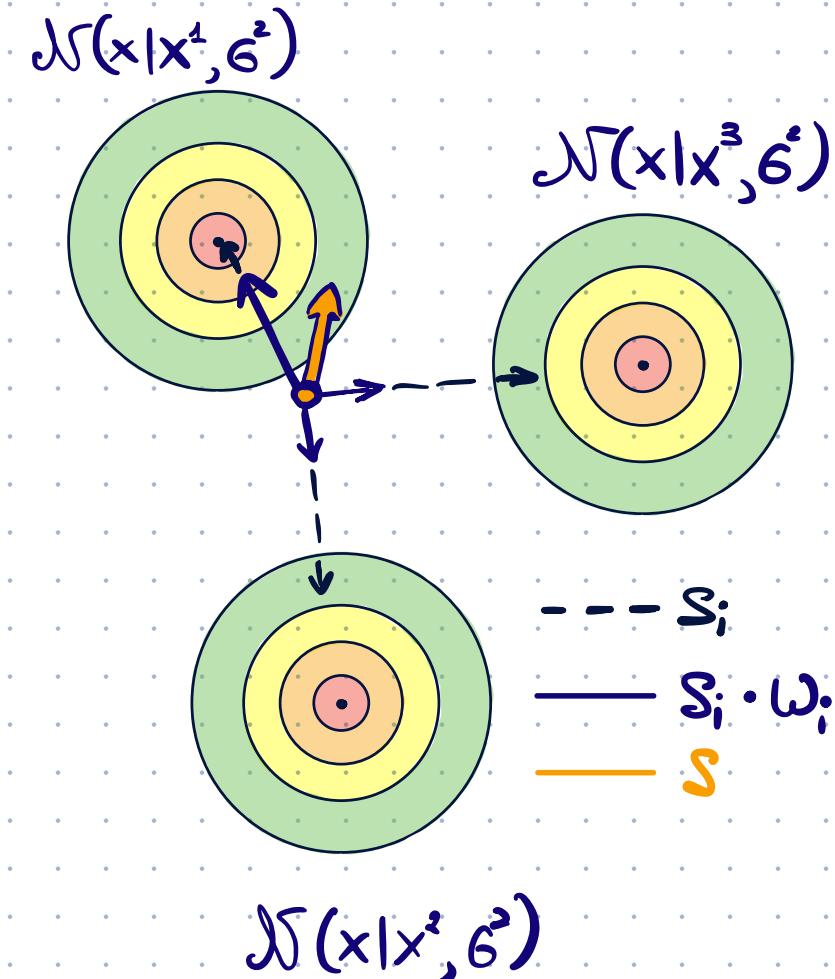
$$\text{Let } P_0(x) = \sum_{i=1}^m \omega_i I\{x = x^i\}$$

$$P_t(x) = \sum_{i=1}^m \omega_i \mathcal{N}(x | x^i, \sigma_t^2 I)$$

$$\frac{\partial}{\partial x} \log P_t(x) =$$

S

$$= \sum_{i=1}^m \underbrace{\frac{\partial}{\partial x} \log P_t(x | x^i)}_{S_i} \cdot \underbrace{P_{0|t}(x^i | x)}_{\omega_i}$$



Learning score function

$$\frac{\partial}{\partial x} \log P_t(x) = E \left[\frac{\partial}{\partial x_t} \log P_{t|0}(X_t | X_0) \mid X_t = x \right]$$

Thm $g^*(x) = E[Y | X=x] = \underset{g}{\operatorname{arg\,min}} E \|g(x) - Y\|^2$

In particular,

$$E \|g(x) - E[Y | X]\|^2 = E \|g(x) - Y\|^2 + \text{const}$$

For score:

$$E \|S_\theta(X_t, t) - \frac{\partial}{\partial x} \log P_t(x_t)\|^2 \rightarrow \min_\theta \Leftrightarrow$$

$$\Leftrightarrow E \|S_\theta(X_t, t) - \frac{\partial}{\partial x} \log P_{t|0}(X_t | X_0)\|^2 \rightarrow \min_\theta$$

L2 projection proof

$$\mathbb{E} \|g(x) - \mathbb{E}[Y|X]\|^2 = \mathbb{E} \|g(x)\|^2 - 2\mathbb{E} \langle g(x), \mathbb{E}[Y|X] \rangle + \text{const}$$

$$\begin{aligned}\mathbb{E} \langle g(x), \mathbb{E}[Y|X] \rangle &= \int \left\langle g(x), \int y \cdot P_{Y|X}(y|x) dy \right\rangle P_X(x) dx \\ &= \iint \langle g(x), y \rangle P_{Y|X}(y|x) P_X(x) dx dy = \mathbb{E} \langle g(x), Y \rangle.\end{aligned}$$

$$\begin{aligned}\mathbb{E} \|g(x) - \mathbb{E}[Y|X]\|^2 &= \mathbb{E} \|g(x)\|^2 - 2\mathbb{E} \langle g(x), Y \rangle + \text{const} \\ &= \mathbb{E} \|g(x) - Y\|^2 + \text{const}\end{aligned}$$



Denoising Score Matching

Loss: $X_0 \sim P_{\text{data}}, \varepsilon \sim \mathcal{N}(0, I), X_t = X_0 + \sigma_t \cdot \varepsilon$

$$\int_0^1 \mathbb{E} \| S_\theta(X_t, t) - \frac{\partial}{\partial x} \log P_t(X_t | X_0) \|^2 dt \rightarrow \min_\theta$$

- $P_{t|0}(X_t | X_0) = \mathcal{N}(X_t | X_0, \sigma_t^2 I) = C \cdot \exp\left(-\frac{1}{2\sigma_t^2} \|X_t - X_0\|^2\right)$

$$\frac{\partial}{\partial x_t} \log P_{t|0}(X_t | X_0) = -\frac{1}{\sigma_t^2} (X_t - X_0) = \frac{X_0}{\sigma_t^2} - \frac{X_t}{\sigma_t^2}$$

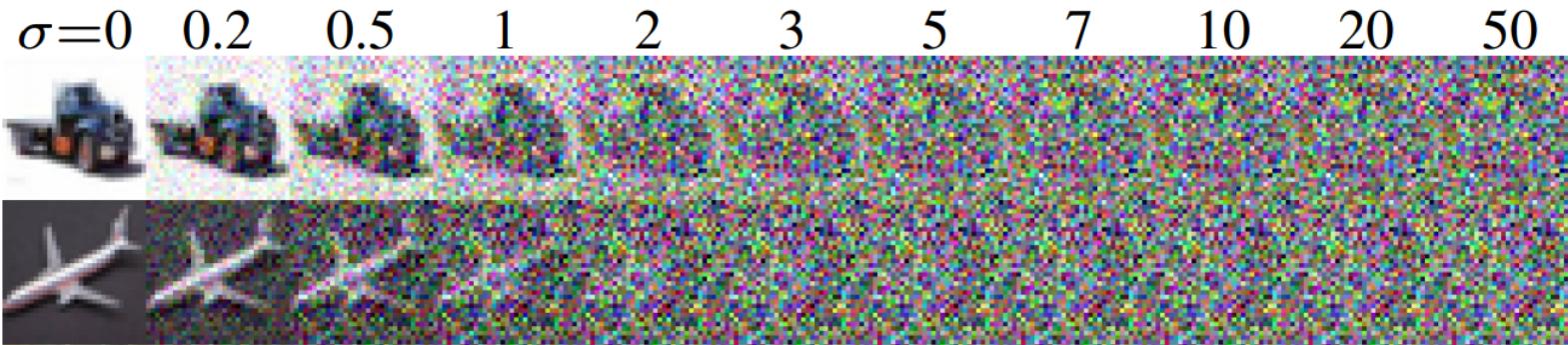
- Parameterize $S_\theta(x_t, t) = \frac{D_\theta(x_t, t)}{\sigma_t^2} - \frac{x_t}{\sigma_t^2}$

Loss: $\int_0^1 \omega_t \mathbb{E} \| D_\theta(x_t, t) - X_0 \|^2 dt \rightarrow \min_\theta$

\nwarrow denoiser

Ideal denoiser

Theoretical optimum: $D^*(X_t, t) = \mathbb{E}[X_0 | X_t]$,
i.e. for $t=1$ $D^*(X_1, 1) = \mathbb{E}[X_0 | X_1] \approx \mathbb{E}X_0$.



$$X_6 = X_0 + 6\epsilon$$



$$D^*(X_6, 6)$$

Sampling

$$dX_t^{(b)} = g^2(t-t) S_\theta(X_t^{(b)}, t-t) dt + g(t-t) d\omega_t$$

$$X_T \sim \mathcal{N}(0, \sigma_T^2)$$

$$\begin{aligned} X_{t+h} &\leftarrow X_t + h g^2(t) S_\theta(X_t, t) + \sqrt{h} g(t) \varepsilon_t \\ &= X_t + h g^2(t) \frac{D_\theta(X_t, t) - X_t}{\sigma_t^2} + \sqrt{h} g(t) \varepsilon_t \\ &= X_t \left(1 - \frac{h g^2(t)}{\sigma_t^2} \right) + D_\theta(X_t, t) \cdot \frac{h g^2(t)}{\sigma_t^2} + \sqrt{h} g(t) \varepsilon_t \end{aligned}$$

Sampling

Remember
this guy?

Sampling:

- Claim: $P_{t-1|t,0}(X_{t-1} | X_t, \hat{X}_0) = \mathcal{N}(X_{t-1} | \mu_t, \Sigma_t)$
 - $\mu_t = X_t \cdot \left(1 - \frac{g^2(t)}{G_t^2}\right) + \hat{X}_0 \cdot \frac{g^2(t)}{G_t^2}, \Sigma_t = \boxed{\frac{G_{t-1}}{G_t^2}} g^2(t) \cdot I$
 - Sample $X_T \sim \mathcal{N}(0, G_T^2)$
 - Iterate
- $$X_{t-1} = X_t \left(1 - \frac{g^2(t)}{G_t^2}\right) + D_\theta(X_t, t) \cdot \frac{g^2(t)}{G_t^2} + g(t) \varepsilon_t$$

This is him now:

$$X_{t-h} =$$

$$= X_t \left(1 - \frac{h g^2(t)}{G_t^2}\right) + D_\theta(X_t, t) \cdot \frac{h g^2(t)}{G_t^2} + \sqrt{h} g(t) \varepsilon_t$$

ODE sampling

$$dY_t^{(b)} = \frac{g^2(t-t)}{2} S_\theta(Y_t^{(b)}, t-t) dt$$

$$Y_{t-h} \leftarrow Y_t + h \cdot \frac{g^2(t)}{2} S_\theta(Y_t, t) =$$

$$= Y_t + h \frac{g^2(t)}{2} \frac{D_\theta(Y_t, t) - Y_t}{G_t^2} =$$

$$= Y_t \left(1 - h \frac{g^2(t)}{2 G_t^2} \right) + h \frac{g^2(t)}{2 G_t^2} D_\theta(Y_t, t)$$

Example: $G_t^2 = t^2$ $g(t) = \sqrt{2t}$

$$Y_{t-h} \leftarrow Y_t \left(1 - \frac{h}{t} \right) + \frac{h}{t} \cdot D_\theta(Y_t, t)$$

Summary

- $\frac{\partial}{\partial x} \log P_t(x) = \bar{E}\left[\frac{\partial}{\partial x} \log P_{t|0}(X_t|X_0) \mid X_t=x\right]$
- $\int_0^t \omega_t \bar{E} \left\| S_\theta(X_t, t) - \frac{\partial}{\partial x} \log P_{t|0}(X_t|X_0) \right\|^2 dt \rightarrow \min_\theta$
 \Leftrightarrow
 $\int_0^t \omega_t \bar{E} \left\| D_\theta(X_t, t) - X_0 \right\|^2 dt \rightarrow \min_\theta$
- SDE sampling $\Leftrightarrow \dot{X}_0 = D(X_t, t); X_{t-1} \sim P_{t-1|t,0}$
- ODE sampling : simple, interpretable

Links

<https://arxiv.org/abs/1907.05600>

(score matching)

<https://arxiv.org/abs/2011.13456>

(SDE diffusion models)

<https://arxiv.org/abs/2206.00364>

(Careful sampling design)

