

Next Week Guest Lectures on Zoom

Score-Based Models and Diffusion Models



Recap: Boltzmann Machine Training

- Objective: maximum likelihood learning (assume $T=1$):
 - Probability of one sample:

$$P(y) = \frac{\exp(\frac{1}{2}y^T W y)}{\underbrace{\sum_{y'} \exp(\frac{1}{2}y'^T W y')}_{Z_\theta}}$$

\mathcal{R}^d

- Maximum log-likelihood:

$$L(W) = \frac{1}{N} \sum_{y \in D} \frac{1}{2} y^T W y - \log \sum_{y'} \exp(\frac{1}{2} y'^T W y')$$

Can we avoid calculating the gradient of normalizing constant ($\nabla_x Z_\theta$)?

Z_θ : normalizing constant

Score Matching

- Score Function
 - Definition:

$$\nabla_x \log p_{data}(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

- Idea: directly fitting the score function:

$$\min_{\theta} \mathbb{E}_{p_{data}} \|\nabla_x \log p_{\theta}(x) - \nabla_x \log p_{data}(x)\|^2$$

- No need to compute $\nabla_x Z_{\theta}$!

- Problem:

- How to compute $\nabla_x \log p_{data}(x)$?

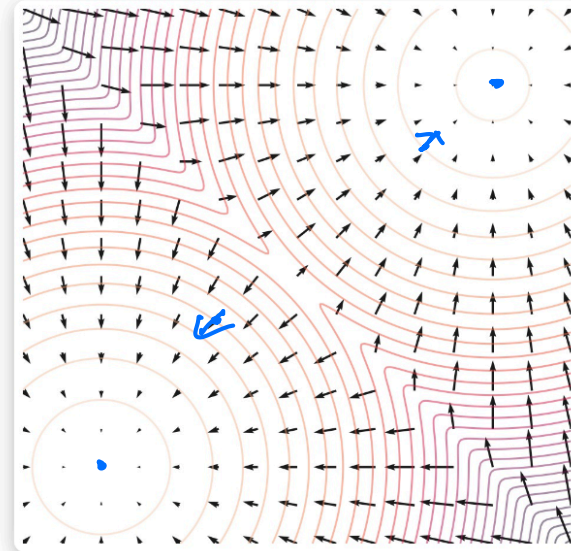
$$\{x_1, \dots, x_n\}$$

model $p_{\theta}(x)$
 \downarrow
 model $\nabla_x \log p_{data}(x)$

$$p(x) = \frac{e^{-f_{\theta}(x)}}{\sum_{x'} e^{-f_{\theta}(x')}}$$

$$\log p(x) = -f_{\theta}(x) - \underbrace{\sum_{x'} f_{\theta}(x')}_{\text{constant}}$$

$$\nabla \log p(x) = -\nabla f_{\theta}(x)$$



Score function (the vector field) and density function (contours) of a mixture of two Gaussians.

Score Matching

$$\begin{aligned} & \mathbb{E}_{P_{\text{data}}} \left\| \nabla_x \log P_{\theta}(x) - \nabla_x \log P_{\text{data}}(x) \right\|^2 \\ &= \mathbb{E}_{P_{\text{data}}} \left\| \nabla_x \log P_{\theta}(x) \right\|^2 + \mathbb{E}_{P_{\text{data}}} \left\| \nabla_x \log P_{\text{data}}(x) \right\|^2 \\ &\quad - 2 \mathbb{E}_{P_{\text{data}}} \left\langle \nabla_x \log P_{\theta}(x), \nabla_x \log P_{\text{data}}(x) \right\rangle \end{aligned}$$

(Integrating by parts)

$$\left(\mathbb{E}_P \left\langle f(x), \nabla_x \log P(x) \right\rangle = - \mathbb{E}_P \left[\text{div} f(x) \right] \right)$$

where $\text{div} f(x) = \sum_i \frac{\partial f_i(x)}{\partial x_i}$

$$\Rightarrow \mathbb{E}_{P_{\text{data}}} \left\langle \nabla_x \log P_{\theta}(x), \nabla_x \log P_{\text{data}}(x) \right\rangle = - \mathbb{E} \left[\text{Tr} \left(\nabla_x^2 \log P_{\theta}(x) \right) \right]$$

$$\Rightarrow \text{loss} \Leftrightarrow \mathbb{E}_{P_{\text{data}}} \left\| \nabla_x \log P_{\theta}(x) \right\|_2^2 - 2 \mathbb{E}_{P_{\text{data}}} \left[\text{Tr} \left(\nabla_x^2 \log P_{\theta}(x) \right) \right]$$

Score Matching

use μ/σ to parameterize

$$\nabla_x \log P_{\theta}(x): S_{\theta}(x): \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$\text{Loss: } \frac{1}{N} \sum_{\text{training } x_i} \left(\| \underbrace{S_{\theta}(x_i)}_{\mathcal{O}(d)} \|_2^2 - 2 \left[\text{Tr} \left(\underbrace{D}_{\mathcal{O}(d)} S_{\theta}(x_i) \right) \right] \right)$$

Sliced Score Matching

A $0 \sim 0 (N \times 1)$

$$L(\theta) = \frac{1}{N} \sum_{x \in D} \underbrace{\|s_\theta(x)\|^2}_{\text{score}} - 2 [\text{Tr}(Ds_\theta(x))]$$

random projection

let $M \in \mathbb{R}^{d \times d}$, v random, $\mathbb{E}[vv^T] = I$, $\mathbb{E}_v[v^T M v] = \mathbb{E}[\text{Tr}(M v v^T)] = \text{Tr}(M)$

Sample v_1, \dots, v_k

$$v_1^T M v_1, \dots, v_k^T M v_k$$

$$k \ll d$$

Score Matching: Langevin Dynamics

x_0

$\epsilon \ll 1$
—

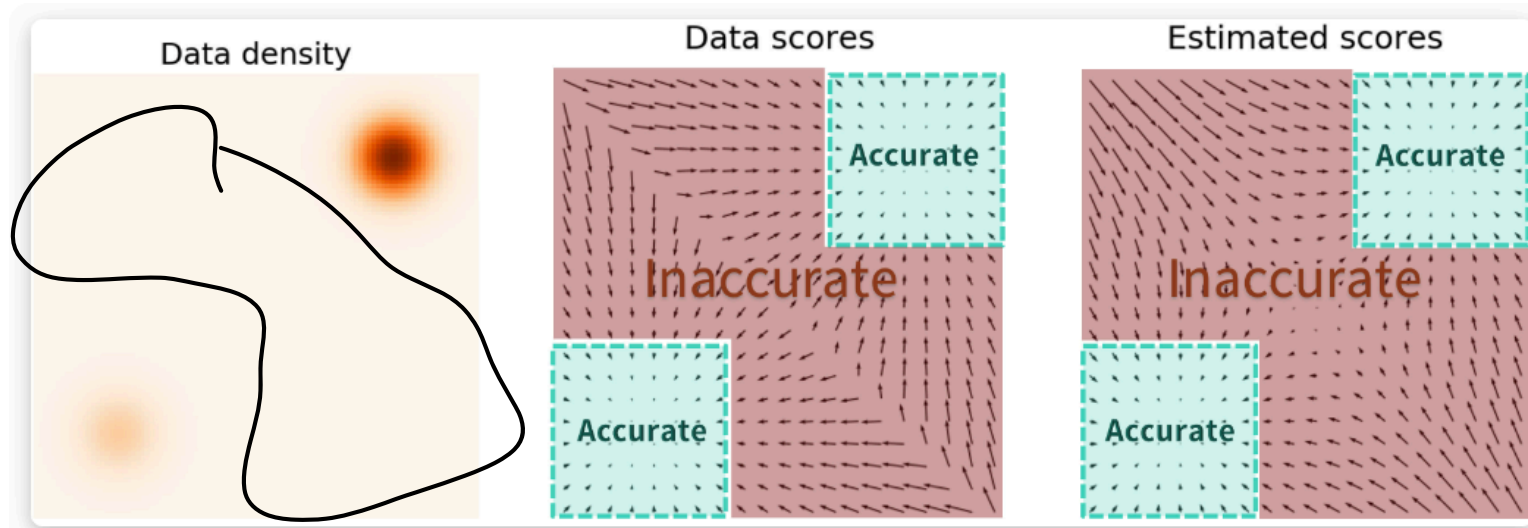
$$\underline{x_{t+1} \leftarrow x_t + \epsilon \nabla_x \log p(x) + \sqrt{2\epsilon} z_t, z_t \sim N(0, I)}$$

Stationary (equilibrium distribution): $p(x)$

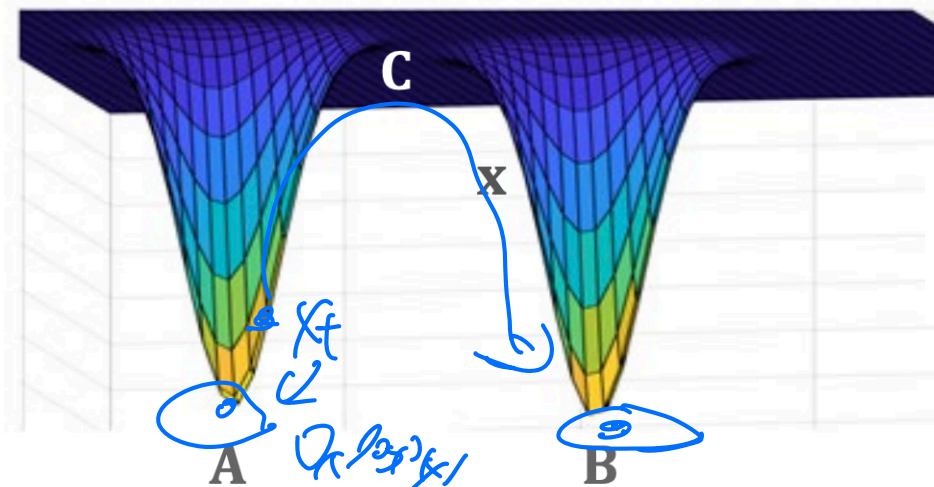
$\{x_1, \dots, x_\infty\} \rightarrow p(x)$

Practical Issues

- Score function estimation is inaccurate in low density regions (few data available).



- Sampling is Slow



Annealing: Denoising Score Matching

- Fit several “smoothed” versions of p_{data} :

- Choose temperatures: $\sigma_1, \sigma_2, \dots, \sigma_T$

- $p_{\sigma_i, data}(x) = p_{data}(x) * N(0, \sigma_i) = \int_{\delta} p_{data}(x - \delta) N(x; \delta, \sigma_i) d\delta$

- Implementation:

- Take a sample x , draw a sample $z \sim N(0, \sigma_i)$, output $x' = x + z$.

$$\sigma_1 > \sigma_2 > \dots > \sigma_{L-1} > \sigma_L$$

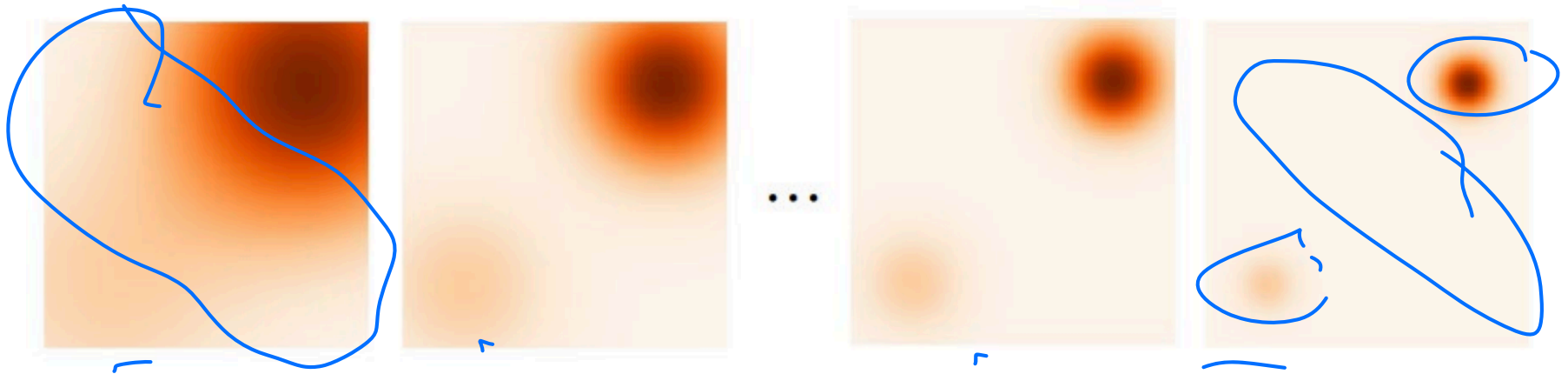


Figure by Stefano Ermon.

Annealing: Denoising Score Matching

$$\arg \min_{\theta} \sum_i \lambda(\sigma_i) \mathbb{E}_{x \sim p_{\sigma_i, data}} \|s_{\theta}(x, i) - \nabla_x \log p_{\sigma_i, data}(x)\|^2$$

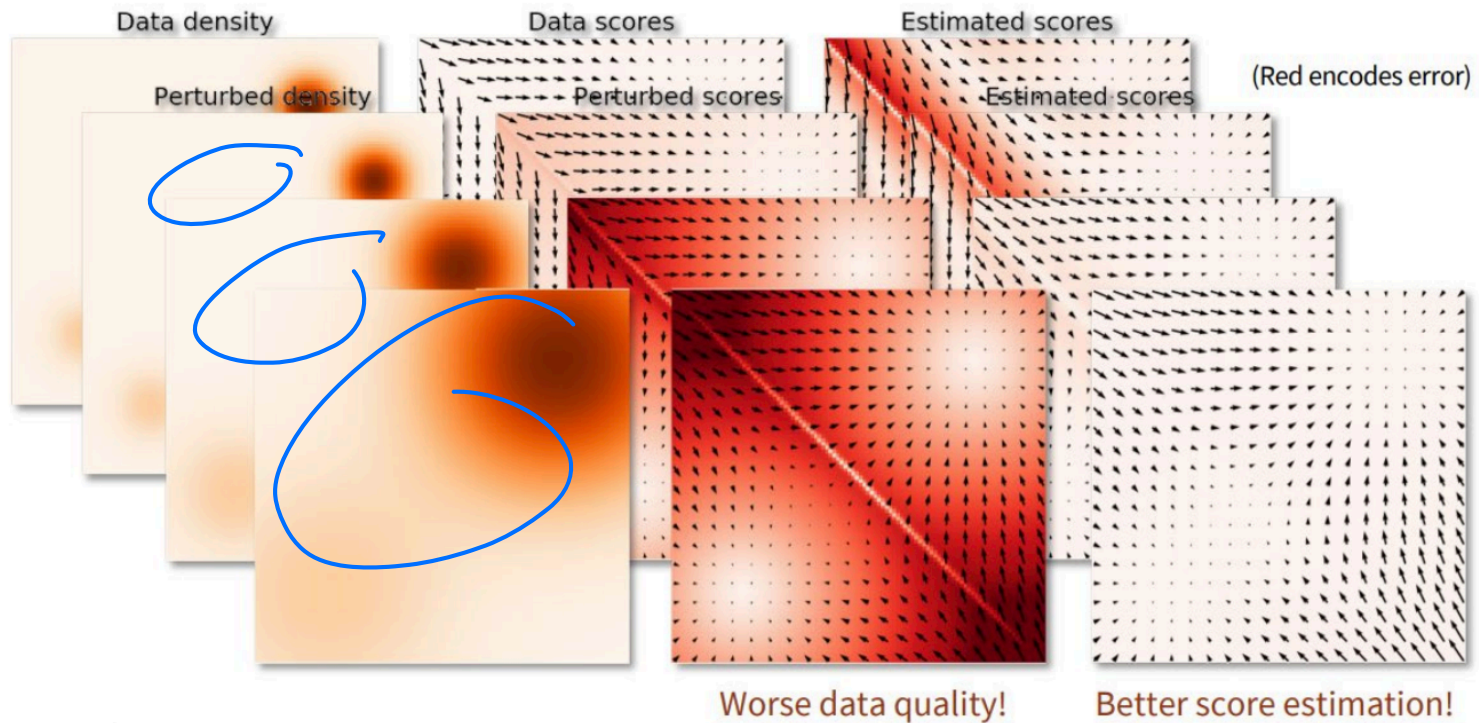


Figure by Stefano Ermon.

Annealed Langevin Dynamics

Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

1: Initialize $\tilde{\mathbf{x}}_0$

2: **for** $i \leftarrow 1$ to L **do**

3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ $\triangleright \alpha_i$ is the step size.

4: **for** $t \leftarrow 1$ to T **do**

5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$

6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$

7: **end for**

8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$

9: **end for**

return $\tilde{\mathbf{x}}_T$

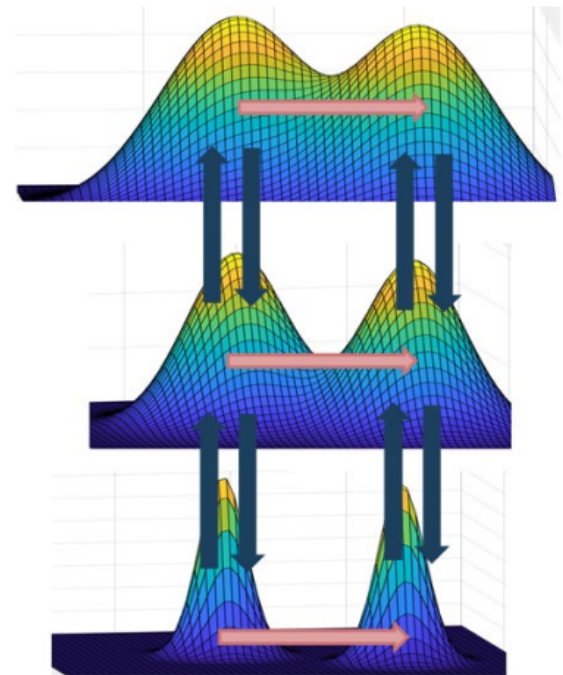


Figure from Song-Ermon '19

Diffusion Models



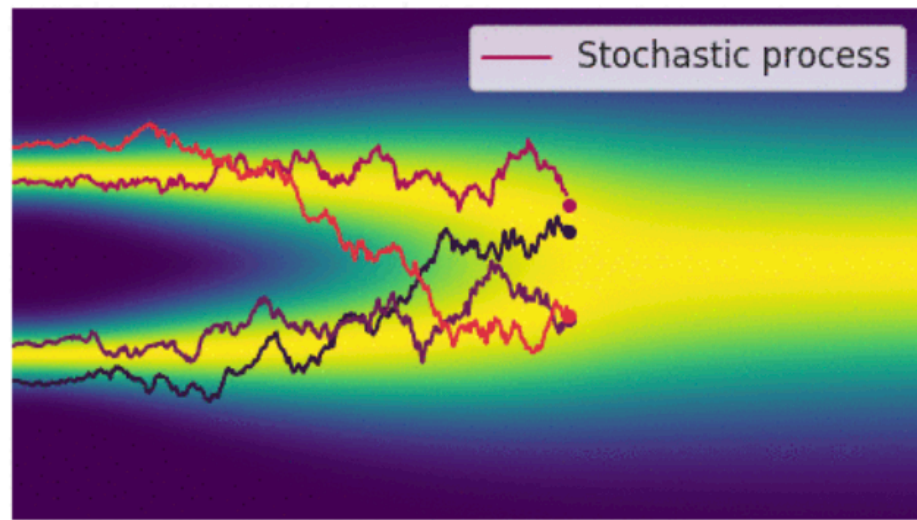
An image generated by Stable Diffusion based on the text prompt "a photograph of an astronaut riding a horse"

Perturbing Data with an SDE

$$\sigma_1 < \sigma_2 \dots \sigma_T$$

- Let the number of noise scales approaches infinity!

$$\sigma_1, \dots, \sigma_T \quad T \text{ large}, \quad X_T = x + \mathcal{N}(0, \sigma_T)$$
$$T \rightarrow \infty, \quad \{\sigma_t\}_t \quad /$$

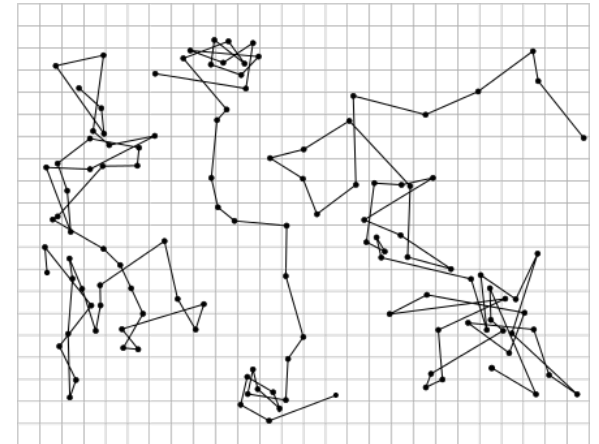


Perturbing data to noise with a continuous-time stochastic process.

Stochastic Differential Equations

$$dx = \underbrace{f(x, t)dt} + \underbrace{g(t)dw}$$

- $x(0)$: real image, $x(T)$: Gaussian noise.
- $f(x,t)$: drift terms. $g(t)$: diffusion coefficient.
- dw : Brownian motion
 - $w(t+u) - w(t) \sim N(0, u)$
- $f(x,t)$ and $g(t)$ are parts of the model.



- Variance Exploding SDE: $dx = \sqrt{\frac{d[\sigma^2(t)]}{dt}} dw.$
 - Variance Preserving SDE: $dx = -\frac{1}{2}\beta(t)xdt + \sqrt{\beta(t)}dw.$
 - $\sigma(t), \beta(t)$ are hyper-parameters.
- Exp [4]

Reversing the SDE

- Reversing the SDE: finding some stochastic process that goes from noise to data.
 - Use to generate data!
- Theorem (Anderson '82): there exists a reversing SDE, and it has a nice form:

$$dx = [f(x, t) - \underbrace{g^2(t) \nabla_x \log p_t(x)}]dt + g(t)dw$$

- Strategy: learn the score function, then solve this reverse SDE.

$$\begin{array}{c} X(\tau) \sim \mathcal{N}(0, \mathbb{I}) \\ \downarrow \\ X(0) \end{array}$$

Reversing the SDE

- Learning the score function: use score matching!

$$\arg \min_{\theta} \sum_i \lambda(\sigma_i) \mathbb{E}_{x \sim p_{\sigma_i, data}} \|s_{\theta}(x, i) - \nabla_x \log p_{\sigma_i, data}(x)\|^2$$

$$\Rightarrow \arg \min_{\theta} \mathbb{E}_{t \sim \text{unif}[0, T]} \mathbb{E}_{p_t(x)} [\lambda(t) \|s_{\theta}(x, t) - \nabla_x \log p_t(x)\|^2]$$

- Use existing techniques: sliced score matching
- No need to tune temperature schedule
 - Still need to choose a forward SDE, $\lambda(\sigma_i)$, etc
 - Typically choose $\lambda(t) \propto 1/\mathbb{E} \left[\|\nabla_{x(t)} \log p(x(t) | x(0))\|^2 \right]$

Sampling by Solving the Reverse SDE

$$X(\tau) \sim \mathcal{N}(\mu, \Sigma), \quad X(0)$$

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)]dt + g(t)dw$$

- Euler-Maruyama discretization:
 - $\Delta x \leftarrow [f(x, t) - g^2(t)s_\theta(x, t)]\Delta t + g(t)\sqrt{\Delta t}z_t$
 - $x \leftarrow x + \Delta x$
 - $t \leftarrow t + \Delta t$
- Other solvers:
 - Runge-Kutta
 - Predictor-corrector (Song et al. '21)

Evaluating Probability by Converting to ODE

- De-randomizing SDE

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)]dt + g(t)dw$$

$X(T)$



$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)]dt, x(T) \sim p_T$$

- Given an initial distribution and an ODE, we can evaluate probability at any time
 - Say given $x(T) \sim p_T$ and $dx = f(x, t)dt$

$$\log p_0(x(0)) = \log p_T(X(T)) + \int_0^T \text{Tr}(Df_\theta(x, t))dt$$

- Solve via ODE.