

Normalizing Flows



Intuition about easy to sample

- Goal: design $p(x)$ such that
 - Easy to sample
 - Tractable likelihood (density function)
- Easy to sample
 - Assume a continuous variable z
 - e.g., Gaussian $z \sim N(0,1)$, or uniform $z \sim \text{Unif}[0,1]$
 - $x = f(z)$, x is also easy to sample

Intuition about tractable density

- Goal: design $f(z; \theta)$ such that
 - Assume z is from an “easy” distribution
 - $p(x) = p(f(z; \theta))$ has tractable likelihood

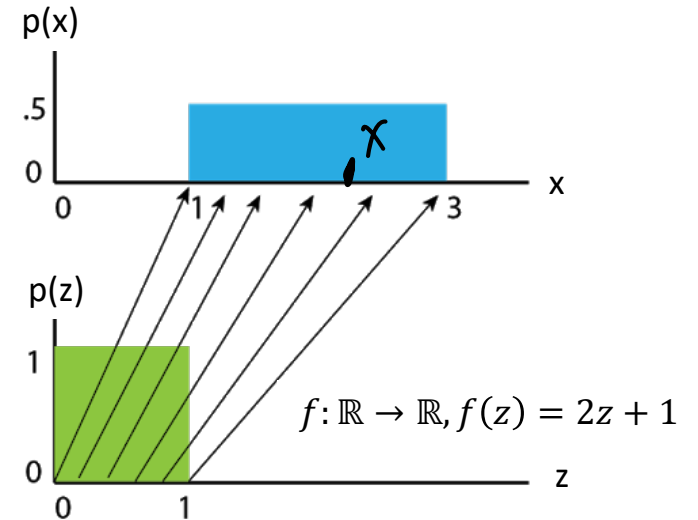
- Uniform: $z \sim \text{Unif}[0,1]$

- Density $p(z) = 1$

- $x = 2z + 1$, then $p(x) = ?$

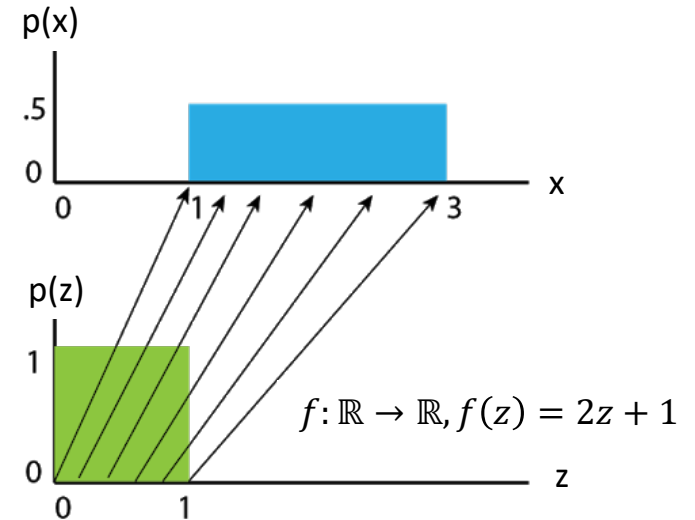
$$f(z) = 2z + 1$$

0.5



Intuition about tractable density

- Goal: design $f(z; \theta)$ such that
 - Assume z is from an “easy” distribution
 - $p(x) = p(f(z; \theta))$ has tractable likelihood
- Uniform: $z \sim \text{Unif}[0,1]$
 - Density $p(z) = 1$
 - $x = 2z + 1$, then $p(x) = 1/2$
 - $x = az + b$, then $p(x) = 1/|a|$ (for $a \neq 0$)
 - $x = f(z), p(x) = p(z) \left| \frac{dz}{dx} \right| = |f'(z)|^{-1} p(z)$
 - Assume $f(z)$ is a bijection



Change of variable

$$\underline{z \leftarrow f^{-1}(x)}$$

- Suppose $x = f(z)$ for some general non-linear $f(\cdot)$
 - The linearized change in volume is determined by the Jacobian of $f(\cdot)$:

$$\frac{\partial f(z)}{\partial z} = \begin{bmatrix} \frac{\partial f_1(z)}{\partial z_1} & \dots & \frac{\partial f_1(z)}{\partial z_d} \\ \dots & \dots & \dots \\ \frac{\partial f_d(z)}{\partial z_1} & \dots & \frac{\partial f_d(z)}{\partial z_d} \end{bmatrix} \quad \begin{matrix} z, x \in \mathbb{R}^d \\ \in \mathbb{R}^{d \times d} \end{matrix}$$

- Given a bijection $f(z) : \mathbb{R}^d \rightarrow \mathbb{R}^d$

- $\underline{z = f^{-1}(x)}$

- $\underline{p(x) = p(f^{-1}(x))} \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| = p(z) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$

- Since $\underline{\frac{\partial f^{-1}}{\partial x} = \left(\frac{\partial f}{\partial x} \right)^{-1}}$ (Jacobian of invertible function)

- $\underline{p(x) = p(z) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| = p(z) \left| \det \left(\frac{\partial f(z)}{\partial z} \right) \right|^{-1}}$

Normalizing Flow

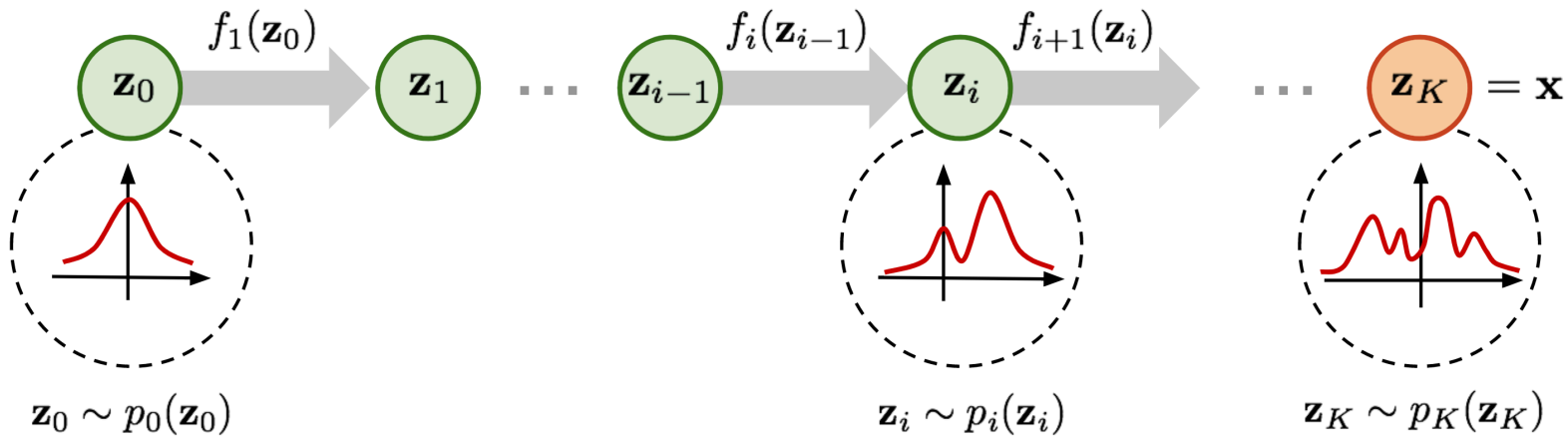
- Idea

- Sample z_0 from an “easy” distribution, e.g., standard Gaussian
- Apply K bijections $z_i = f_i(z_{i-1})$ f_1, \dots, f_K
- The final sample $x = f_K(z_K)$ has tractable density

- Normalizing Flow

- $z_0 \sim N(0, I)$, $z_i = f_i(z_{i-1})$, $x = z_K$ where $x, z_i \in \mathbb{R}^d$ and f_i is invertible
- Every reversible function produces a normalized density function

$$p(z_i) = p(z_{i-1}) \left| \det \left(\frac{\partial f_i}{\partial z_{i-1}} \right) \right|^{-1} \quad p(x) = \prod \left| \det(\cdot) \right|^{-1} \cdot p(z_0)$$



Normalizing Flow

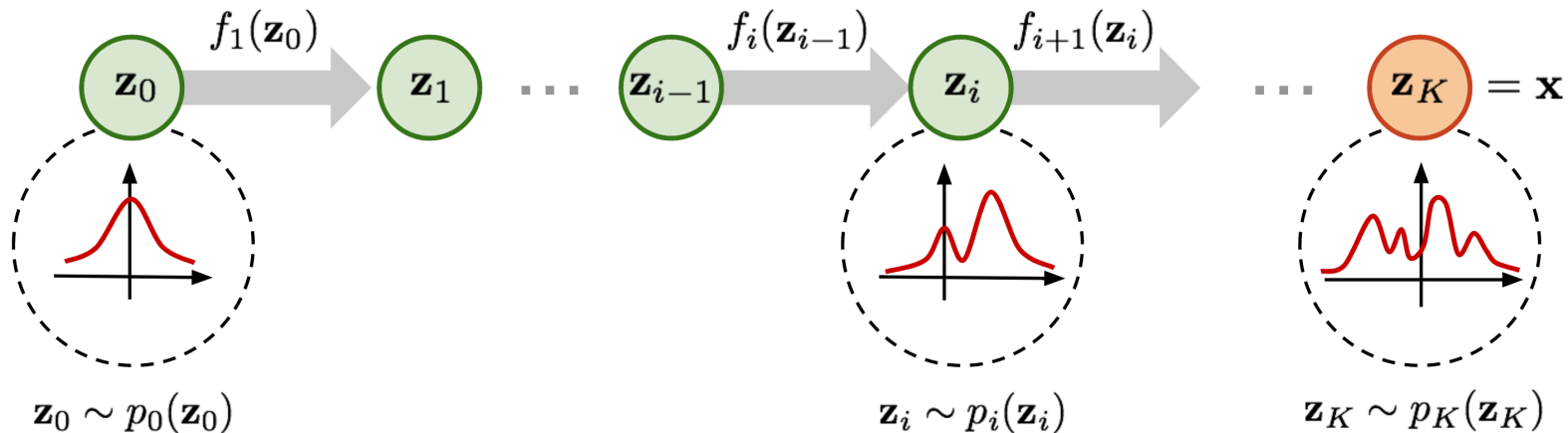
f_1, f_2, \dots, f_K

- Generation is trivial
 - Sample z_0 then apply the transformations
- Log-likelihood

$$\bullet \log p(x) = \log p(z_{K-1}) - \log \left| \det \left(\frac{\partial f_K}{\partial z_{K-1}} \right) \right|$$

$$\bullet \log p(x) = \log p(z_0) - \sum_i \log \left| \det \left(\frac{\partial f_i}{\partial z_{i-1}} \right) \right|$$

$O(d^3)!!!$



Normalizing Flow

- Naive flow model requires extremely expensive computation
 - Computing determinant of $d \times d$ matrices
- Idea:
 - Design a good bijection $f_i(z)$ such that the determinant is easy to compute

Plannar Flow

$A \in \mathbb{R}^{d \times d}$ $u, v \in \mathbb{R}^d$, $uv \in \mathbb{R}^{d \times d}$

- Technical tool: Matrix Determinant Lemma:

- $\det(A + uv^T) = (1 + v^T A^{-1} u) \det A$

- Model:

- $f_\theta(z) = z + u \odot h(w^T z + b)$

- $h(\cdot)$ chosen to be $\tanh(\cdot)$ ($0 < h'(\cdot) < 1$)

- $\theta = [u, w, b]$, $\det \left(\frac{\partial f}{\partial z} \right) = \det(I + h'(w^T z + b)uw^T) = \underbrace{1 + h'(w^T z + b)u^T w}$

- Computation in $O(d)$ time

- Remarks:

- $u^T w > -1$ to ensure invertibility

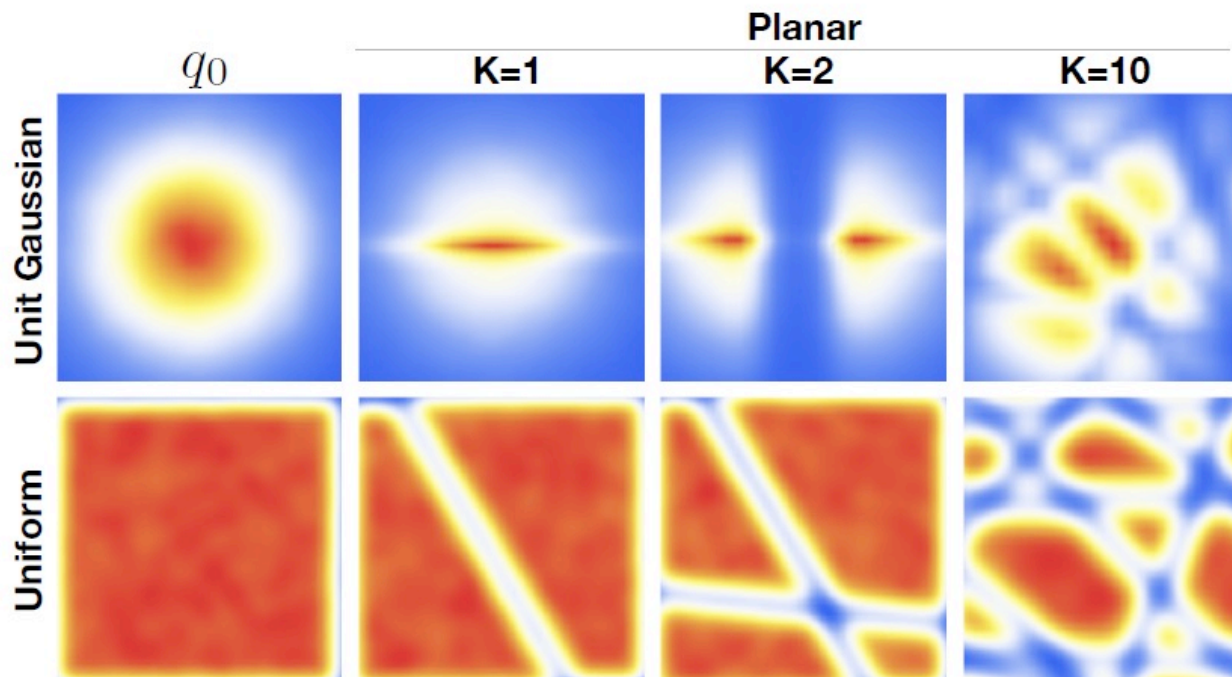
- Require normalization on u and w

\odot : pointwise product
 $A = I$

$O(d)$

Planar Flow (Rezende & Mohamed, '16)

- $f_{\theta}(z) = z + uh(w^{\top}z + b)$
- 10 planar transformations can transform simple distributions into a more complex one



Extensions

- Other flow models uses triangular Jacobian (NICE, Dinh et al. '14)
- Invertible 1x1 convolutions (Kingma et al. '18)
- Auto-regressive flow:
 - WaveNet (Deepmind '16)
 - PixelCNN (Deepmind '16)

Summary

- Pros:

- Easy to sample by transforming from a simple distribution
- Easy to evaluate the probability
- Easy training (MLE)

- Con

- Most restricted neural network structure
- Trade expressiveness for tractability

Score-Based Models and Diffusion Models



Recap: Boltzmann Machine Training

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

$$P(y) = \frac{\exp(\frac{1}{2}y^\top W y)}{\sum_{y'} \exp(\frac{1}{2}y'^\top W y')} \quad \neq$$

- Maximum log-likelihood:

$$L(W) = \frac{1}{N} \sum_{y \in D} \frac{1}{2} y^\top W y - \log \sum_{y'} \exp(\frac{1}{2}y'^\top W y')$$

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

Score Matching

- Score Function
 - Definition:

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

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

- Idea: directly fitting the score function:

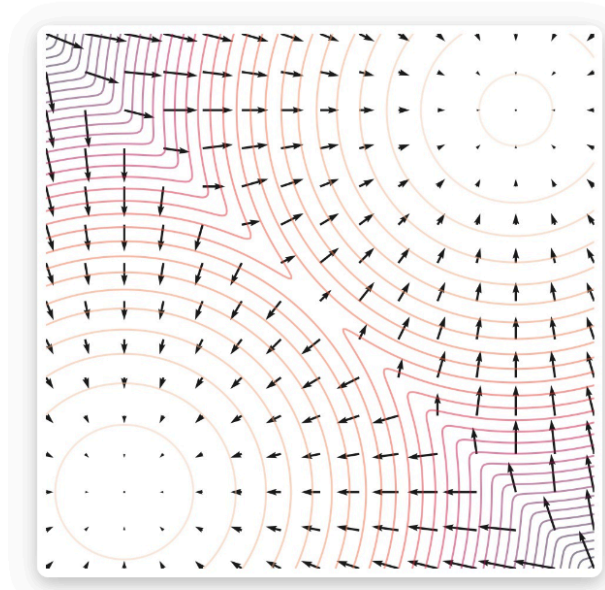
- $\min_{\theta} \mathbb{E}_{p_{data}} \|\nabla_x \log p_\theta(x) - \nabla_x \log p_{data}(x)\|^2$

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

- Problem:

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

$$\{x_1, \dots, x_N\}$$



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

Score Matching

$$\begin{aligned} & \mathbb{E}_{p_{\text{data}}} \left\| \nabla_x \log p_{\theta}(x) - \nabla_x \log p_{\text{data}}(x) \right\|_2^2 \\ &= \mathbb{E}_{p_{\text{data}}} \left\| \nabla_x \log p_{\theta}(x) \right\|_2^2 + \mathbb{E}_{p_{\text{data}}} \left\| \nabla_x \log p_{\text{data}}(x) \right\|_2^2 \\ &\quad - 2 \mathbb{E}_{p_{\text{data}}} \left\langle \nabla_x \log p_{\theta}(x), \nabla_x \log p_{\text{data}}(x) \right\rangle \end{aligned}$$

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

$$\Rightarrow \mathbb{E}_{p_{\text{data}}} \left\langle \nabla_x \log p_{\theta}(x), \nabla_x \log p_{\text{data}}(x) \right\rangle$$

$$= - \mathbb{E} \left[\text{Tr} \left(\nabla_x^2 \log p_{\theta}(x) \right) \right]$$

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

Score Matching

Use NN to parameterize

$\{x_1, \dots, x_N\}$

$S_\theta(x)$

$$\text{loss} : \frac{1}{N} \sum_{i=1}^N$$

$$\| \underbrace{S_\theta(x_i)}_{O(d)} \|_2^2 - 2 \left[\text{Tr} \left(\underbrace{I}_{O(d^2)} S_\theta(x_i) \right) \right]$$

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

Sliced Score Matching

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

Random projection

Let $M \in \mathbb{R}^{d \times d}$, v random, $\mathbb{E}[vv^T] = I$

$$\begin{aligned} \mathbb{E}_v [v^T M v] &= \mathbb{E} [\text{Tr}(M v v^T)] \\ &= \text{Tr}(M \cdot \mathbb{E}[v v^T]) \\ &= \text{Tr}(M) \end{aligned}$$

Sample v_1, \dots, v_K , for large K

$$\frac{1}{K} \sum_{i=1}^K v_i^T M v_i \rightarrow \text{Tr}(M)$$

Score Matching: Langevin Dynamics

x_0 arbitrary

$\epsilon \ll 1$

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

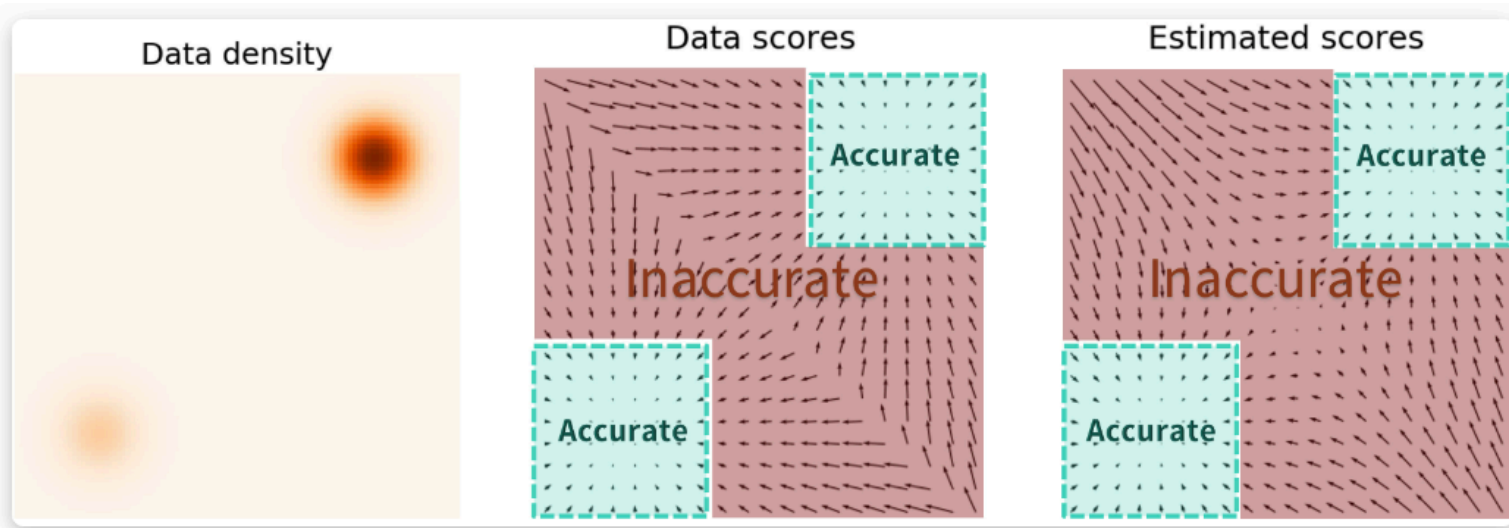
Stationary (equilibrium distribution): $p(x)$

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

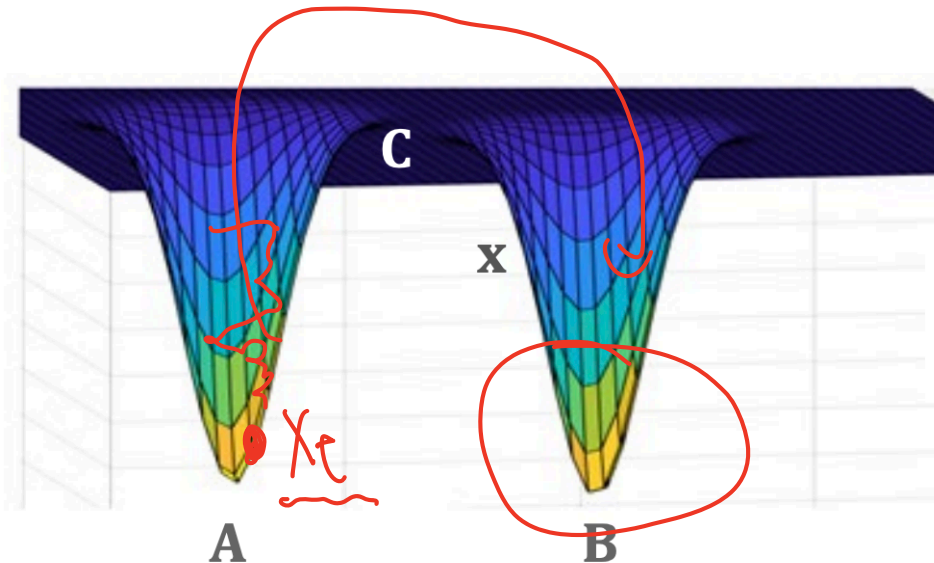
Practical Issues

$$\nabla_x \log p(x)$$

- 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 $x' = x + z$.

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

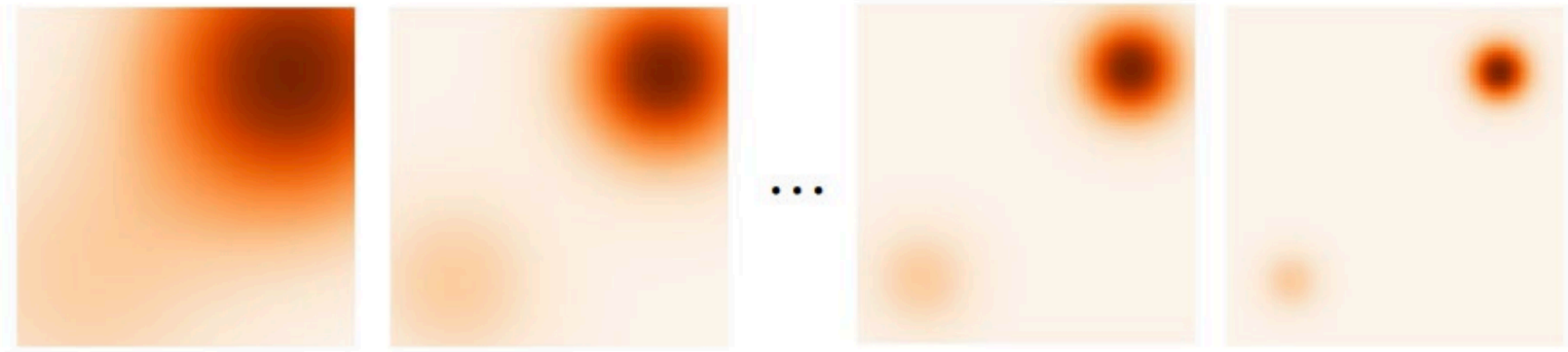


Figure by Stefano Ermon.

Annealing: Denoising Score Matching

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

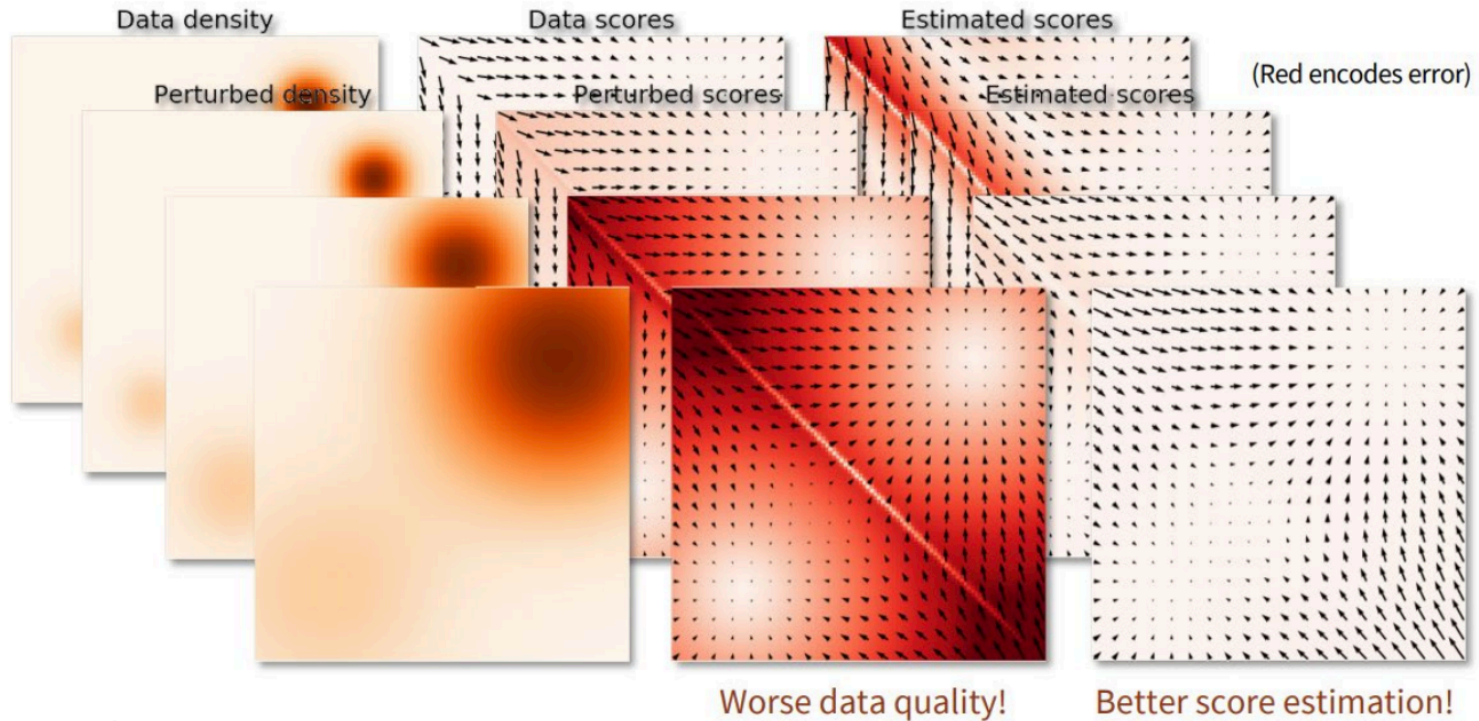


Figure by Stefano Ermon.

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

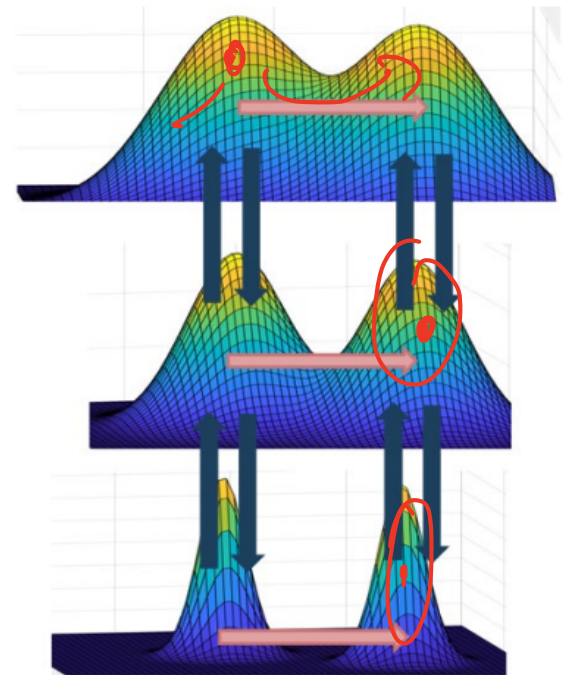


Figure from Song-Ermon '19

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

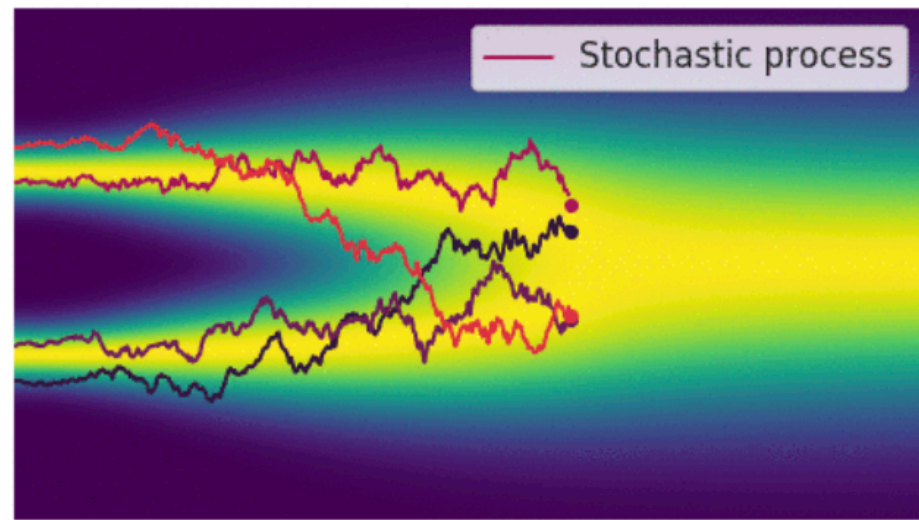
$\{\sigma_t\}$

$\sigma_1, \sigma_2, \dots, \sigma_T$

- Let the number of noise scales approaches infinity!

$Z \rightarrow f(Z)$

$X_t = X_{t-1} + \mathcal{N}(0, \sigma_t)$
 $T \rightarrow \infty \Rightarrow$ gaussian noise



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

Stochastic Differential Equations

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

- $x(0)$: real image, $x(T)$: Gaussian noise.
- $f(x,t)$: drift terms. $g(t)$: diffusion coefficient.

- dw : Brownian motion

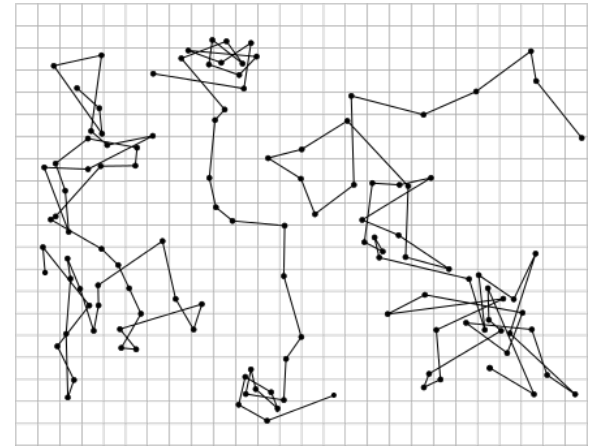
- $w(t+u) - w(t) \sim N(0, u)$

- $f(x,t)$ 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.



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:

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

$x(0)$

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

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 $\lambda(t) \propto 1/\mathbb{E} \left[\|\nabla_{x(t)} \log p(x(t) \mid x(0))\|^2 \right]$

Sampling by Solving the Reverse SDE

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

- 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

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

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

Gaussian

- 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.