Normalizing Flows

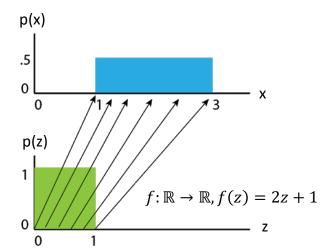


Intuition about easy to sample

- 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) = ?

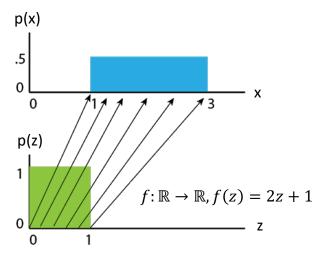


Intuition about tractable density

- 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 $a \ne 0$)

•
$$x = f(z), p(z) \left| \frac{dz}{dx} \right| = |f'(z)|^{-1} p(z)$$

• Assume f(z) is a bijection



Change of variable

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

$$\frac{\partial f(z)}{\partial z} = \begin{bmatrix} \frac{\partial f_z(x)}{\partial z_1} & \dots & \frac{\partial f_1(z)}{\partial z_d} \\ \dots & \dots & \dots \\ \frac{\partial f_d(z)}{\partial z_1} & \dots & \frac{\partial f_d(z)}{\partial z_d} \end{bmatrix}$$

- Given a bijection $f(z): \mathbb{R}^d \to \mathbb{R}^d$
 - $\bullet \ z = f^{-1}(x)$

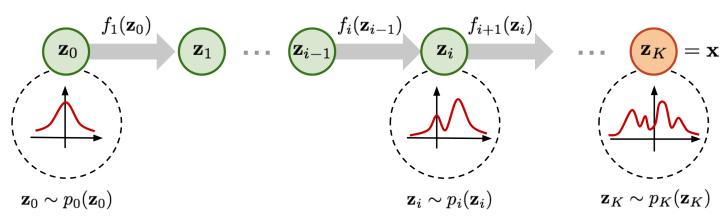
$$p(x) = p(f^{-1}(x)) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| = p(z) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

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

Normalizing Flow

- Idea
 - Sample z_0 from an "easy" distribution, e.g., standard Gaussian
 - Apply K bijections $z_i = f_i(z_{i-1})$
 - The final sample $x = f_K(z_K)$ has tractable desnity
- 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 revertible function produces a normalized density function

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



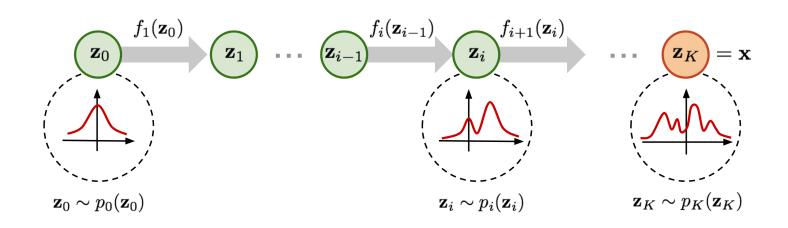
Normalizing Flow

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

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

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

$$O(d^3)!!!!$$



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

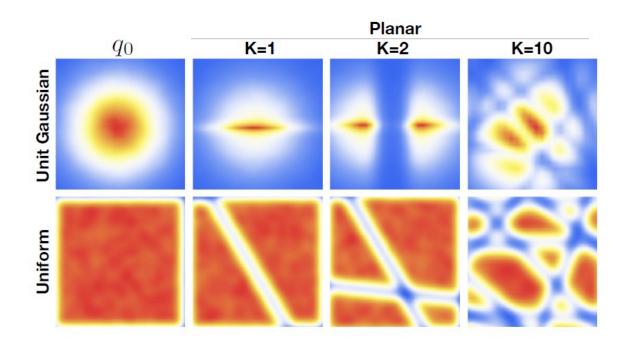
- Technical tool: Matrix Determinant Lemma:
 - $\det(A + uv^{\mathsf{T}}) + (1 + v^{\mathsf{T}}A^{-1}u) \det A$
- Model:
 - $f_{\theta}(z) + z + u \odot h(w^{\mathsf{T}}z + b)$
 - $h(\cdot)$ chosen to be $tanh(\cdot)(0 < h'(\cdot) < 1)$

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

- Computation in O(d) time
- Remarks:
 - $u^{\mathsf{T}}w > -1$ to ensure invertibility
 - Require normalization on u and w

Planar Flow (Rezende & Mohamed, '16)

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



Extensions

- Other flow models uses triangular Jacobian
 - Suppose $x_i = f_i(z)$ only depends on $z_{\leq i}$

Score-Based Modelsand Diffusion Models



Recap: Boltzmann Machine Training

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

$$P(y) = \frac{\exp(\frac{1}{2}y^{\mathsf{T}}Wy)}{\sum_{y'} \exp(y'^{\mathsf{T}}Wy')}$$

Maximum log-likelihood:

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

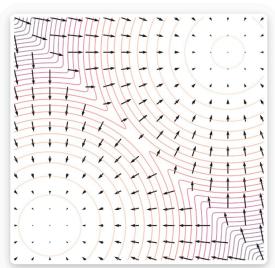
Can we avoid calculating the gradient of normalizing constant ($\nabla_x Z_\theta$)?

Score Matching

- Score Function
 - Definition:

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

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



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

Score Matching

Score Matching

Sliced Score Matching

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

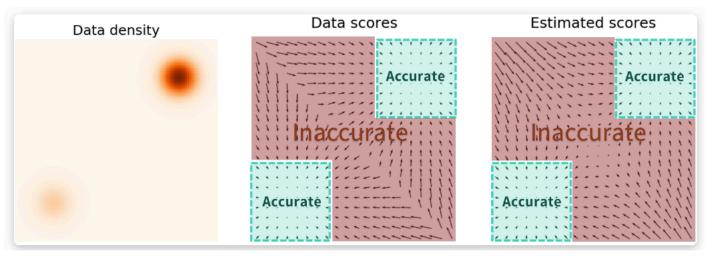
Score Matching: Langevin Dynamics

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

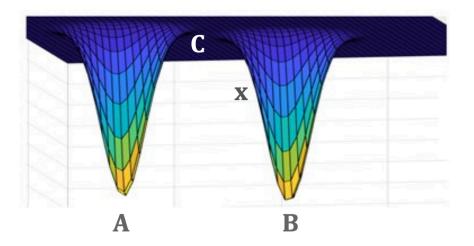
Stationary (equilibrium distribution): p(x)

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

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

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

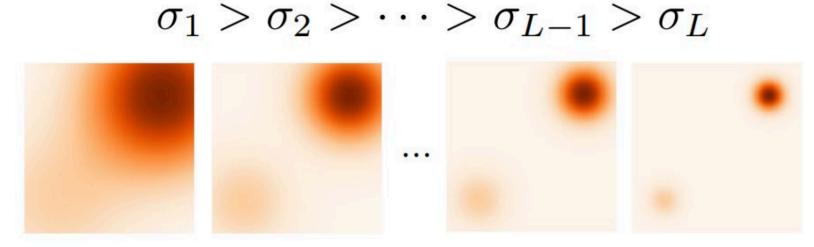


Figure by Stefano Ermon.

Annealing: Denoising Score Matching

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

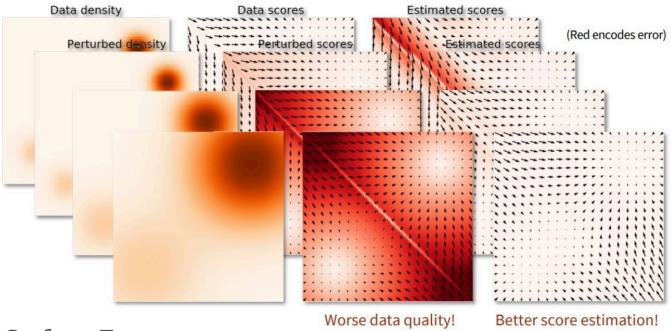


Figure by Stefano Ermon.

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
         \alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2 \Rightarrow \alpha_i is the step size.
  3:
          for t \leftarrow 1 to T do
  4:
                           Draw \mathbf{z}_t \sim \mathcal{N}(0, I)
  5:
                          \tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \ \mathbf{z}_t
  6:
                 end for
  7:
                 \tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T
  9: end for
         return \tilde{\mathbf{x}}_T
```

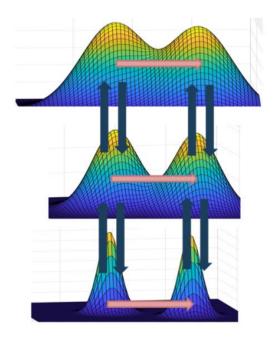


Figure from Song-Ermon '19

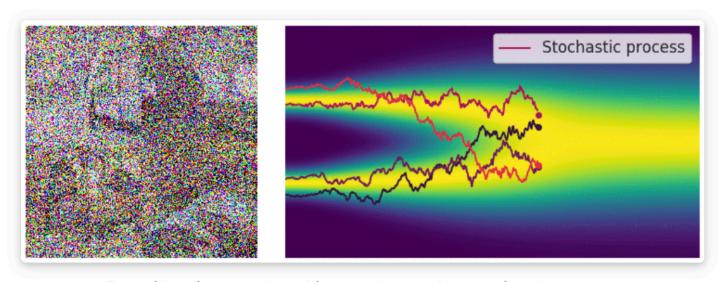
Diffusion Models



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

Perturbing Data with an SDE

• Let the number of noise scales approaches infinity!



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

Stochastic Differential Equations

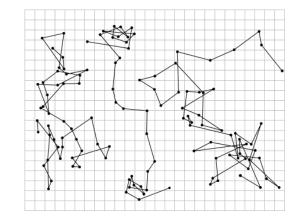
$$dx = f(x, t)dt + g(t)dw$$

- 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.

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:

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

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

Reversing the SDE

Learning the score function: use score matching!

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

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

- 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

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

$$dx = [f(x,t) - g^{2}(t) \nabla_{x} \log p_{t}(x)]dt, x(T) \sim p_{T}$$

- 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

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

Solve via ODE.