

Deep Generative Models

11. Energy-Based Models



- 국가수리과학연구소 산업수학혁신센터 김민중

Recap. of Energy-based model

- Energy-based models: $\frac{1}{Z(\theta)} \exp(f_\theta(x))$
 - $Z(\theta)$ is intractable, so no access to likelihood
 - Comparing the probability of two points is easy

$$\frac{p_\theta(x)}{p_\theta(x')} = \exp(f_\theta(x) - f_\theta(x'))$$

- Maximum likelihood training:

$$\max_{\theta} [f_\theta(x_{train}) - \log Z(\theta)]$$

- Contrastive divergence:

$$\nabla_{\theta} f_{\theta}(x_{train}) - \nabla_{\theta} \log Z(\theta) \approx \nabla_{\theta} f_{\theta}(x_{train}) - \nabla_{\theta} f_{\theta}(x_{sample})$$

- where $x_{sample} \sim p_{\theta}(x) = \frac{\exp(f_{\theta}(x))}{Z(\theta)}$

Sampling from EBMs: MH-MCMC

- Metropolis-Hastings Markov Chain Monte Carlo (MCMC):
 - Initialize \mathbf{x}^0 randomly, $t = 0$
 - $\mathbf{x}' = \mathbf{x}^t + noise$
 - $\mathbf{x}^{t+1} = \mathbf{x}'$ if $f_\theta(\mathbf{x}') \geq f_\theta(\mathbf{x}^t)$
 - If $f_\theta(\mathbf{x}') < f_\theta(\mathbf{x}^t)$, set $\mathbf{x}^{t+1} = \mathbf{x}'$ with probability $\exp(f_\theta(\mathbf{x}') - f_\theta(\mathbf{x}^t))$, otherwise set $\mathbf{x}^{t+1} = \mathbf{x}^t$
- In theory, when $T \rightarrow \infty$, then $\mathbf{x}^T \sim p_\theta$
- In practice, we need many iterations, and it converges slowly

Sampling from EBMs: unadjusted Langevin MCMC

- Unadjusted Langevin MCMC
 - Initialize $\mathbf{x}^0 \sim \pi(\mathbf{x})$
 - Repeat for $t = 0, \dots, T - 1$
 - $\mathbf{z} \sim N(0, I)$
 - $\mathbf{x}^{t+1} \sim \mathbf{x}^t + \epsilon \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})|_{\mathbf{x}=\mathbf{x}^t} + \sqrt{2\epsilon} \mathbf{z}$
- Properties:
 - If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, then we have \mathbf{x}^T converges to a sample from p_{θ}
 - $\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x})$ for continuous energy-based models
 - Convergence slows down as dimensionality grows

Sampling from EBMs: unadjusted Langevin MCMC

- Remark: **Sampling converges slowly in high dimensional spaces and is very expensive, yet we need sampling for each training iteration in contrastive divergence**

Today's topic

- Goal: Training method without sampling
 - Score Matching
 - Noise Contrastive Estimation

Score function

- Energy-based models

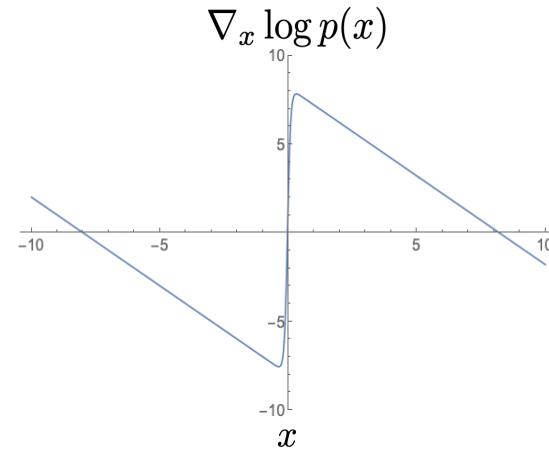
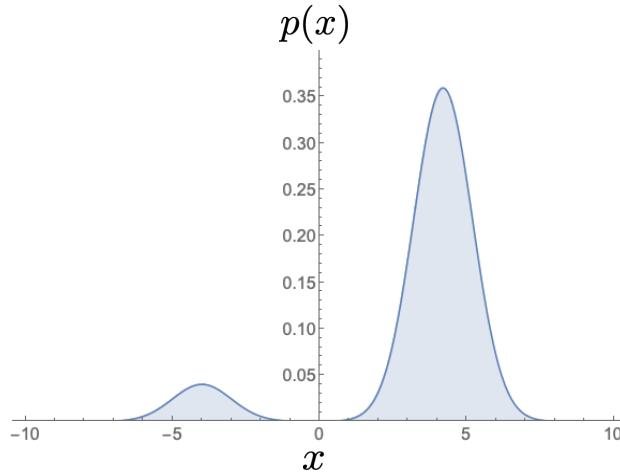
$$\frac{1}{Z(\theta)} \exp(f_\theta(x)), \quad \log p_\theta(x) = f_\theta(x) - \log Z(\theta)$$

- (Stein) Score function

$$s_\theta(x) := \nabla_x \log p_\theta(x) = \nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_{=0} = \nabla_x f_\theta(x)$$

Score matching

- Observation
 - $s_\theta(x)$ is independent of the partition function $Z(\theta)$



Score matching

- Fisher divergence between $p(\mathbf{x})$ and $q(\mathbf{x})$:

$$D_F(p, q) := \frac{1}{2} E_{\mathbf{x} \sim p} [\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x})\|_2^2]$$

- Score matching
 - minimizing the Fisher divergence between $p_{data}(\mathbf{x})$ and the EBM $p_\theta(\mathbf{x}) \propto \exp[f_\theta(\mathbf{x})]$

$$\begin{aligned} & \frac{1}{2} E_{\mathbf{x} \sim p_{data}} [\|\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x})\|_2^2] \\ &= \frac{1}{2} E_{\mathbf{x} \sim p_{data}} [\|\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x}) - \nabla_{\mathbf{x}} f_\theta(\mathbf{x})\|_2^2] \end{aligned}$$

Score matching

$$\frac{1}{2} E_{x \sim p_{data}} [\|\nabla_x \log p_{data}(x) - \nabla_x f_\theta(x)\|_2^2]$$

- How to deal with $\nabla_x \log p_{data}(x)$ given only samples?
- Integration by parts (univariate case)

$$\frac{1}{2} E_{x \sim p_{data}} \left[(\nabla_x \log p_{data}(x) - \nabla_x f_\theta(x))^2 \right]$$

$$= \frac{1}{2} \int p_{data}(x) \left[(\nabla_x \log p_{data}(x) - \nabla_x f_\theta(x))^2 \right] dx$$

$$= \frac{1}{2} \int p_{data}(x) (\nabla_x \log p_{data}(x))^2 dx + \frac{1}{2} \int p_{data}(x) (\nabla_x f_\theta(x))^2 dx$$

$$- \int p_{data}(x) \nabla_x \log p_{data}(x) \nabla_x f_\theta(x) dx$$

Score matching

$$\begin{aligned} & - \int p_{data}(x) \nabla_x \log p_{data}(x) \nabla_x f_\theta(x) dx \\ &= - \int p_{data}(x) \frac{1}{p_{data}(x)} \nabla_x p_{data}(x) \nabla_x f_\theta(x) dx \\ &= \underbrace{-p_{data}(x) \nabla_x f_\theta(x) \Big|_{x=-\infty}^{x=\infty}}_{=0} + \int p_{data}(x) \nabla_x^2 f_\theta(x) dx \\ &= \int p_{data}(x) \nabla_x^2 f_\theta(x) dx \end{aligned}$$

- Note that we need to assume p_{data} decays sufficiently rapidly, $p_{data}(x) \rightarrow 0$ when $x \rightarrow \pm\infty$

Score matching

- Univariate score matching

$$\begin{aligned} & \frac{1}{2} E_{x \sim p_{data}} \left[(\nabla_x \log p_{data}(x) - \nabla_x f_\theta(x))^2 \right] = \\ &= \underbrace{\frac{1}{2} \int p_{data}(x) (\nabla_x \log p_{data}(x))^2 dx + \frac{1}{2} \int p_{data}(x) (\nabla_x f_\theta(x))^2 dx}_{\text{const. wrt } \theta} \\ &+ \int p_{data}(x) \nabla_x^2 f_\theta(x) dx \\ &= E_{x \sim p_{data}} \left[\frac{1}{2} (\nabla_x f_\theta(x))^2 + \nabla_x^2 f_\theta(x) \right] + \text{const.} \end{aligned}$$

Score matching

- Multivariate score matching

$$\begin{aligned} & \frac{1}{2} E_{\mathbf{x} \sim p_{data}} [\|\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x}) - \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x})\|_2^2] \\ &= E_{\mathbf{x} \sim p_{data}} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} f_{\theta}(\mathbf{x})\|_2^2 + \text{tr} \left(\underbrace{\nabla_{\mathbf{x}}^2 f_{\theta}(\mathbf{x})}_{= \text{Hessian of } f_{\theta}(\mathbf{x})} \right) \right] + \text{const.} \end{aligned}$$

Score matching

- Sample a mini-batch of datapoints $\{x_1, x_2, \dots, x_n\} \sim p_{data}(x)$
- Estimate the score matching loss with the empirical mean

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \|\nabla_x \log p_\theta(x_i)\|_2^2 + \text{tr}(\nabla_x^2 \log p_\theta(x_i)) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \|\nabla_x f_\theta(x_i)\|_2^2 + \text{tr}(\nabla_x^2 f_\theta(x_i)) \right] \end{aligned}$$

- Stochastic gradient descent
- No need to sample from the EBM
- **Remark:** Computing the trace of Hessian $\text{tr}(\nabla_x^2 f_\theta(x))$ is in general very expensive for large models
- Denoising score matching (Vincent 2010) and sliced score matching (Song et al. 2019)

Recap

- Distances used for training energy-based models

- KL divergence = maximum likelihood

$$\nabla_{\theta} [f_{\theta}(x_{data}) - f_{\theta}(x_{sample})] \quad (\text{contrastive divergence})$$

- Fisher divergence = score matching

$$\min_{\theta} E_{x \sim p_{data}} [\|\nabla_x \log p_{data}(x) - \nabla_x f_{\theta}(x)\|_2^2]$$

$$\Leftrightarrow \min_{\theta} E_{x \sim p_{data}} \left[\frac{1}{2} \|\nabla_x f_{\theta}(x)\|_2^2 + \text{tr}(\nabla_x^2 f_{\theta}(x)) \right]$$

Noise contrastive estimation

- Learning an energy-based model by contrasting it with a noise distribution
 - Data distribution: $p_{data}(\mathbf{x})$
 - Noise distribution: $p_n(\mathbf{x})$ analytically tractable and easy to sample from
 - Training a discriminator(neural network) $D_\theta(\mathbf{x}) \in [0,1]$ to distinguish between data samples and noise samples
$$\max_\theta E_{\mathbf{x} \sim p_{data}} [\log D_\theta(\mathbf{x})] + E_{\mathbf{x} \sim p_n} [\log(1 - D_\theta(\mathbf{x}))]$$
- What is the optimal discriminator $D_{\theta^*}(\mathbf{x})$?

$$D_{\theta^*}(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_n(\mathbf{x})}$$

Noise contrastive estimation

- Let the discriminator be parametrized by $D_\theta(\mathbf{x}) = \frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_n(\mathbf{x})}$
- Then the optimal discriminator $D_{\theta^*}(\mathbf{x})$ satisfies

$$D_{\theta^*}(\mathbf{x}) = \frac{p_{\theta^*}(\mathbf{x})}{p_{\theta^*}(\mathbf{x}) + p_n(\mathbf{x})} = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_n(\mathbf{x})}$$

- By training the discriminator, we are implicitly learning $p_{\theta^*} \approx p_{data}$. Particularly suitable for cases where $p_\theta(\mathbf{x})$ is defined up to a normalization constant (EBMs)
- Equivalently,

$$p_{\theta^*}(\mathbf{x}) = \frac{p_n(\mathbf{x})D_{\theta^*}(\mathbf{x})}{1 - D_{\theta^*}(\mathbf{x})} = p_{data}(\mathbf{x})$$

- Classifier is used to correct density estimates from p_n . It can be used to improve a base generative model (Boosted Generative Models, Grover et al., 2018)

Noise contrastive estimation for training EBMs

- Energy-based models: $\frac{1}{Z(\theta)} \exp(f_\theta(x))$
- Modeling $Z(\theta)$ with an additional trainable parameter Z that is not explicitly constrained to satisfy $Z = \int e^{f_\theta(x)} dx$
- Let

$$p_{\theta,Z}(x) = \frac{\exp(f_\theta(x))}{Z}$$

- With noise contrastive estimation, the optimal parameters θ^*, Z^* are

$$p_{\theta^*,Z^*}(x) = \frac{\exp(f_{\theta^*}(x))}{Z^*} = p_{data}(x)$$

- The optimal parameter Z^*

$$\int \frac{\exp(f_{\theta^*}(x))}{Z^*} dx = \int p_{data}(x) dx = 1 \Rightarrow Z^* = \int \exp(f_{\theta^*}(x)) dx$$

Noise contrastive estimation for training EBMs

- The discriminator $D_{\theta,Z}(\mathbf{x})$ for probabilistic model $p_{\theta,Z}(\mathbf{x})$ is

$$D_{\theta,Z}(\mathbf{x}) = \frac{\frac{e^{f_\theta(\mathbf{x})}}{Z}}{\frac{e^{f_\theta(\mathbf{x})}}{Z} + p_n(\mathbf{x})} = \frac{e^{f_\theta(\mathbf{x})}}{e^{f_\theta(\mathbf{x})} + Zp_n(\mathbf{x})}$$

- Noise contrastive estimation training

$$\max_{\theta,Z} E_{\mathbf{x} \sim p_{data}} [\log D_{\theta,Z}(\mathbf{x})] + E_{\mathbf{x} \sim p_n} [\log (1 - D_{\theta,Z}(\mathbf{x}))]$$

- Equivalently,

$$\max_{\theta,Z} E_{\mathbf{x} \sim p_{data}} [f_\theta(\mathbf{x}) - \log (e^{f_\theta(\mathbf{x})} + Zp_n(\mathbf{x}))] + E_{\mathbf{x} \sim p_n} [\log (Zp_n(\mathbf{x})) - \log (e^{f_\theta(\mathbf{x})} + Zp_n(\mathbf{x}))]$$

- Log-sum-exp trick for numerical stability:

$$\log (e^{f_\theta(\mathbf{x})} + Zp_n(\mathbf{x})) = \text{logsumexp}(f_\theta(\mathbf{x}), \log Z + \log p_n(\mathbf{x}))$$

Noise contrastive estimation for training EBMs

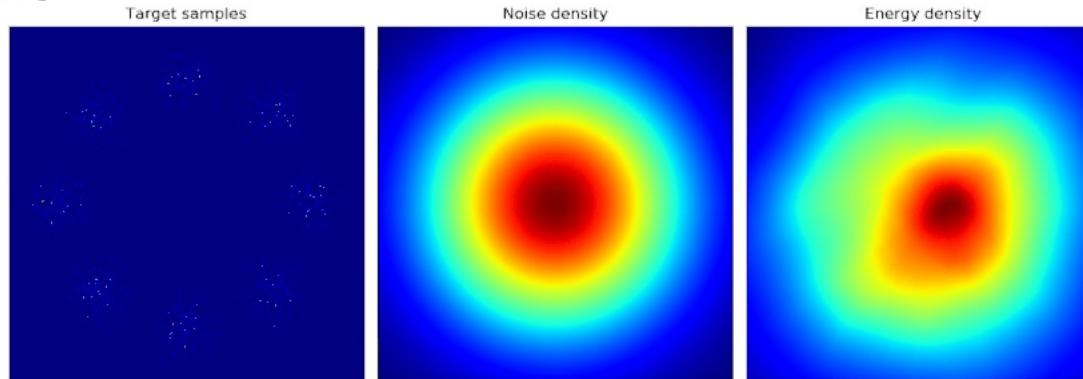
- Sample a mini-batch of datapoints $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \sim p_{data}(\mathbf{x})$
- Sample a mini-batch of datapoints $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(n)} \sim p_n(\mathbf{y})$
- Estimate the NCE loss

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n [f_\theta(\mathbf{x}^{(i)}) - \text{logsumexp}(f_\theta(\mathbf{x}^{(i)}), \log Z + \log p_n(\mathbf{x}^{(i)})) + \log Z \\ & + p_n(\mathbf{y}^{(i)}) - \text{logsumexp}(f_\theta(\mathbf{y}^{(i)}), \log Z + \log p_n(\mathbf{y}^{(i)}))] \end{aligned}$$

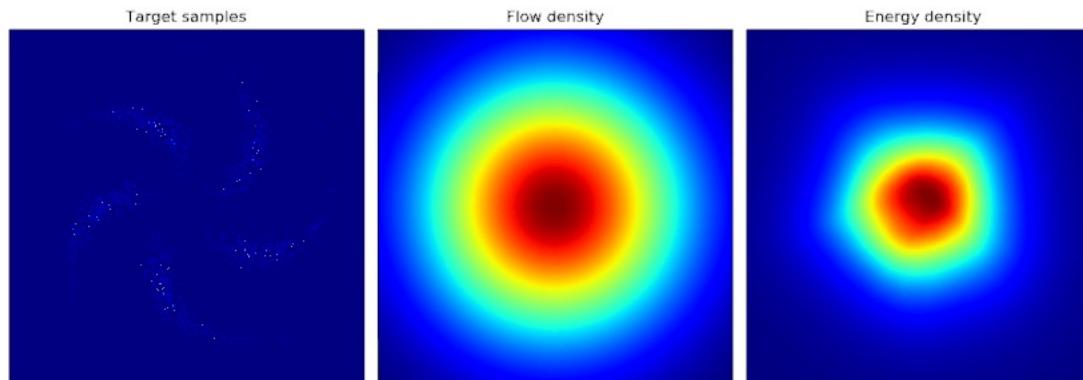
- Stochastic gradient ascent with respect to θ and Z
- No need to sample from the EBM

Noise contrastive estimation

- 8 Gaussians



- Pinwheel



Source: <https://github.com/lifeitech/nce>

Comparing NCE and GAN

- **Similarities**
 - Both involve training a discriminator to perform binary classification with a cross-entropy loss
 - Both are likelihood-free (recall likelihood not tractable in EBM)
- **Differences**
 - GAN requires adversarial training or minimax optimization for training, while NCE does not
 - NCE requires the likelihood of the noise distribution for training, while GAN only requires efficient sampling from the prior
 - NCE trains an energy-based model, while GAN trains a deterministic sample generator

Flow contrastive estimation (Gao et al. 2020)

- **Observations**
 - We need to both evaluate the probability of $p_n(x)$, and sample from it efficiently
 - We hope to make the classification task as hard as possible, i.e., $p_n(x)$ should be close to $p_{data}(x)$ (but not exactly the same)

Flow contrastive estimation (Gao et al. 2020)

- Flow contrastive estimation

- Parameterize the noise distribution with a normalizing flow model $p_{n,\phi}(x)$
- Parameterize the discriminator $D_{\theta,Z,\phi}(x)$ as

$$D_{\theta,Z,\phi}(x) = \frac{\frac{e^{f_\theta(x)}}{Z}}{\frac{e^{f_\theta(x)}}{Z} + p_{n,\phi}(x)} = \frac{e^{f_\theta(x)}}{e^{f_\theta(x)} + Zp_{n,\phi}(x)}$$

- Train the flow model to minimize $D_{JS}(p_{data}, p_{n,\phi})$

$$\min_{\phi} \max_{\theta, Z} E_{x \sim p_{data}} [\log D_{\theta,Z,\phi}(x)] + E_{x \sim p_{n,\phi}} [\log (1 - D_{\theta,Z,\phi}(x))]$$

Flow contrastive estimation (Gao et al. 2020)



Figure 3: Synthesized examples from the Glow model learned by FCE. From left to right panels are from SVHN, CIFAR-10 and CelebA datasets, respectively. The image size is 32×32 .

Source: Gao et al. 2020

Conclusion

- Energy-based models are very flexible probabilistic models with intractable partition functions
- Sampling is hard and typically requires iterative MCMC approaches. Computing the likelihood is hard
- Comparing the likelihood/probability of two different points is tractable
- Maximum likelihood training by contrastive divergence. Requires sampling for each training iteration
- Sampling-free training: score matching, noise contrastive estimation (can provide an estimate of the partition function)
- Reference: How to Train Your Energy-Based Models by Yang Song and Durk Kingma

Thanks
