

Deep Generative Models

10. Energy-Based Models



- 국가수리과학연구소 산업수학혁신센터 김민중

Recap

- Model Families
 - Autoregressive Models: $p_{\theta}(\mathbf{x}) = \prod_{i=1}^d p_{\theta}(x_i | \mathbf{x}_{<i})$
 - Variational Autoencoders: $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
 - Normalizing Flow Models: $p_{\mathbf{X}}(\mathbf{x}; \theta) = p_{\mathbf{Z}}\left(\mathbf{f}_{\theta}^{-1}(\mathbf{x})\right) \left| \det\left(\frac{\partial \mathbf{f}_{\theta}^{-1}(\mathbf{x})}{\partial \mathbf{x}}\right) \right|$
- All the above families are trained by minimizing KL divergence $D(p_{data} \parallel p_{\theta})$ or equivalently maximizing likelihoods (or approximations)

Recap

- Generative Adversarial Networks (GANs)

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}} [\log D_{\phi}(x)] + E_{z \sim p_z} [\log (1 - D_{\phi}(G_{\theta}(z)))]$$

- Two sample tests
- (approximately) optimize f -divergences and the Wasserstein distance
- Very flexible model architectures
- But likelihood is intractable, training is unstable, hard to evaluate, and has mode collapse issues

Plan

- Energy-based models (EBMs).
 - Very flexible model architectures
 - Stable training
 - Relatively high sample quality
 - Flexible composition

Parameterizing probability distributions

- Probability distributions $p(\mathbf{x})$ are a key building block in generative modeling
 - **non-negative:** $p(\mathbf{x}) \geq 0$
 - **sum-to-one:** $\sum_{\mathbf{x}} p(\mathbf{x})$ (or $\int p(\mathbf{x}) d\mathbf{x} = 1$ for continuous variables)
- Condition of non-negative function $p_{\theta}(\mathbf{x})$ is not difficult
- Given any function $f_{\theta}(\mathbf{x})$, we can choose
 - $g_{\theta}(\mathbf{x}) = f_{\theta}(\mathbf{x})^2$
 - $g_{\theta}(\mathbf{x}) = \exp f_{\theta}(\mathbf{x})$
 - $g_{\theta}(\mathbf{x}) = |f_{\theta}(\mathbf{x})|$
 - $g_{\theta}(\mathbf{x}) = \log(1 + \exp f_{\theta}(\mathbf{x}))$
 - etc.
 - In general, g_{θ} is not a normalized function of p_{θ}

Parameterizing probability distributions

- Probability distributions $p(\mathbf{x})$ are a key building block in generative modeling
 - **non-negative:** $p(\mathbf{x}) \geq 0$
 - **sum-to-one:** $\sum_{\mathbf{x}} p(\mathbf{x})$ (or $\int p(\mathbf{x})d\mathbf{x} = 1$ for continuous variables)
- Sum-to-one is key



Parameterizing probability distributions

- Probability distributions $p(\mathbf{x})$ are a key building block in generative modeling
 - **non-negative**: $p(\mathbf{x}) \geq 0$
 - **sum-to-one**: $\sum_{\mathbf{x}} p(\mathbf{x})$ (or $\int p(\mathbf{x}) d\mathbf{x} = 1$ for continuous variables)
- Total “**volume**” is fixed: increasing $p(\mathbf{x}_{train})$ guarantees that \mathbf{x}_{train} becomes relatively more likely
- Problem:
 - $g_{\theta}(\mathbf{x})$ might not sum-to-one
 - $\sum_{\mathbf{x}} g_{\theta}(\mathbf{x}) =: Z(\theta) \neq 1$ in general, so $g_{\theta}(\mathbf{x})$ is not a valid probability mass function or density

Energy-based model

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- I.e., $g_{\theta}(\mathbf{x}) = \exp(f_{\theta}(\mathbf{x}))$
- The volume/normalization constant

$$\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}$$

- is also called the partition function

Energy-based model

- **Why exponential** (and not e.g. $f_{\theta}(\mathbf{x})^2$)?
 - Want to capture very large variations in probability. log-probability is the natural scale we want to work with
Otherwise need highly non-smooth f_{θ} .
 - Exponential families. Many common distributions can be written in this form
 - These distributions arise under general assumptions in statistical physics (maximum entropy, second law of thermodynamics)
 - $-f_{\theta}$ is called the energy, hence the name
 - Intuitively, configurations \mathbf{x} with low energy (high $f_{\theta}(\mathbf{x})$) are more likely

Idea

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- I.e., $p_{\theta}(\mathbf{x}) \propto \exp(f_{\theta}(\mathbf{x}))$
- Given \mathbf{x}, \mathbf{x}' evaluating $p_{\theta}(\mathbf{x})$ or $p_{\theta}(\mathbf{x}')$ requires $Z(\theta)$
- However, the ratio

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

- does not involve $Z(\theta)$

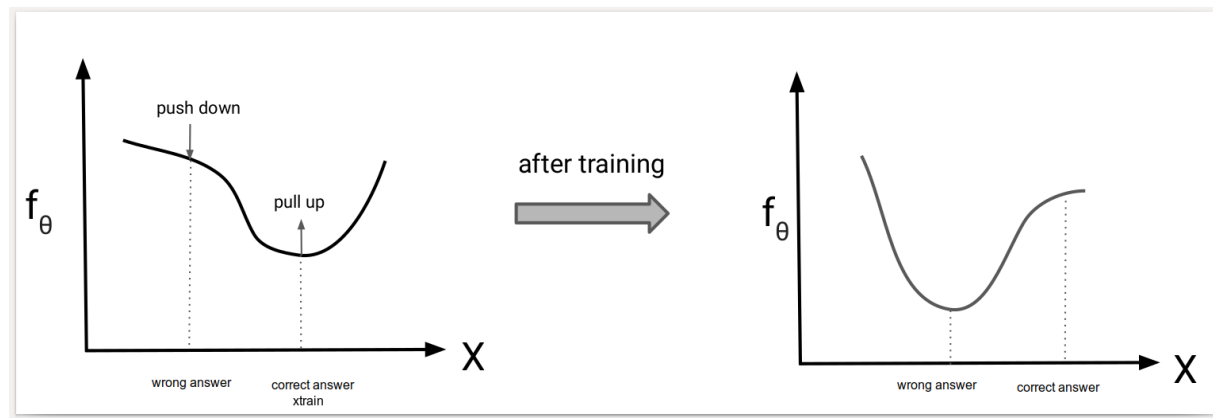
Energy-based model

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- Pros:
 - Extreme flexibility: can use pretty much any function f_{θ}
- Cons:
 - Sampling from p_{θ} is difficult
 - Evaluating and optimizing likelihood p_{θ} is hard (learning is hard)
 - No feature learning (but can add latent variables)
- Curse of dimensionality: The fundamental issue is that computing $Z(\theta)$ numerically (when no analytic solution is available) scales exponentially in the number of dimensions of \mathbf{x}
- Nevertheless, some tasks do not require knowing $Z(\theta)$

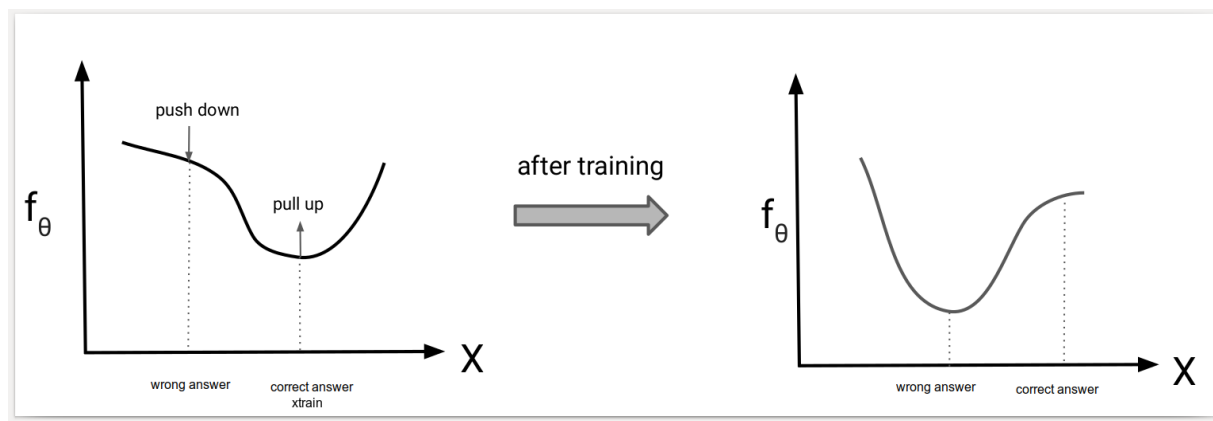
Training intuition

- **Goal:** maximize $\frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}_{train}))$
- **Intuition:** because the model is not normalized, increasing the un-normalized log-probability $f_{\theta}(\mathbf{x}_{train})$ by changing θ does not guarantee that \mathbf{x}_{train} becomes relatively more likely (compared to the rest)
- We also need to consider the effect on other “wrong points” and try to “push them down” to also make $Z(\theta)$ small



Contrastive Divergence

- Goal: maximize $\frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}_{train}))$
- Instead of evaluating $Z(\theta)$ exactly, use a Monte Carlo estimate
- **Contrastive divergence algorithm**
 - Sample $\mathbf{x}_{sample} \sim p_{\theta}$ and maximize $f_{\theta}(\mathbf{x}_{train}) - f_{\theta}(\mathbf{x}_{sample})$
 - Take step on $\nabla_{\theta} (f_{\theta}(\mathbf{x}_{train}) - f_{\theta}(\mathbf{x}_{sample}))$
 - Make training data more likely than typical sample from the model



Contrastive Divergence

- Maximize log-likelihood: $f_{\theta}(\mathbf{x}_{train}) - \log Z(\theta)$
- Gradient of log-likelihood:

$$\begin{aligned}\nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \nabla_{\theta} \log Z(\theta) &= \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \frac{\nabla_{\theta} Z(\theta)}{Z(\theta)} \\&= \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \frac{1}{Z(\theta)} \int \nabla_{\theta} \exp(f_{\theta}(\mathbf{x})) d\mathbf{x} \\&= \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \frac{1}{Z(\theta)} \int \exp(f_{\theta}(\mathbf{x})) \nabla_{\theta} f_{\theta}(\mathbf{x}) d\mathbf{x} \\&= \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \int \frac{\exp(f_{\theta}(\mathbf{x}))}{Z(\theta)} \nabla_{\theta} f_{\theta}(\mathbf{x}) d\mathbf{x} \\&= \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - E_{\mathbf{x} \sim p_{\theta}}[\nabla_{\theta} f_{\theta}(\mathbf{x})]\end{aligned}$$

Contrastive Divergence

- Maximize log-likelihood: $\log p_{\theta}(\mathbf{x}_{train}) = f_{\theta}(\mathbf{x}_{train}) - \log Z(\theta)$
- Gradient of log-likelihood:
$$\nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \nabla_{\theta} \log Z(\theta) = \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - E_{\mathbf{x} \sim p_{\theta}}[\nabla_{\theta} f_{\theta}(\mathbf{x})]$$
$$\approx \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \nabla_{\theta} f_{\theta}(\mathbf{x}_{sample})$$
- where $\mathbf{x}_{sample} \sim p_{\theta}(\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x}))}{Z(\theta)}$
- How to sample?

Sampling from Energy-based model

$$p_{\theta}(\mathbf{x}) = \frac{1}{\int \exp(f_{\theta}(\mathbf{x})) d\mathbf{x}} \exp(f_{\theta}(\mathbf{x})) = \frac{1}{Z(\theta)} \exp(f_{\theta}(\mathbf{x}))$$

- No direct way to sample like in autoregressive or flow models
- **Main issue:** cannot easily compute how likely each possible sample is
- However, we can easily compare two samples \mathbf{x}, \mathbf{x}'
- Use an iterative approach called Markov Chain Monte Carlo:
 - Initialize \mathbf{x}^0 randomly, $t = 0$
 - Let $\mathbf{x}' = \mathbf{x}^t + \text{noise}$
 - If $f_{\theta}(\mathbf{x}') \geq f_{\theta}(\mathbf{x}^t)$, let $\mathbf{x}^{t+1} = \mathbf{x}'$
 - Else let $\mathbf{x}^{t+1} = \mathbf{x}'$ with probability $\exp(f_{\theta}(\mathbf{x}') - f_{\theta}(\mathbf{x}^t))$
- Works in theory, but can take a very long time to converge

Sampling from Energy-based model

- For any continuous distribution $p_\theta(\mathbf{x})$, suppose we can compute its gradient (the **score function**) $\nabla_{\mathbf{x}} \log p_\theta(\mathbf{x})$
- Let $\pi(\mathbf{x})$ be a prior distribution that is easy to sample
- **Langevin MCMC**
 - Initialize $\mathbf{x}^0 \sim \pi(\mathbf{x})$ from prior distribution
 - Repeat $\mathbf{x}^{t+1} \sim \mathbf{x}^t + \epsilon \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}^t) + \sqrt{2\epsilon} \mathbf{z}$ for $t = 0, \dots, T - 1$ where $\mathbf{z} \sim N(0, I)$
 - If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, then we have $\mathbf{x}^T \sim p_\theta$
- Note that for energy-based models, the score function is tractable

$$\nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} f_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log Z(\theta) = \nabla_{\mathbf{x}} f_\theta(\mathbf{x})$$

Training on Energy-based model

- Define the function $f_{\theta}(\mathbf{x})$ parametrized by θ
- Find \mathbf{x}_{sample} that makes $f_{\theta}(\mathbf{x})$ relatively more likely using Langevin MCMC
 - $\mathbf{x}^{t+1} \sim \mathbf{x}^t + \epsilon \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x}_t) + \sqrt{2\epsilon} \mathbf{z}$ for $t = 0, \dots, T - 1$ where $\mathbf{z} \sim N(0, I)$ where ϵ is the step size
- Update the parameter θ

$$\theta^{t+1} = \theta^t + \eta \nabla_{\theta} \left(f_{\theta}(\mathbf{x}_{train}) - f_{\theta}(\mathbf{x}_{sample}) \right)$$

Modern Energy-based model



Figure 1: **Synthesis by short-run MCMC**: Generating synthesized examples by running 100 steps of Langevin dynamics initialized from uniform noise for CelebA (64×64).

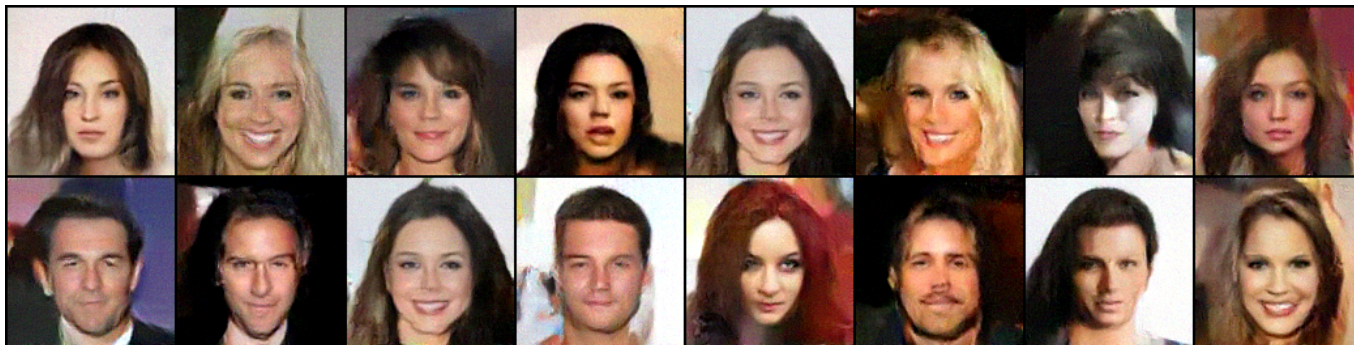


Figure 2: **Synthesis by short-run MCMC**: Generating synthesized examples by running 100 steps of Langevin dynamics initialized from uniform noise for CelebA (128×128).

Source: Nijkamp et al. 2019

Recap. of Energy-based model

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

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

- Maximum likelihood training:

$$\max_{\theta} [f_{\theta}(\mathbf{x}_{train}) - \log Z(\theta)]$$

- Contrastive divergence:

$$\nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \nabla_{\theta} \log Z(\theta) \approx \nabla_{\theta} f_{\theta}(\mathbf{x}_{train}) - \nabla_{\theta} f_{\theta}(\mathbf{x}_{sample})$$

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

Sampling from Energy-based model

- Trained model f_θ is given
- Let $\pi(\mathbf{x})$ be a prior distribution that is easy to sample
- **Langevin MCMC**
 - Initialize $\mathbf{x}^0 \sim \pi(\mathbf{x})$ from prior distribution
 - Repeat $\mathbf{x}^{t+1} \sim \mathbf{x}^t + \epsilon \nabla_{\mathbf{x}} f_\theta(\mathbf{x}^t) + \sqrt{2\epsilon} \mathbf{z}$ for $t = 0, \dots, T - 1$ where $\mathbf{z} \sim N(0, I)$

Thanks
