Deep Generative Models

5. Latent variable models



• 국가수리과학연구소 산업수학혁신센터 김민중

Recap

- Autoregressive models:
 - Chain rule-based factorization is fully general
 - Compact representation via conditional independence and/or neural parameterizations
- Autoregressive models Pros
 - Easy to evaluate likelihoods
 - Easy to train
- Autoregressive models Cons
 - Requires an ordering
 - Generation is sequential
 - Cannot learn features in an unsupervised way

Plan for this lecture

- Latent Variable Models
 - Mixture models
 - Variational autoencoder
 - Variational inference and learning

Latent Variable Models: Assumption

- Observable variables $x \in \mathbb{R}^d$
- Latent variables $z \in \mathbb{R}^h$ (unobservable)

$$p_{data}(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$$

or
$$= \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

Latent Variable Models: Motivation

- Lots of variability in images x due to gender, eye color, hair color, pose, etc.
- However, unless images are annotated, these factors of variation are not explicitly available (latent)
- Idea: explicitly model these factors using latent variables z



Latent Variable Models

- Only shaded variables *x* are observed in the data (pixel values)
- Latent variables z correspond to high level features
 - If z chosen properly, p(x|z) could be much simpler than p(x)
 - If we had trained this model, then we could identify features via p(z|x), e.g., p(EyeColor = Blue|x)
- Challenge: Very difficult to specify these conditionals by hand



Why Latent Variable Models?

- Curse of Dimensionality
 - Sparse and peaky
- Model $p(\mathbf{x}, \mathbf{z})$ instead of $p(\mathbf{x})$
- Maximizing $p_{\theta}(\mathbf{x}, \mathbf{z})$ and Maximizing $p_{\theta}(\mathbf{x})$
 - Easy to sample a new data
 - Neural network can learn and model a mapping effectively. E.g., $p_{\theta}(\mathbf{x}|\mathbf{z})$

Mixture of Gaussians: a Shallow Latent Variable Model

- Mixture of Gaussians. Bayes net: $z \rightarrow x$
 - $z = Categorical(z|\gamma_1, \cdots, \gamma_K)$
 - $p(\boldsymbol{x}|\boldsymbol{z}=\boldsymbol{k}) = N(\boldsymbol{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)$



- Generative Process
 - Pick a mixture component k by sampling z
 - Generate a data point by sampling from that Gaussian

Mixture of Gaussians: a Shallow Latent Variable Model

- Mixture of Gaussians. Bayes net: $z \rightarrow x$
 - $z = Categorical(z|\gamma_1, \cdots, \gamma_K)$
 - $p(\boldsymbol{x}|\boldsymbol{z}=\boldsymbol{k}) = N(\boldsymbol{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)$



• Clustering: The posterior $p(z|\mathbf{x})$ identifies the mixture component

Mixture Models

• Alternative motivation: Combine simple models into a more complex and expressive one

$$p(\mathbf{x}) = \sum_{Z} p(\mathbf{x}, z) = \sum_{Z} p(z)p(\mathbf{x}|z) = \sum_{Z} p(z = k)N(\mathbf{x}|\mu_k, \Sigma_k)$$

Example: Variational Autoencoder



- A Mixture of an infinite number of Gaussians
 - $\boldsymbol{z} = N(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}), \, \boldsymbol{z} \in \mathbb{R}^h$
 - $p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}|\mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$ where $\mu_{\theta}, \Sigma_{\theta}$ are neural networks
 - Even though $p(\mathbf{x}|\mathbf{z})$ is simple, the marginal $p(\mathbf{x})$ is very complex/flexible
- Hope that after training, z will correspond to meaningful latent factors of variation (features)
- Unsupervised representation learning
- Features can be computed via $p(\mathbf{z}|\mathbf{x})$

Example: Variational Autoencoder



• A Mixture of an infinite number of Gaussians

•
$$\boldsymbol{z} = N(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}), \, \boldsymbol{z} \in \mathbb{R}^h$$

- $p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}|\mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$ where $\mu_{\theta}, \Sigma_{\theta}$ are neural networks
 - $\mu_{\theta}(\mathbf{z}) = \sigma(A\mathbf{z} + \mathbf{c}) = (\sigma(\mathbf{a}_1^T\mathbf{z} + c_1), \sigma(\mathbf{a}_2^T\mathbf{z} + c_2))$

•
$$\Sigma_{\theta}(\mathbf{z}) = diag(\exp \sigma(B\mathbf{z} + \mathbf{d})) =$$

 $\begin{pmatrix} \exp \sigma(\mathbf{b}_1^T \mathbf{z} + d_1) & 0 \\ 0 & \exp \sigma(\mathbf{b}_2^T \mathbf{z} + d_2) \end{pmatrix}$
• $\theta = (A, B, \mathbf{c}, \mathbf{d})$

Recap

- Latent Variable Models
 - Allow us to define complex models p(x) in terms of simpler building blocks p(x|z)
 - Natural for unsupervised learning tasks (clustering, unsupervised representation learning, etc.)
 - No free lunch: much more difficult to learn compared to fully observed, autoregressive models

Marginal Likelihood

- Let *X* denote observed random variables, and *Z* the unobserved hidden or latent variables
- Suppose we have a model for the joint distribution

 $p_{\theta}(X,Z)$ • What is the probability $p_{\theta}(X = \bar{x})$ of observing a training data point \bar{x} ?

$$p_{\theta}(\bar{x}) = \sum_{Z} p_{\theta}(X = \bar{x}, Z = Z) = \sum_{Z} p_{\theta}(\bar{x}, Z)$$

Partially observed data

- Suppose we have a model for the joint distribution
- $p_{\theta}(X,Z)$ • We have a dataset *D*, where for each datapoint the *X* variables are observed (e.g., pixel values) and the latent variables *Z* are never observed. $D = \{x^{(1)}, \dots, x^{(N)}\}$
- Maximum likelihood learning

$$\log \prod_{\boldsymbol{x} \in D} p_{\theta}(\boldsymbol{x}) = \sum_{\boldsymbol{x} \in D} \log p_{\theta}(\boldsymbol{x}) = \sum_{\boldsymbol{x} \in D} \log \sum_{\boldsymbol{x}} p_{\theta}(\boldsymbol{x}, \boldsymbol{z})$$

- Evaluating $\log \sum_{z} p_{\theta}(x, z)$ can be intractable. For continuous variables, $\log \int p_{\theta}(x, z) dz$ is often intractable
- Gradients ∇_{θ} also hard to compute
- Need approximations: One gradient evaluation per training data point $x \in D$, so approximation needs to be cheap

Attempt: Importance Sampling

 Likelihood function p_θ(X) for partially observed data is hard to compute:

$$p_{\theta}(\boldsymbol{x}) = \sum_{\boldsymbol{z} \in \mathcal{Z}} p_{\theta}(\boldsymbol{x}, \boldsymbol{z}) = \sum_{\boldsymbol{z} \in \mathcal{Z}} \frac{q(\boldsymbol{z})}{q(\boldsymbol{z})} p_{\theta}(\boldsymbol{x}, \boldsymbol{z}) = E_{\boldsymbol{z} \sim q(\boldsymbol{z})} \left[\frac{p_{\theta}(\boldsymbol{x}, \boldsymbol{z})}{q(\boldsymbol{z})} \right]$$

- Monte Carlo to the rescue:
 - Sample $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \cdots, \mathbf{z}^{(K)}$ from $q(\mathbf{z})$
 - Approximate expectation with sample average

$$p_{\theta}(\boldsymbol{x}) \approx \frac{1}{K} \sum_{j=1}^{K} \frac{p_{\theta}(\boldsymbol{x}, \boldsymbol{z}^{(j)})}{q(\boldsymbol{z}^{(j)})}$$

• What is a good choice for q(z)? Intuitively, frequently sample z that are likely given x under $p_{\theta}(x, z)$.

Estimating log-likelihoods

- Monte Carlo to the rescue:
 - Sample $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \cdots, \mathbf{z}^{(K)}$ from $q(\mathbf{z})$
 - Approximate expectation with sample average(unbiased)

$$p_{\theta}(\boldsymbol{x}) \approx \frac{1}{K} \sum_{j=1}^{K} \frac{p_{\theta}(\boldsymbol{x}, \boldsymbol{z}^{(j)})}{q(\boldsymbol{z}^{(j)})}$$

• Recall that for training, we need the log-likelihood $\log p_{\theta}(x)$

$$\log p_{\theta}(\boldsymbol{x}) \approx \log \left(\frac{1}{K} \sum_{j=1}^{K} \frac{p_{\theta}(\boldsymbol{x}, \boldsymbol{z}^{(j)})}{q(\boldsymbol{z}^{(j)})} \right) \approx \log \left(\frac{p_{\theta}(\boldsymbol{x}, \boldsymbol{z}')}{q(\boldsymbol{z}')} \right)$$

• However, it is clear that

$$\log\left(E_{\boldsymbol{z}\sim q(\boldsymbol{z})}\left[\frac{p_{\theta}(\boldsymbol{x},\boldsymbol{z})}{q(\boldsymbol{z})}\right]\right) \neq E_{\boldsymbol{z}\sim q(\boldsymbol{z})}\left[\log\left(\frac{p_{\theta}(\boldsymbol{x},\boldsymbol{z})}{q(\boldsymbol{z})}\right)\right]$$

Evidence Lower Bound

Log-Likelihood function

$$\log \sum_{\mathbf{z} \in \mathcal{Z}} p_{\theta}(\mathbf{x}, \mathbf{z}) = \log \sum_{\mathbf{z} \in \mathcal{Z}} \frac{q(\mathbf{z})}{q(\mathbf{z})} p_{\theta}(\mathbf{x}, \mathbf{z}) = \log \left(E_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \right)$$

- log is a concave function $\log(px + (1-p)x') \ge p\log(x) + (1-p)\log(x')$
- By Jensen Inequality(for concave function),

$$\log(E_{z \sim q(z)}[f(z)]) = \log\left(\sum_{z} q(z)f(z)\right) \ge \sum_{z} q(z)\log f(z)$$

Deep Generative Models

mjgim@nims.re.kr

NIMS & AJOU University

Evidence Lower Bound

Log-Likelihood function

$$\log \sum_{\mathbf{z} \in \mathcal{Z}} p_{\theta}(\mathbf{x}, \mathbf{z}) = \log \sum_{\mathbf{z} \in \mathcal{Z}} \frac{q(\mathbf{z})}{q(\mathbf{z})} p_{\theta}(\mathbf{x}, \mathbf{z}) = \log \left(E_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \right)$$

- log is a concave function $\log(px + (1-p)x') \ge p \log(x) + (1-p) \log(x')$
- By Jensen Inequality(for concave function),

$$\log(E_{\boldsymbol{z}\sim q(\boldsymbol{z})}[f(\boldsymbol{z})]) = \log\left(\sum_{\boldsymbol{z}} q(\boldsymbol{z})f(\boldsymbol{z})\right) \ge \sum_{\boldsymbol{z}} q(\boldsymbol{z})\log f(\boldsymbol{z})$$

• Choosing $f(\mathbf{z}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})}$ $\log\left(E_{\mathbf{z}\sim q(\mathbf{z})}\left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})}\right]\right) \ge \sum_{\mathbf{z}} q(\mathbf{z})\log\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})}$

Variational inference

- Suppose q(z) is any probability distribution over the hidden variables
- Evidence lower bound (ELBO) holds for any q(z)

$$\log p_{\theta}(\mathbf{x}) \geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})}$$
$$= \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \sum_{\mathbf{z}} q(\mathbf{z}) \log q(\mathbf{z})$$
$$= \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) + H(q)$$

• Claim: Equality holds if $q(z) = p_{\theta}(z|x)$

Variational inference

• We derived this lower bound that holds holds for any choice of q(z)

$$\log p_{\theta}(x) \geq \sum_{z} q(z) \log \frac{p_{\theta}(x, z)}{q(z)}$$
• If $q(z) = p_{\theta}(z|x)$, then the bound becomes

$$\sum_{z} p_{\theta}(z|x) \log \frac{p_{\theta}(x, z)}{p_{\theta}(z|x)} = \sum_{z} p_{\theta}(z|x) \log \frac{p_{\theta}(z|x)p_{\theta}(x)}{p_{\theta}(z|x)}$$

$$= \sum_{z} p_{\theta}(z|x) \log p_{\theta}(x) = \log p_{\theta}(x) \sum_{z} p_{\theta}(z|x) = \log p_{\theta}(x)$$
• What if the posterior $p_{\theta}(z|x)$ is intractable to compute? How

• What if the posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ is intractable to compute? How loose is the bound?

Variational inference(continued)

 Suppose q(z) is any probability distribution over the hidden variables. A little bit of algebra reveals

$$D(q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) = -\sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) + \log p_{\theta}(\mathbf{x}) - H(q) \ge 0$$

Rearranging, we re-derived the Evidence lower bound (ELBO)

$$\log p_{\theta}(\boldsymbol{x}) \geq \sum_{\boldsymbol{z}} q(\boldsymbol{z}) \log p_{\theta}(\boldsymbol{x}, \boldsymbol{z}) + H(q)$$

- Equality holds if $q(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})$ because $D(q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) = 0$
- In general, $\log p_{\theta}(\mathbf{x}) = \text{ELBO} + D(q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}))$

The Evidence Lower bound

- What if the posterior $p_{\theta}(\boldsymbol{z}|\boldsymbol{x})$ is intractable to compute?
- Suppose $q_{\phi}(z)$ is a (tractable) probability distribution over the hidden variables parameterized by ϕ (variational parameters)
 - For example, a Gaussian with mean and covariance specified by ϕ

$$q_{\phi}(\mathbf{z}) = N(\mathbf{z}|\boldsymbol{\mu}_{\phi}, \boldsymbol{\Sigma}_{\phi})$$

• Variational inference: pick ϕ so that $q_{\phi}(z)$ is as close as possible to $p_{\theta}(z|x)$

The Evidence Lower bound



- The better $q_{\phi}(z)$ can approximate the posterior $p_{\theta}(z|x)$, the smaller $D\left(q_{\phi}(z) \parallel p_{\theta}(z|x)\right)$ we can achieve, the closer ELBO will be to $\log p_{\theta}(x)$
- Next: jointly optimize over θ and ϕ to maximize the ELBO over a dataset

Summary

- Latent Variable Models Pros:
 - Easy to build flexible models
 - Suitable for unsupervised learning
- Latent Variable Models Cons
 - Hard to evaluate likelihoods
 - Hard to train via maximum-likelihood
 - Fundamentally, the challenge is that posterior inference $p_{\theta}(\boldsymbol{z}|\boldsymbol{x})$ is hard. Typically requires variational approximations
- Alternative: give up on KL-divergence and likelihood (GANs)

Thanks