



Deep Generative Models

Qinliang Su (苏勤亮)

Sun Yat-sen University

suqliang@mail.sysu.edu.cn

Outline

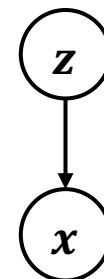
- Deep Generative Model
- Learning under the VB-EM Framework
- Estimating the Gradient Naively
- Estimating the Gradient using Re-parameterization Trick
- Amortizing the Inference

Generative Models

- Describing the data generation process by a joint pdf

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

where \mathbf{x} and \mathbf{z} denote data and latent variable



Obviously, generative model is a kind of latent-variable models

- $p(\mathbf{z})$ and $p(\mathbf{x}|\mathbf{z})$ are chosen by taking following factors into account

1) Data compatibility 2) Modeling flexibilities 3) Training easiness

- Examples

PCA: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ and $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$

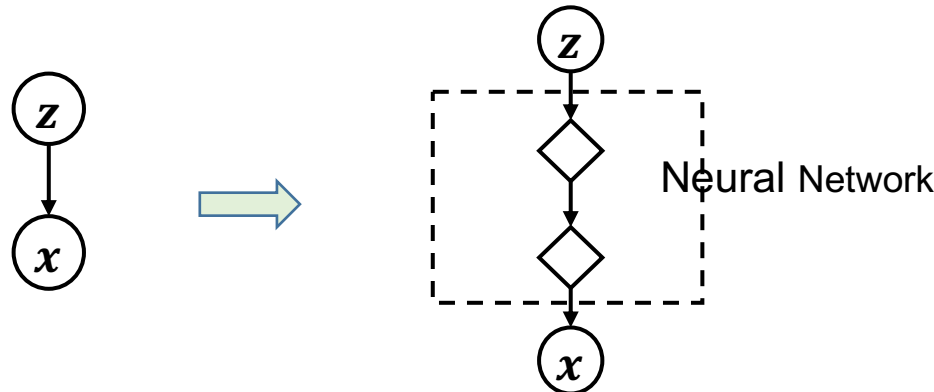
GMM: $p(\mathbf{z}) = \text{Cat}(\mathbf{z}; \boldsymbol{\pi})$ and $p(\mathbf{x}|\mathbf{z}) = \prod_{m=1}^M [\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)]^{z_m}$

Deep Generative Models

- To increase modeling ability, a deep neural network is introduced between z and x . Then, the joint pdf becomes

$$p(\mathbf{x}, \mathbf{z}) = \underbrace{p(\mathbf{x} | T(\mathbf{z}))}_{p(\mathbf{x} | \mathbf{z})} p(\mathbf{z})$$

where $T(\cdot)$ represents a neural network



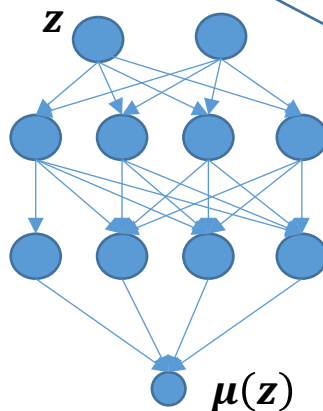
- Comparing to generative models (GM), deep GMs (DGM) let the *conditional pdf rely on a neural-network-transformed variable $T(z)$*

- **Example:** To model images, we can specify the joint pdf as

$$p(\mathbf{x}, \mathbf{z}) = \underbrace{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}(\mathbf{z}), \mathbf{I})}_{p(\mathbf{x}|\mathbf{z})} \underbrace{\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})}_{p(\mathbf{z})}$$

where $\boldsymbol{\mu}(\mathbf{z})$ denotes the output of a neural network, *e.g.*,

$$\boldsymbol{\mu}(\mathbf{z}) = \mathbf{W}_3 a(\mathbf{W}_2 a(\mathbf{W}_1 \mathbf{z} + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_3$$



In PCA, $\boldsymbol{\mu}(\mathbf{z}) = \mathbf{W}\mathbf{z} + \mathbf{b}$

Denote the parameters in NNs as $\boldsymbol{\theta}$, *i.e.*, $\boldsymbol{\theta} \triangleq \{\mathbf{W}_\ell, \mathbf{b}_\ell\}_{\ell=1}^3$

Outline

- Deep Generative Model
- Learning under the VB-EM Framework
- Estimating the Gradient Naively
- Estimating the Gradient using Re-parameterization Trick
- Amortizing the Inference

Learning under VB-EM Framework

- DGM is a latent-variable model, hence can be trained with the EM algorithms

Key steps in EM

- 1) The posteriori distribution $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$
- 2) Deriving the expectation $\mathbb{E}_{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})}[\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})]$
- 3) Maximization

- But due to the existence of neural networks, *the exact posterior $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$ and expectation $\mathbb{E}[\cdot]$ are difficult to obtain*
- Thus, we resort to VB-EM algorithm by *using a simple distribution to approximate the true posterior $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$*

- VB-EM seeks to maximize the lower bound of log-likelihood

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \int q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

where $q(\mathbf{z}|\mathbf{x}; \boldsymbol{\phi})$ denotes the approximate posterior

- Substituting $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ into $\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$ gives

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) &= \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) d\mathbf{z} - \int_{\mathbf{z}} q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \ln \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} d\mathbf{z} \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \end{aligned}$$

- In the subsequent, we consider a concrete example, in which $p_{\theta}(\mathbf{x}|\mathbf{z})$ and $q_{\phi}(\mathbf{z}|\mathbf{x})$ are set as diagonal Gaussian form

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\theta}(\mathbf{z}), \mathbf{I}) \quad q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\lambda}, \text{diag}(\boldsymbol{\eta}^2))$$

where $\boldsymbol{\mu}_{\theta}(\cdot)$ represents a neural network function with parameter denoted as $\boldsymbol{\theta}$; and $\boldsymbol{\phi} = \{\boldsymbol{\lambda}, \boldsymbol{\eta}\}$ denotes the posterior parameter

- To train the model, we can maximize the lower bound $\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$ w.r.t. the model and posterior parameter $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$. To this end, what we need is the gradient

$$\frac{\partial \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\theta}} \quad \text{and} \quad \frac{\partial \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\phi}}$$

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

- How to obtain the gradient of $\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi})$ w.r.t. $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$?
 - 1) Since $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ are both Gaussian, the close-form expression of $KL(q_{\boldsymbol{\phi}}||p)$ can be easily obtained, and so does its gradient
 - 2) Due to the existence of neural networks, the close-form expression of $\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]$ cannot be obtained. So, its gradient expression cannot be obtained directly
- Now, the only problem left is how to estimate the derivatives of $\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]$ w.r.t. $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$, i.e.,

$$\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\theta}}$$

$$\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\phi}}$$

Outline

- Deep Generative Model
- Learning under the VB-EM Framework
- **Estimating the Gradient Naively**
- Estimating the Gradient using Re-parameterization Trick
- Amortizing the Inference

Estimating $\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})] / \partial \theta$

$$\begin{aligned} \frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \theta} &= \frac{\partial}{\partial \theta} \int q_{\phi}(\mathbf{z}|\mathbf{x}) \ln p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z})}{\partial \theta} d\mathbf{z} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z})}{\partial \theta} \right] \end{aligned}$$

where $\ln p_{\theta}(\mathbf{x}|\mathbf{z}) = C - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_{\theta}(\mathbf{z})\|^2$, and C is a constant

- Due to the Gaussian form of $p_{\theta}(\mathbf{x}|\mathbf{z})$, $\frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z})}{\partial \theta}$ can be obtained by using the BP algorithm
- The only obstacle left is about **how to evaluate the expectation**

$$\frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \theta} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z})}{\partial \theta} \right]$$

- Due to the existence of neural network, it is impossible to derive a close-form expression for the expectation, even if $\frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z})}{\partial \theta}$ is known
- Instead, we can approximate the expectation with its empirical average, that is,

$$\frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \theta} \approx \frac{1}{K} \sum_{k=1}^K \frac{\partial \ln p_{\theta}(\mathbf{x}|\mathbf{z}^{(k)})}{\partial \theta}$$

where $\mathbf{z}^{(k)}$ are samples drawn from $q_{\phi}(\mathbf{z}|\mathbf{x})$

- From $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \lambda, \text{diag}(\boldsymbol{\eta}^2))$, it is easy to draw samples from the Gaussian distribution **under a given value of λ and $\boldsymbol{\eta}^2$** , e.g., the parameter value λ_t and $\boldsymbol{\eta}_t$ at the t -th step

- In summary, given the model parameter θ_t and posterior parameter values λ_t and η_t , we can estimate the gradient $\left. \frac{\partial \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\ln p_\theta(\mathbf{x}|\mathbf{z})]}{\partial \theta} \right|_{\theta=\theta_t}$ with the following steps
 - Draw K samples $\{\mathbf{z}^{(k)}\}_{k=1}^K$ from the posterior $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \lambda_t, \text{diag}(\eta_t^2))$
 - Use BP to evaluate the gradient $\left. \frac{\partial \ln p_\theta(\mathbf{x}|\mathbf{z}^{(k)})}{\partial \theta} \right|_{\theta=\theta_t}$ at model parameter θ_t and latent value $\mathbf{z}^{(k)}$
 - Estimate the gradient by empirical average $\frac{1}{K} \sum_{k=1}^K \left. \frac{\partial \ln p_\theta(\mathbf{x}|\mathbf{z}^{(k)})}{\partial \theta} \right|_{\theta=\theta_t}$

Question: To estimate $\frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \theta}$, can we first estimate the expectation $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]$ by its empirical average at the given model parameter θ_t and posterior parameter λ_t and η_t as

$$\frac{1}{K} \sum_{k=1}^K \ln p_{\theta_t}(\mathbf{x}|\mathbf{z}^{(k)}),$$

and then compute the gradient of the empirical average w.r.t θ ?

- How to estimate the gradient if we want to use this method?
 - Estimate the expression of $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]$ for all possible values of θ

$$\frac{1}{K} \sum_{k=1}^K \ln p_{\theta}(\mathbf{x}|\mathbf{z}^{(k)}),$$

rather than the value only at $\theta = \theta_t$

- Obviously, it will lead to the same gradient expression as the previous method

Estimating $\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})] / \partial \phi$

$$\begin{aligned} \frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \phi} &= \frac{\partial}{\partial \phi} \int q_{\phi}(\mathbf{z}) \ln p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \int \frac{\partial q_{\phi}(\mathbf{z}|\mathbf{x})}{\partial \phi} \ln p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \frac{\partial \ln q_{\phi}(\mathbf{z})}{\partial \phi} \ln p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{\partial \ln q_{\phi}(\mathbf{z}|\mathbf{x})}{\partial \phi} \ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \end{aligned}$$

- Similarly, we can also use the empirical average to approximate the expectation, that is,

$$\frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \phi} \approx \frac{1}{K} \sum_{k=1}^K \frac{\partial \ln q_{\phi}(\mathbf{z}^{(k)}|\mathbf{x})}{\partial \phi} \ln p_{\theta}(\mathbf{x}|\mathbf{z}^{(k)})$$

Question: Can we first use the samples $\{\mathbf{z}^{(k)}\}_{k=1}^K$ drawn from the posterior $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\lambda}_t, \text{diag}(\boldsymbol{\eta}_t^2))$ to estimate the expression of $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]$ w.r.t. $\boldsymbol{\phi}$, and then use the expression to derive the gradient?

➤ Then answer is no. **Why?**

- For the method proposed above, it has been observed that

the variance of the estimated derivative $\frac{1}{K} \sum_{k=1}^K \frac{\partial \ln q_{\phi}(\mathbf{z}|\mathbf{x})}{\partial \boldsymbol{\phi}} \ln p_{\theta}(\mathbf{x}|\mathbf{z})$ is very large, making the estimate **unreliable**

- Thus, a new method is needed to estimate the derivative
$$\frac{\partial \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\phi}}$$

Outline

- Deep Generative Model
- Learning under the VB-EM Framework
- Estimating the Gradient Naively
- Estimating the Gradient using Re-parameterization Trick
- Amortizing the Inference

Re-parameterization Trick

Re-parameterization Trick: For any sample $\mathbf{z}^{(i)}$ drawn from the distribution $q_{\phi}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\lambda}, \text{diag}(\boldsymbol{\eta}^2))$, it can be represented as

$$\mathbf{z}^{(i)} = \boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)}$$

where $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_M]$, $\boldsymbol{\eta} = [\eta_1, \dots, \eta_M]$ and $\boldsymbol{\epsilon}^{(i)} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$ (i.e., standard Gaussian noise)

To see this, we can prove that if $\boldsymbol{\epsilon}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then the sample $\mathbf{z}^{(i)} = \boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)}$ follows the distribution $\mathcal{N}(\mathbf{z}; \boldsymbol{\lambda}, \text{diag}(\boldsymbol{\eta}^2))$, that is,

$$\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z})$$

- What is the key differences between Re-parameterization trick and traditional sampling methods?

Remark: We can also use MCMC to draw samples from $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\lambda}, \text{diag}(\boldsymbol{\eta}^2))$

Samples drawn with the re-parameterized method **include the unknown parameters $\boldsymbol{\lambda}$ and $\boldsymbol{\eta}$ explicitly**, while the traditional sampling methods cannot

- It can be seen that the re-parameterized sample

$$\mathbf{z}^{(i)} = \boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)}$$

can separate the parameters from the model randomness

Estimating Derivatives with the Re-parameterization Trick

- Using the re-parameterization trick, the expectation $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})]$ can be estimated as

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})] \approx \frac{1}{K} \sum_{i=1}^K \ln p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)})$$

where

$$\mathbf{z}^{(i)} = \boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)}$$

- Substituting it into $\ln p_{\theta}(\mathbf{x}|\mathbf{z}) = C - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_{\theta}(\mathbf{z})\|^2$ gives

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\ln p_{\theta}(\mathbf{x}|\mathbf{z})] \approx \frac{1}{K} \sum_{i=1}^K \left(C - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_{\theta}(\boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)})\|^2 \right)$$

where $\boldsymbol{\epsilon}^{(i)}$ is a random noise from standard Gaussian $\mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$

- Therefore, the derivatives $\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\theta}}$ and $\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\phi}}$ can be estimated from the approximate function

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{K} \sum_{i=1}^K \left(C - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\boldsymbol{\lambda} + \boldsymbol{\eta} \cdot \boldsymbol{\epsilon}^{(i)})\|^2 \right)$$

- That is,

$$\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\theta}} \approx \frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\theta}}$$

$$\frac{\partial \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\ln p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}{\partial \boldsymbol{\phi}} \approx \frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\phi}}$$

- $\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\theta}}$ and $\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi})}{\partial \boldsymbol{\phi}}$ can be evaluated with BP algorithm with the automatic tools

Outline

- Deep Generative Model
- Learning under the VB-EM Framework
- Estimating the Gradient Naively
- Estimating the Gradient using Re-parameterization Trick
- Amortizing the Inference

Extending to Large Datasets

- So far, only one training example is considered
- When a training set containing N examples $\mathcal{X} = \{\mathbf{x}_n\}_{n=1}^N$ is considered, the training objective becomes

$$\mathcal{L}(\mathcal{X}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{N} \sum_{n=1}^N \left(\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}_n)} [\ln p_{\boldsymbol{\theta}}(\mathbf{x}_n|\mathbf{z})] - KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}_n) || p(\mathbf{z})) \right)$$

- To optimize the objective, it is better to let each \mathbf{x}_n has its own $\boldsymbol{\phi}_n = \{\boldsymbol{\lambda}_n, \boldsymbol{\eta}_n\}$, that is,

$$q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n; \boldsymbol{\lambda}_n, \text{diag}(\boldsymbol{\eta}_n^2))$$

- With the re-parameterization trick, a sample $\mathbf{z}_n^{(k)}$ from $q_{\phi}(\mathbf{z}_n|\mathbf{x}_n)$ can be represented as

$$\mathbf{z}_n^{(k)} = \boldsymbol{\lambda}_n + \boldsymbol{\eta}_n \cdot \boldsymbol{\epsilon}_n^{(k)}$$

- Then, the objective $\mathcal{L}(\mathcal{X}; \boldsymbol{\theta}, \boldsymbol{\phi})$ can be approximated as

$$\mathcal{L} \approx \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \left(C - \frac{1}{2} \left\| \mathbf{x}_n - \boldsymbol{\mu}_{\boldsymbol{\theta}} \left(\boldsymbol{\lambda}_n + \boldsymbol{\eta}_n \cdot \boldsymbol{\epsilon}_n^{(k)} \right) \right\|^2 \right) - \frac{1}{N} \sum_{n=1}^N KL(q_{\phi}(\mathbf{z}_n|\mathbf{x}_n) || p(\mathbf{z}_n))$$

- Training complexity
 - $\boldsymbol{\phi}_n$ is only updated *on* the data sample $\mathbf{x}^{(n)}$, while $\boldsymbol{\theta}$ will be updated on all samples
 - Because $\boldsymbol{\phi}_n$ is updated much less frequent than $\boldsymbol{\theta}$, to ensure $q_{\phi}(\mathbf{z}_n|\mathbf{x}_n)$ is a good approximate to true posterior $p_{\boldsymbol{\theta}_t}(\mathbf{z}_n|\mathbf{x}_n)$, the parameter $\boldsymbol{\phi}_n$ has to be updated much more times than $\boldsymbol{\theta}$

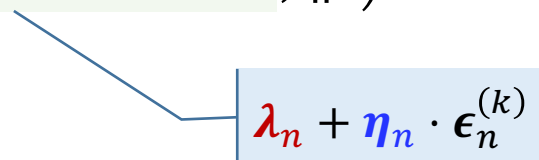
Amortizing the Burden of Inference

- Instead of learning λ_n and η_n directly, we set λ_n and η_n as the outputs of neural networks, that is,

$$\lambda_n = g_{\phi_1}(x_n) \quad \eta_n = g_{\phi_2}(x_n)$$

where $g_{\phi_\ell}(\cdot)$ represents neural networks parameterized by ϕ_ℓ

- Then, the objective $\mathcal{L}(\mathcal{X}; \theta, \phi)$ can be written as

$$\begin{aligned} \mathcal{L} \approx & \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \left(C - \frac{1}{2} \left\| x_n - \mu_\theta \left(g_{\phi_1}(x_n) + g_{\phi_2}(x_n) \cdot \epsilon_n^{(k)} \right) \right\|^2 \right) \\ & - \frac{1}{N} \sum_{n=1}^N KL(q_\phi(z_n | x_n) || p(z)) \end{aligned}$$


- Here, the parameters to be optimized are $\{\theta, \phi_1, \phi_2\}$

- Differences between the *non-amortized* and *amortized* methods
 - For the non-amortized method, a separate parameter $\{\lambda_n, \eta_n\}$ should be learned for each data example x_n
 - For the amortized method, common deep neural networks are learned for all data examples $\{x_n\}_{i=1}^N$
- **Question:** Learning a DNN is much more expensive than learning one $\{\lambda_i, \eta_i\}$. Then, why we choose to learn DNNs for inference?
 - The computation burden of learning a DNN is **shouldered by all examples**, while the complexity of learning $\{\lambda_n, \eta_n\}$ is undertaken only by the n -th example x_n
 - The training complexity amortized on each data sample is low

Examining the Training Process

- By looking at the training objective

$$\mathcal{L} \approx \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \left(C - \frac{1}{2} \left\| \mathbf{x}_n - \mu_{\theta} \left(g_{\phi_1}(\mathbf{x}_n) + g_{\phi_2}(\mathbf{x}_n) \cdot \epsilon_n^{(k)} \right) \right\|^2 \right)$$

Diagram illustrating the components of the training objective:

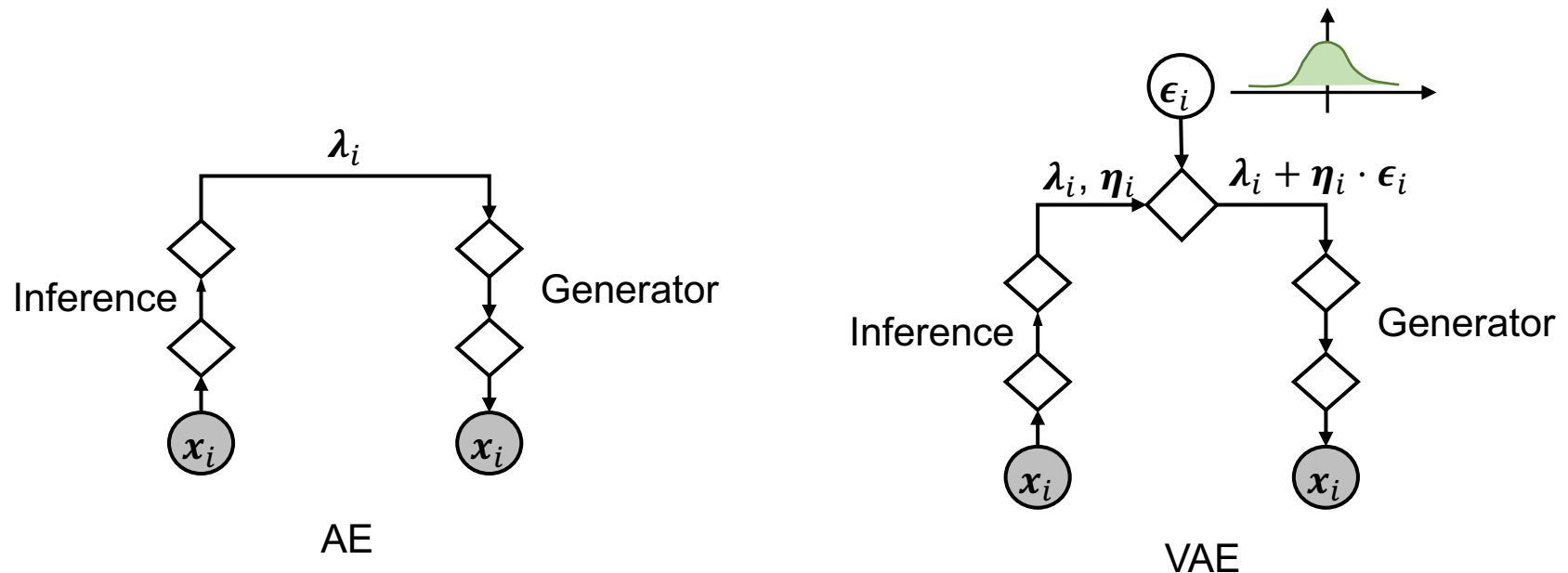
- μ_{θ} is labeled as the **decoder**.
- $g_{\phi_1}(\mathbf{x}_n)$ is labeled as the **encoder**.
- $\epsilon_n^{(k)}$ is labeled as **noise**.

$$- \frac{1}{N} \sum_{n=1}^N KL(q_{\phi}(\mathbf{z}|\mathbf{x}_n) || p(\mathbf{z}))$$

we see that the proposed model is learning an encoder $g_{\phi_{\ell}}(\cdot)$ and a decoder $\mu_{\theta}(\cdot)$ for reconstruction, with two additional features

- 1) An additional KL regularizer
- 2) Imposing noise on the latent code $g_{\phi_1}(\mathbf{x}_n)$

- The auto-encoder (AE) is to build an encoder and decoder such that the reconstruction loss is minimized



- Comparing to AE, VAE differs in that
 - it includes a regularization term $KL(q||p)$
 - it adds some Gaussian noise in the latent code

Performance

- Samples drawn from VAE trained on MNIST



Samples from training dataset



Samples generated by VAE

Improving the Quality of Generated Samples

- Improve the inference accuracy
 - ❖ Importance weighted auto-encoder (IWAE)
 - ❖ Normalizing flow
 - ❖ Inverse autoregressive flow
 - ❖ Implicit model
- Adversarial training
 - ❖ Various GANs...
- Energy-based models
- Diffusion models

Generated Examples



Samples from training dataset



Samples generated by models

