

# Variational Autoencoder

Deep Learning (ENGG\*6600\*01)

School of Engineering,  
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghoghj  
Summer 2023

## Variational Inference

# Variational Inference

- Consider a dataset  $\{\mathbf{x}_i\}_{i=1}^n$ . Assume that every data point  $\mathbf{x}_i \in \mathbb{R}^d$  is generated from a latent variable  $\mathbf{z}_i \in \mathbb{R}^p$ . This latent variable has a prior distribution  $\mathbb{P}(\mathbf{z}_i)$ . According to Bayes' rule, we have:

$$\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i) = \frac{\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) \mathbb{P}(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i)}. \quad (1)$$

- Let  $\mathbb{P}(\mathbf{z}_i)$  be an arbitrary distribution denoted by  $q(\mathbf{z}_i)$ . Suppose the parameter of conditional distribution of  $\mathbf{z}_i$  on  $\mathbf{x}_i$  is denoted by  $\theta$ ; hence,  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta)$ . Therefore, we can say:

$$\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta) = \frac{\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta) \mathbb{P}(\mathbf{z}_i | \theta)}{\mathbb{P}(\mathbf{x}_i | \theta)}. \quad (2)$$

# Variational Inference

- Consider the Kullback-Leibler (KL) divergence [1] between the prior probability of the latent variable and the posterior of the latent variable:

$$\begin{aligned}\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})) &\stackrel{(a)}{=} \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta}))) d\mathbf{z}_i \\&\stackrel{(2)}{=} \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})) - \log(\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta})) + \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta}))) d\mathbf{z}_i \\&\stackrel{(b)}{=} \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})) - \log(\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta}))) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})),\end{aligned}$$

where (a) is for definition of KL divergence and (b) is because  $\log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta}))$  is independent of  $\mathbf{z}_i$  and comes out of integral and  $\int d\mathbf{z}_i = 1$ .

- Hence:

$$\log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})). \quad (3)$$

# Variational Inference

- We found:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})).$$

- We define the **Evidence Lower Bound (ELBO)** as:

$$\mathcal{L}(q, \boldsymbol{\theta}) := -\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})). \quad (4)$$

So:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) + \mathcal{L}(q, \boldsymbol{\theta}).$$

- Therefore:

$$\mathcal{L}(q, \boldsymbol{\theta}) = \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) - \underbrace{\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}))}_{\geq 0}. \quad (5)$$

- As the second term is negative with its minus, the ELBO is a lower bound on the log likelihood of data:

$$\mathcal{L}(q, \boldsymbol{\theta}) \leq \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})). \quad (6)$$

The likelihood  $\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})$  is also referred to as the **evidence**.

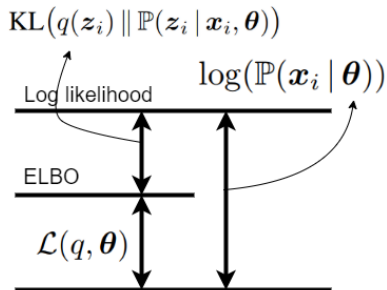
- Note that this lower bound gets tight when:

$$\begin{aligned} \mathcal{L}(q, \boldsymbol{\theta}) \approx \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) &\implies 0 \leq \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) \stackrel{\text{set}}{=} 0 \\ &\implies q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}). \end{aligned} \quad (7)$$

# Variational Inference

- We found:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) + \mathcal{L}(q, \boldsymbol{\theta}).$$



# Expectation Maximization in Variational Inference

- According to MLE, we want to maximize the log-likelihood of data. According to Eq. (6):

$$\mathcal{L}(q, \theta) \leq \log(\mathbb{P}(\mathbf{x}_i | \theta)),$$

maximizing the ELBO will also maximize the log-likelihood.

- The Eq. (6) holds for any prior distribution  $q$ . We want to find the best distribution to maximize the lower bound.
- Hence, EM for variational inference is performed iteratively as:

$$\text{E-step: } q^{(t)} := \arg \max_q \mathcal{L}(q, \theta^{(t-1)}), \quad (8)$$

$$\text{M-step: } \theta^{(t)} := \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta), \quad (9)$$

where  $t$  denotes the iteration index.

# Expectation Maximization in Variational Inference

- **E-step in EM for Variational Inference:** The E-step is:

$$\begin{aligned}\max_q \mathcal{L}(q, \theta^{(t-1)}) &\stackrel{(5)}{=} \max_q \log(\mathbb{P}(\mathbf{x}_i | \theta^{(t-1)})) + \max_q (-\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}))) \\ &= \max_q \log(\mathbb{P}(\mathbf{x}_i | \theta^{(t-1)})) + \min_q \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)})).\end{aligned}$$

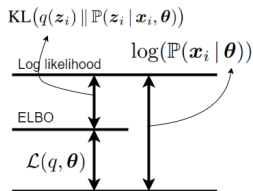
- The second term is always non-negative; hence, its minimum is zero:

$$\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)})) \stackrel{\text{set}}{=} 0 \implies q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}),$$

which was already found in Eq. (7). Thus, the E-step assigns:

$$q^{(t)}(\mathbf{z}_i) \leftarrow \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}). \quad (10)$$

- In other words, in the figure, it pushes the middle line toward the above line by maximizing the ELBO.





# Expectation Maximization in Variational Inference

- **M-step in EM for Variational Inference:** The M-step is:

$$\begin{aligned}\max_{\theta} \mathcal{L}(q^{(t)}, \theta) &\stackrel{(4)}{=} \max_{\theta} (-\text{KL}(q^{(t)}(z_i) \parallel \mathbb{P}(x_i, z_i \mid \theta))) \\ &\stackrel{(a)}{=} \max_{\theta} \left[ - \int q^{(t)}(z_i) \log\left(\frac{q^{(t)}(z_i)}{\mathbb{P}(x_i, z_i \mid \theta)}\right) dz_i \right] \\ &= \max_{\theta} \int q^{(t)}(z_i) \log(\mathbb{P}(x_i, z_i \mid \theta)) dz_i - \max_{\theta} \int q^{(t)}(z_i) \log(q^{(t)}(z_i)) dz_i,\end{aligned}$$

where (a) is for definition of KL divergence.

- The second term is constant w.r.t.  $\theta$ . Hence:

$$\begin{aligned}\max_{\theta} \mathcal{L}(q^{(t)}, \theta) &= \max_{\theta} \int q^{(t)}(z_i) \log(\mathbb{P}(x_i, z_i \mid \theta)) dz_i \\ &\stackrel{(a)}{=} \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(x_i, z_i \mid \theta)],\end{aligned}$$

where (a) is because of definition of expectation. Thus, the M-step assigns:

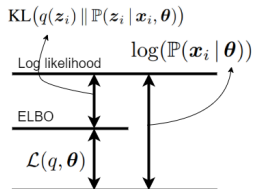
$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(x_i, z_i \mid \theta)]. \quad (11)$$

# Expectation Maximization in Variational Inference

- We found:

$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta)].$$

- In other words, in the figure, it pushes the above line higher.



- The E-step and M-step together somehow play a **game** where the E-step tries to reach the middle line (or the ELBO) to the log-likelihood and the M-step tries to increase the above line (or the log-likelihood). This procedure is done repeatedly so the two steps help each other improve to higher values.
- To summarize, the EM in variational inference is:

$$q^{(t)}(z_i) \leftarrow \mathbb{P}(z_i | \mathbf{x}_i, \theta^{(t-1)}), \quad (12)$$

$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta)]. \quad (13)$$

# Expectation Maximization in Variational Inference

- It is noteworthy that, in variational inference, sometimes, the parameter  $\theta$  is absorbed into the latent variable  $\mathbf{z}_i$ .
- According to the chain rule, we have:

$$\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i, \theta) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta) \mathbb{P}(\mathbf{z}_i | \theta) \mathbb{P}(\theta).$$

- Considering the term  $\mathbb{P}(\mathbf{z}_i | \theta) \mathbb{P}(\theta)$  as one probability term, we have:

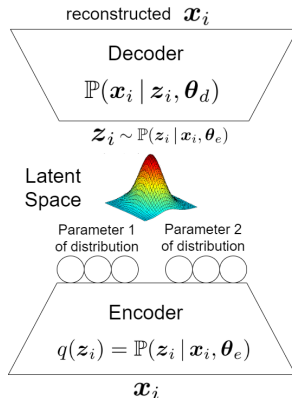
$$\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) \mathbb{P}(\mathbf{z}_i),$$

where the parameter  $\theta$  disappears because of absorption.

## **Variational Autoencoder**

# Variational Autoencoder

- **Variational Autoencoder (VAE)** (2014) [2] applies variational inference, i.e., maximizes the ELBO, but in an **autoencoder** setup and makes it differentiable for the **backpropagation** training [3].
- As this figure shows, VAE includes an **encoder** and a **decoder**, each of which can have several network layers. A **latent space** is learned between the encoder and decoder. The **latent variable  $z_i$**  is **sampled** from the latent space. The input of encoder in VAE is the **data point  $x_i$**  and the output of decoder in VAE is its **reconstruction  $x_i$** .



# Encoder of Variational Autoencoder

- The encoder of VAE models the distribution  $q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$  where the parameters of distribution  $\theta_e$  are the weights of encoder layers in VAE.
- The input and output of encoder are  $\mathbf{x}_i \in \mathbb{R}^d$  and  $\mathbf{z}_i \in \mathbb{R}^p$ , respectively.
- As the figure depicts, the output neurons of encoder are supposed to determine the parameters of the conditional distribution  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$ . If this conditional distribution has  **$m$  number of parameters**, we have  **$m$  sets of output neurons from the encoder**, denoted by  $\{\mathbf{e}_j\}_{j=1}^m$ . The dimensionality of these sets may differ depending on the size of the parameters.
- For example, let the latent space be  $p$ -dimensional, i.e.,  $\mathbf{z}_i \in \mathbb{R}^p$ . If the distribution  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$  is a **multivariate Gaussian distribution**, we have two sets of output neurons for encoder where one set has  $p$  neurons for the **mean** of this distribution  $\mu_{\mathbf{z}|\mathbf{x}} = \mathbf{e}_1 \in \mathbb{R}^p$  and the other set has  $(p \times p)$  neurons for the **covariance** of this distribution  $\Sigma_{\mathbf{z}|\mathbf{x}} =$  matrix form of  $\mathbf{e}_2 \in \mathbb{R}^{p \times p}$ . If the covariance matrix is **diagonal**, the second set has  $p$  neurons rather than  $(p \times p)$  neurons. In this case, we have  $\Sigma_{\mathbf{z}|\mathbf{x}} = \text{diag}(\mathbf{e}_2) \in \mathbb{R}^{p \times p}$ .
- Any distribution with any number of parameters can be chosen for  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$  but the **multivariate Gaussian with diagonal covariance** is very well-used:

$$q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) = \mathcal{N}(\mathbf{z}_i | \mu_{\mathbf{z}|\mathbf{x}}, \Sigma_{\mathbf{z}|\mathbf{x}}). \quad (14)$$

- Let the network **weights** for the output sets of encoder,  $\{\mathbf{e}_j\}_{j=1}^m$ , be denoted by  $\{\theta_{e,j}\}_{j=1}^m$ . As the input of encoder is  $\mathbf{x}_i$ , the  $j$ -th output set of encoder can be written as  $\mathbf{e}_j(\mathbf{x}_i, \theta_{e,j})$ . In the case of multivariate Gaussian distribution for the latent space, the parameters are  $\mu_{\mathbf{z}|\mathbf{x}} = \mathbf{e}_1(\mathbf{x}_i, \theta_{e,1})$  and  $\Sigma_{\mathbf{z}|\mathbf{x}} = \text{diag}(\mathbf{e}_2(\mathbf{x}_i, \theta_{e,2}))$ .

# Sampling the Latent Variable

- When the **data point  $x_i$  is fed as input** to the encoder, the parameters of the conditional distribution  $q(z_i)$  are obtained; hence, the **distribution of latent space, which is  $q(z_i)$ , is determined corresponding to the data point  $x_i$ .**
- Now, in the latent space, we **sample** the corresponding latent variable from the distribution of latent space:

$$z_i \sim q(z_i) = \mathbb{P}(z_i | x_i, \theta_e). \quad (15)$$

- This latent variable is **fed as input to the decoder** which is explained in the following.

# Decoder of Variational Autoencoder

- As the figure shows, the decoder of VAE models the conditional distribution  $\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}_d)$  where  $\boldsymbol{\theta}_d$  are the weights of decoder layers in VAE.
- The input and output of decoder are  $\mathbf{z}_i \in \mathbb{R}^p$  and  $\mathbf{x}_i \in \mathbb{R}^d$ , respectively. The output neurons of decoder are supposed to either generate the reconstructed data point or determine the parameters of the conditional distribution  $\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}_d)$ .
- The former is more common.
- In the latter case, if this conditional distribution has  $l$  number of parameters, we have  $l$  sets of output neurons from the decoder, denoted by  $\{\mathbf{d}_j\}_{j=1}^l$ . The dimensionality of these sets may differ depending the size of every parameters. The example of multivariate Gaussian distribution also can be mentioned for the decoder.
- Let the network **weights** for the output sets of decoder,  $\{\mathbf{d}_j\}_{j=1}^l$ , be denoted by  $\{\boldsymbol{\theta}_{d,j}\}_{j=1}^l$ . As the input of decoder is  $\mathbf{z}_i$ , the  $j$ -th output set of decoder can be written as  $\mathbf{d}_j(\mathbf{z}_i, \boldsymbol{\theta}_{d,j})$ .



# Training Variational Autoencoder with Expectation Maximization

- We use EM for training the VAE. Recall Eqs. (8) and (9) for EM in variational inference:

$$\text{E-step: } q^{(t)} := \arg \max_q \mathcal{L}(q, \theta^{(t-1)}),$$

$$\text{M-step: } \theta^{(t)} := \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta).$$

- Inspired by that, VAE uses EM for training where the ELBO is a function of encoder weights  $\theta_e$ , decoder weights  $\theta_d$ , and data point  $\mathbf{x}_i$ :

$$\text{E-step: } \theta_e^{(t)} := \arg \max_q \mathcal{L}(\theta_e, \theta_d^{(t-1)}, \mathbf{x}_i), \quad (16)$$

$$\text{M-step: } \theta_d^{(t)} := \arg \max_q \mathcal{L}(\theta_e^{(t)}, \theta_d, \mathbf{x}_i). \quad (17)$$

# Training Variational Autoencoder with Expectation Maximization

- We had:

$$\text{E-step: } \theta_e^{(t)} := \arg \max_q \mathcal{L}(\theta_e, \theta_d^{(t-1)}, \mathbf{x}_i),$$

$$\text{M-step: } \theta_d^{(t)} := \arg \max_q \mathcal{L}(\theta_e^{(t)}, \theta_d, \mathbf{x}_i).$$

- We can simplify this iterative optimization algorithm by **alternating optimization** [4] where we take a step of gradient ascent optimization in every iteration. We consider mini-batch stochastic gradient ascent and take training data in batches where  $b$  denotes the mini-batch size. Hence, the optimization is:

$$\text{E-step: } \theta_e^{(t)} := \theta_e^{(t-1)} + \eta_e \frac{\partial \sum_{i=1}^b \mathcal{L}(\theta_e, \theta_d^{(t-1)}, \mathbf{x}_i)}{\partial \theta_e}, \quad (18)$$

$$\text{M-step: } \theta_d^{(t)} := \theta_d^{(t-1)} + \eta_d \frac{\partial \sum_{i=1}^b \mathcal{L}(\theta_e^{(t)}, \theta_d, \mathbf{x}_i)}{\partial \theta_d}, \quad (19)$$

where  $\eta_e$  and  $\eta_d$  are the learning rates for  $\theta_e$  and  $\theta_d$ , respectively.

# Training Variational Autoencoder with Expectation Maximization

- Eqs. (4) and (12) were:

$$\begin{aligned}\mathcal{L}(q, \theta) &:= -\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \theta)), \\ q^{(t)}(\mathbf{z}_i) &\leftarrow \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \theta^{(t-1)}).\end{aligned}$$

- The ELBO is simplified as:

$$\begin{aligned}\sum_{i=1}^b \mathcal{L}(q, \theta) &\stackrel{(4)}{=} -\sum_{i=1}^b \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \theta_d)) \\ &\stackrel{(12)}{=} -\sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \theta_e) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \theta_d)).\end{aligned}\tag{20}$$

- Note that the parameter of  $\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \theta_d)$  is  $\theta_d$  because  $\mathbf{z}_i$  is generated after the encoder and before the decoder.
- There are different ways for **approximating the KL divergence** in Eq. (20) [5, 6]. We can simplify the ELBO in at least two different ways which are explained in the following.

# Simplification Type 1

- We continue the simplification of ELBO:

$$\begin{aligned}\sum_{i=1}^b \mathcal{L}(q, \theta) &= - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)) \\ &= - \sum_{i=1}^b \mathbb{E}_{\mathbf{z}_i \sim q^{(t-1)}(\mathbf{z}_i)} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)} \right) \right] \\ &= - \sum_{i=1}^b \mathbb{E}_{\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)} \right) \right].\end{aligned}\tag{21}$$

- This expectation can be approximated using **Monte Carlo approximation** [7] where we draw  $\ell$  samples  $\{\mathbf{z}_{i,j}\}_{j=1}^{\ell}$ , corresponding to the  $i$ -th data point, from the conditional distribution distribution as:

$$\mathbf{z}_{i,j} \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e), \quad \forall j \in \{1, \dots, \ell\}.\tag{22}$$

- Monte Carlo approximation [7], in general, approximates expectation as:

$$\mathbb{E}_{\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} [f(\mathbf{z}_i)] \approx \frac{1}{\ell} \sum_{j=1}^{\ell} f(\mathbf{z}_{i,j}),\tag{23}$$

where  $f(\mathbf{z}_i)$  is a function of  $\mathbf{z}_i$ .

# Simplification Type 1

- We had:

$$\sum_{i=1}^b \mathcal{L}(q, \theta) = - \sum_{i=1}^b \mathbb{E}_{\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)} \right) \right].$$

- Here, the approximation is:

$$\begin{aligned} \sum_{i=1}^b \mathcal{L}(q, \theta) &\approx \sum_{i=1}^b \tilde{\mathcal{L}}(q, \theta) \\ &= - \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log \left( \frac{\mathbb{P}(\mathbf{z}_{i,j} | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_{i,j} | \theta_d)} \right) \\ &= \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \left[ \log (\mathbb{P}(\mathbf{x}_i, \mathbf{z}_{i,j} | \theta_d)) - \log (\mathbb{P}(\mathbf{z}_{i,j} | \mathbf{x}_i, \theta_e)) \right]. \end{aligned} \quad (24)$$

## Simplification Type 2

- We can simplify the ELBO using another approach:

$$\begin{aligned}\sum_{i=1}^b \mathcal{L}(q, \theta) &= - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \| \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)) \\&= - \sum_{i=1}^b \int \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)} \right) d\mathbf{z}_i \\&= - \sum_{i=1}^b \int \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d) \mathbb{P}(\mathbf{z}_i)} \right) d\mathbf{z}_i \\&= - \sum_{i=1}^b \int \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{z}_i)} \right) d\mathbf{z}_i \\&\quad + \sum_{i=1}^b \int \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d)) d\mathbf{z}_i \\&= - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \| \mathbb{P}(\mathbf{z}_i)) \\&\quad + \sum_{i=1}^b \mathbb{E}_{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d)) \right].\end{aligned}\tag{25}$$

## Simplification Type 2

- We found:

$$\sum_{i=1}^b \mathcal{L}(q, \theta) = - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \| \mathbb{P}(\mathbf{z}_i)) + \sum_{i=1}^b \mathbb{E}_{\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d)) \right].$$

- The second term in the above equation can be estimated using **Monte Carlo approximation** [7] where we draw  $\ell$  samples  $\{\mathbf{z}_{i,j}\}_{j=1}^{\ell}$  from  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$ :

$$\begin{aligned} \sum_{i=1}^b \mathcal{L}(q, \theta) &\approx \sum_{i=1}^b \tilde{\mathcal{L}}(q, \theta) \\ &= - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \| \mathbb{P}(\mathbf{z}_i)) + \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_{i,j}, \theta_d)). \end{aligned} \quad (26)$$

## Simplification Type 2

- We had:

$$\sum_{i=1}^b \mathcal{L}(q, \theta) \approx - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \parallel \mathbb{P}(\mathbf{z}_i)) + \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log(\mathbb{P}(\mathbf{x}_i | \mathbf{z}_{i,j}, \theta_d)).$$

- The first term in the above equation can be **converted to expectation** and then computed using **Monte Monte Carlo approximation** [7] again, where we draw  $\ell$  samples  $\{\mathbf{z}_{i,j}\}_{j=1}^{\ell}$  from  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$ :

$$\begin{aligned} \sum_{i=1}^b \mathcal{L}(q, \theta) &\approx \sum_{i=1}^b \tilde{\mathcal{L}}(q, \theta) \\ &= - \sum_{i=1}^b \mathbb{E}_{\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{z}_i)} \right) \right] + \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log(\mathbb{P}(\mathbf{x}_i | \mathbf{z}_{i,j}, \theta_d)) \\ &\approx - \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log(\mathbb{P}(\mathbf{z}_{i,j} | \mathbf{x}_i, \theta_e)) - \log(\mathbb{P}(\mathbf{z}_{i,j})) + \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log(\mathbb{P}(\mathbf{x}_i | \mathbf{z}_{i,j}, \theta_d)). \end{aligned} \tag{27}$$

- In case we have some families of distributions, such as **Gaussian** distributions, for  $\mathbb{P}(\mathbf{z}_{i,j} | \mathbf{x}_i, \theta_e)$  and  $\mathbb{P}(\mathbf{z}_{i,j})$ , the first term in Eq. (26) can be **computed analytically**. In the following, we simply Eq. (26) further for Gaussian distributions.



# Simplification Type 2 for Special Case of Gaussian Distributions

- We can compute the KL divergence in the **first term** of Eq. (26) analytically for **univariate or multivariate Gaussian** distributions. For this, we need two following lemmas (see our tutorial paper [8] for proof).

## Lemma

*The KL divergence between two univariate Gaussian distributions  $p_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $p_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  is:*

$$KL(p_1 \| p_2) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}. \quad (28)$$

## Lemma

*The KL divergence between two multivariate Gaussian distributions  $p_1 \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $p_2 \sim \mathcal{N}(\mu_2, \Sigma_2)$  with dimensionality  $p$  is:*

$$KL(p_1 \| p_2) = \frac{1}{2} \left( \log\left(\frac{|\Sigma_2|}{|\Sigma_1|}\right) - p + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^\top \Sigma_2^{-1}(\mu_2 - \mu_1) \right). \quad (29)$$

# Simplification Type 2 for Special Case of Gaussian Distributions

- Consider the case in which we have:

$$\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}_e) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}), \quad (30)$$

$$\mathbb{P}(\mathbf{z}_i) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}}), \quad (31)$$

where  $\mathbf{z}_i \in \mathbb{R}^p$ . Note that the parameters  $\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}$  and  $\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}$  are trained in neural network while the parameters  $\mathbb{P}(\mathbf{z}_{i,j})$  can be set to  $\boldsymbol{\mu}_{\mathbf{z}} = \mathbf{0}$  and  $\boldsymbol{\Sigma}_{\mathbf{z}} = \mathbf{I}$  (inspired by the prior distribution of  $\mathbf{z}$  in factor analysis).

- According to Lemma 2, the approximation of ELBO, i.e. Eq. (26), can be simplified to:

$$\begin{aligned} \sum_{i=1}^b \mathcal{L}(q, \boldsymbol{\theta}) &\approx \sum_{i=1}^b \tilde{\mathcal{L}}(q, \boldsymbol{\theta}) \\ &= - \sum_{i=1}^b \frac{1}{2} \left( \log \left( \frac{|\boldsymbol{\Sigma}_{\mathbf{z}}|}{|\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}|} \right) - p + \text{tr}(\boldsymbol{\Sigma}_{\mathbf{z}}^{-1} \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}) + (\boldsymbol{\mu}_{\mathbf{z}} - \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}})^\top \boldsymbol{\Sigma}_{\mathbf{z}}^{-1} (\boldsymbol{\mu}_{\mathbf{z}} - \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}) \right) \\ &\quad + \sum_{i=1}^b \frac{1}{\ell} \sum_{j=1}^{\ell} \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_{i,j}, \boldsymbol{\theta}_d)). \end{aligned} \quad (32)$$

# Training Variational Autoencoder with Approximations

- We can train VAE with EM, where Monte Carlo approximations are applied to ELBO. The Eqs. (18) and (19):

$$\text{E-step: } \boldsymbol{\theta}_e^{(t)} := \boldsymbol{\theta}_e^{(t-1)} + \eta_e \frac{\partial \sum_{i=1}^b \mathcal{L}(\boldsymbol{\theta}_e, \boldsymbol{\theta}_d^{(t-1)}, \mathbf{x}_i)}{\partial \boldsymbol{\theta}_e},$$

$$\text{M-step: } \boldsymbol{\theta}_d^{(t)} := \boldsymbol{\theta}_d^{(t-1)} + \eta_d \frac{\partial \sum_{i=1}^b \mathcal{L}(\boldsymbol{\theta}_e^{(t)}, \boldsymbol{\theta}_d, \mathbf{x}_i)}{\partial \boldsymbol{\theta}_d},$$

are replaced by the following equations:

$$\text{E-step: } \boldsymbol{\theta}_e^{(t)} := \boldsymbol{\theta}_e^{(t-1)} + \eta_e \frac{\partial \sum_{i=1}^b \tilde{\mathcal{L}}(\boldsymbol{\theta}_e, \boldsymbol{\theta}_d^{(t-1)}, \mathbf{x}_i)}{\partial \boldsymbol{\theta}_e}, \quad (33)$$

$$\text{M-step: } \boldsymbol{\theta}_d^{(t)} := \boldsymbol{\theta}_d^{(t-1)} + \eta_d \frac{\partial \sum_{i=1}^b \tilde{\mathcal{L}}(\boldsymbol{\theta}_e^{(t)}, \boldsymbol{\theta}_d, \mathbf{x}_i)}{\partial \boldsymbol{\theta}_d}, \quad (34)$$

where the approximated ELBO was introduced in previous sections.

# The Reparameterization Trick

- Sampling the  $\ell$  samples for the latent variables, i.e. Eq. (15):

$$\mathbf{z}_i \sim q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e),$$

**blocks the gradient flow** because computing the **derivatives** through  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$  by chain rule gives a **high variance estimate of gradient**.

- In order to overcome this problem, we use the **reparameterization technique** (2014) [2, 9, 10]. In this technique, instead of sampling  $\mathbf{z}_i \sim \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$ , we assume  $\mathbf{z}_i$  is a random variable but is a **deterministic function of another random variable**  $\epsilon_i$  as follows:

$$\mathbf{z}_i = g(\epsilon_i, \mathbf{x}_i, \theta_e), \quad (35)$$

where  $\epsilon_i$  is a stochastic variable sampled from a distribution as:

$$\epsilon_i \sim \mathbb{P}(\epsilon). \quad (36)$$

# The Reparameterization Trick

- The Eqs. (21) and (25):

$$\sum_{i=1}^b \mathcal{L}(q, \theta) = - \sum_{i=1}^b \mathbb{E}_{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta_d)} \right) \right],$$

$$\sum_{i=1}^b \mathcal{L}(q, \theta) = - \sum_{i=1}^b \text{KL}(\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e) \| \mathbb{P}(\mathbf{z}_i)) + \sum_{i=1}^b \mathbb{E}_{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)} \left[ \log (\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d)) \right],$$

both contain an expectation of a function  $f(\mathbf{z}_i)$ . Using this technique, this expectation is replaced as:

$$\mathbb{E}_{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}[f(\mathbf{z}_i)] \rightarrow \mathbb{E}_{\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)}[f(g(\epsilon_i, \mathbf{x}_i, \theta_e))]. \quad (37)$$

- Using the reparameterization technique, the **encoder**, which implemented  $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta_e)$ , is replaced by  $g(\epsilon_i, \mathbf{x}_i, \theta_e)$  where in the latent space between encoder and decoder, we have  $\epsilon_i \sim \mathbb{P}(\epsilon)$  and  $\mathbf{z}_i = g(\epsilon_i, \mathbf{x}_i, \theta_e)$ .
- A simple example for the reparameterization technique is when  $\mathbf{z}_i$  and  $\epsilon_i$  are univariate Gaussian variables:

$$\mathbf{z}_i \sim \mathcal{N}(\mu, \sigma^2),$$

$$\epsilon_i \sim \mathcal{N}(0, 1),$$

$$\mathbf{z}_i = g(\epsilon_i) = \mu + \sigma \epsilon_i.$$

- For some more advanced reparameterization techniques, the reader can refer to [11].

# Training Variational Autoencoder with Backpropagation

- In practice, VAE is trained by **backpropagation** [9] where the backpropagation algorithm [3] is used for training the weights of network.
- Recall that in **training VAE with EM**, the **encoder and decoder are trained separately using the E-step and the M-step of EM**, respectively.
- However, in **training VAE with backpropagation**, the **whole network is trained together and not in separate steps**.
- Suppose the whole weights of VAE are denoted by  $\theta := \{\theta_e, \theta_d\}$ . Backpropagation trains VAE using the **mini-batch stochastic gradient descent** with the **negative ELBO**,  $\sum_{i=1}^b -\tilde{\mathcal{L}}(\theta, \mathbf{x}_i)$ , as the loss function:

$$\theta^{(t)} := \theta^{(t-1)} - \eta \frac{\partial \sum_{i=1}^b -\tilde{\mathcal{L}}(\theta, \mathbf{x}_i)}{\partial \theta}, \quad (38)$$

where  $\eta$  is the learning rate. Note that we are **minimizing** here because neural networks usually minimize the loss function.

# The Test Phase in Variational Autoencoder

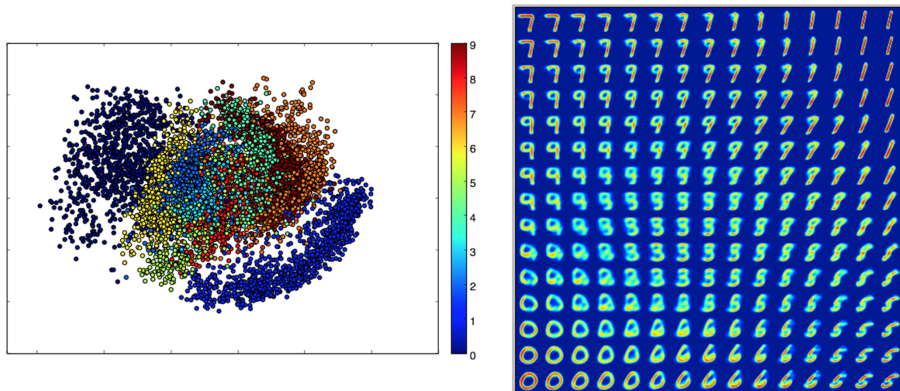
- In the test phase, we **feed the test data point  $x_i$  to the encoder** to determine the **parameters of the conditional distribution of latent space**, i.e.,  $\mathbb{P}(z_i | x_i, \theta_e)$ .
- Then, from this distribution, we **sample** the latent variable  $z_i$  from the latent space and generate the corresponding **reconstructed data point  $x_i$**  by the **decoder**.
- As you see, VAE is a **generative model** which generates data points [12].

# Blurry Images Generated by VAE

- One of the problems of VAE is generating **blurry images** when data points are images. This blurry artifact may be because of several following reasons:
  - ▶ sampling for the **Monte Carlo approximations**
  - ▶ **lower bound approximation** by ELBO
  - ▶ **restrictions on the family of distributions** where usually simple Gaussian distributions are used.
- Note that **generative adversarial networks** [13] usually generate clearer images; therefore, some works have combined variational and adversarial inferences [14] for using the advantages of both models.



# Simulation on MNIST Digit Dataset



Credit of image: <https://blog.keras.io/building-autoencoders-in-keras.html>

# Acknowledgment

- Some slides are based on our tutorial paper: “Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey” [8]
- Some slides of this slide deck are inspired by teachings of deep learning course at the Carnegie Mellon University (you can see their YouTube channel).
- Variational autoencoder in Keras:
  - ▶ <https://blog.keras.io/building-autoencoders-in-keras.html>
  - ▶ <https://keras.io/examples/generative/vae/>

# References

- [1] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [2] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *International Conference on Learning Representations*, 2014.
- [3] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [4] P. Jain and P. Kar, "Non-convex optimization for machine learning," *Foundations and Trends® in Machine Learning*, vol. 10, no. 3-4, pp. 142–336, 2017.
- [5] J. R. Hershey and P. A. Olsen, "Approximating the Kullback Leibler divergence between Gaussian mixture models," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing*, vol. 4, pp. IV–317, IEEE, 2007.
- [6] J. Duchi, "Derivations for linear algebra and optimization," tech. rep., Berkeley, California, 2007.
- [7] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karay, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," *arXiv preprint arXiv:2011.00901*, 2020.
- [8] B. Ghojogh, A. Ghodsi, F. Karay, and M. Crowley, "Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey," *arXiv preprint arXiv:2101.00734*, 2021.

# References (cont.)

- [9] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," in *International Conference on Machine Learning*, 2014.
- [10] M. Titsias and M. Lázaro-Gredilla, "Doubly stochastic variational Bayes for non-conjugate inference," in *International conference on machine learning*, pp. 1971–1979, 2014.
- [11] M. Figurnov, S. Mohamed, and A. Mnih, "Implicit reparameterization gradients," *Advances in Neural Information Processing Systems*, vol. 31, pp. 441–452, 2018.
- [12] A. Y. Ng and M. I. Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes," in *Advances in neural information processing systems*, pp. 841–848, 2002.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- [14] L. Mescheder, S. Nowozin, and A. Geiger, "Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks," in *International Conference on Machine Learning*, 2017.