Variational Autoencoder

Deep Learning (ENGG\*6600\*01)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Summer 2023

(aton f
 (2) observed

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

$$\mathbb{P}(\boldsymbol{z}_i \mid \boldsymbol{x}_i) = \frac{\mathbb{P}(\boldsymbol{x}_i \mid \boldsymbol{z}_i) \mathbb{P}(\boldsymbol{z}_i)}{\mathbb{P}(\boldsymbol{x}_i)}.$$
 (1)

Let P(z<sub>i</sub>) be an arbitrary distribution denoted by q(z<sub>i</sub>). Suppose the parameter of conditional distribution of z<sub>i</sub> on x<sub>i</sub> is denoted by θ; hence, P(z<sub>i</sub> | x<sub>i</sub>) = P(z<sub>i</sub> | x<sub>i</sub>, θ). Therefore, we can say:

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

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

• We found:

4

$$\boldsymbol{\mathsf{K}} \quad \log(\mathbb{P}(\boldsymbol{x}_i \mid \boldsymbol{\theta})) = \mathsf{KL}\big(q(\boldsymbol{z}_i) \parallel \mathbb{P}(\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{\theta})\big) - \mathsf{KL}\big(q(\boldsymbol{z}_i) \parallel \mathbb{P}(\boldsymbol{x}_i, \boldsymbol{z}_i \mid \boldsymbol{\theta})\big).$$

• We define the Evidence Lower Bound (ELBO) as:

 $\mathcal{L}(q,\theta) =$ 

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

M

×

So:

$$\log(\mathbb{P}(\boldsymbol{x}_i \mid \boldsymbol{\theta})) = \mathsf{KL}(q(\boldsymbol{z}_i) \parallel \mathbb{P}(\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{\theta})) + \mathcal{L}(q, \boldsymbol{\theta})$$

 $\log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) - \left(\mathsf{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta}))\right).$ 

Therefore:

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

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

>0

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

Note that this lower bound gets tight when:

$$\underbrace{\mathcal{L}(q,\theta) \bigotimes \log(\mathbb{P}(\mathbf{x}_i \mid \theta))}_{\Longrightarrow} \Longrightarrow \underbrace{0 \leq}_{\substack{\{\mathbf{z}_i\} = \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \theta).}} \underbrace{\mathsf{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \theta))}_{(z_i \mid z_i, \theta)} \underbrace{\overset{\text{set}}{=} 0}_{(z_i) = \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \theta).}$$
(7)

(5)

We found:



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

$$\mathcal{L}(q, \theta) \leq \widehat{\log(\mathbb{P}(x_i \mid \theta))},$$

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:

E-step: 
$$q^{(t)} := \arg \max_{q} (\mathcal{L}(q, \theta^{(t-1)})),$$
 (8)  
M-step:  $\theta^{(t)} := \arg \max_{\theta} (\mathcal{L}(q^{(t)}, \theta)),$  (9)

where t denotes the iteration index.

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

$$\max_{q} \mathcal{L}(q, \theta^{(t-1)}) \stackrel{(5)}{=} \max_{q} \log(\mathbb{P}(\mathbf{x}_{i} | \theta^{(t-1)})) + \max_{q} \mathcal{O}(\mathbb{K}L(q(\mathbf{z}_{i}) || \mathbb{P}(\mathbf{z}_{i} | \mathbf{x}_{i}, \theta^{(t-1)}))))$$
$$= \max_{q} \log(\mathbb{P}(\mathbf{x}_{i} | \theta^{(t-1)})) + \min_{q} \mathcal{K}L(q(\mathbf{z}_{i}) || \mathbb{P}(\mathbf{z}_{i} | \mathbf{x}_{i}, \theta^{(t-1)}))).$$

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

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

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

$$\underbrace{q^{(t)}(\boldsymbol{z}_i) \leftarrow \mathbb{P}(\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{\theta}^{(t-1)}).}_{(10)}$$

 In other words, in the figure, it <u>pushes the middle line</u> toward the above line by maximizing the <u>ELBO</u>.

$$\mathsf{KL}(q(\boldsymbol{z}_i) \parallel \mathbb{P}(\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{\theta}))$$

$$\mathsf{Log}_{\mathsf{likelihood}} \log(\mathbb{P}(\boldsymbol{x}_i \mid \boldsymbol{\theta}))$$

$$\mathsf{L}(q, \boldsymbol{\theta})$$

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

$$\max_{\theta} \mathcal{L}(q^{(t)}, \theta) \stackrel{(4)}{=} \max_{\theta} \left( -\operatorname{KL}(q^{(t)}(z_i) || \mathbb{P}(x_i, z_i | \theta)) \right)$$

$$\stackrel{(a)}{=} \max_{\theta} \left[ -\int_{z_i} q^{(t)}(z_i) \log(\frac{q^{(t)}(z_i)}{\mathbb{P}(x_i, z_i | \theta)}) dz_i \right]$$

$$= \max_{\theta} \int_{z_i} q^{(t)}(z_i) \log(\mathbb{P}(x_i, z_i | \theta)) dz_i - \max_{\theta} \int_{z_i} q^{(t)}(z_i) \log(q^{(t)}(z_i)) dz_i,$$

where (a) is for definition of KL divergence.

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

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{q}^{(t)}, \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}} \int \overline{\boldsymbol{q}^{(t)}(\boldsymbol{z}_i)} \log(\mathbb{P}(\boldsymbol{x}_i, \boldsymbol{z}_i \mid \boldsymbol{\theta})) d\boldsymbol{z}_i$$

$$\stackrel{(a)}{=} \max_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{q}^{(t)}(\boldsymbol{z}_i)} [\log \mathbb{P}(\boldsymbol{x}_i, \boldsymbol{z}_i \mid \boldsymbol{\theta})],$$

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

$$\underbrace{\begin{pmatrix} \boldsymbol{\theta}^{(t)} \leftarrow \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\sim q^{(t)}(\boldsymbol{z}_i)} \big[ \log \mathbb{P}(\boldsymbol{x}_i, \boldsymbol{z}_i \mid \boldsymbol{\theta}) \big].$$
(11)

We found:

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

• 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:

$$\sum q^{(t)}(\boldsymbol{z}_i) \leftarrow \mathbb{P}(\boldsymbol{z}_i \,|\, \boldsymbol{x}_i, \boldsymbol{\theta}^{(t-1)}), \tag{12}$$

$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(\boldsymbol{z}_i)} \big[ \log \mathbb{P}(\boldsymbol{x}_i, \boldsymbol{z}_i \mid \theta) \big].$$
(13)

 It is noteworthy that, in variational inference, sometimes, the parameter θ is absorbed into the latent variable z<sub>i</sub>.



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

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

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 <u>encoder</u> and a <u>decoder</u>, each of which can have several network layers. A latent space is learned <u>between the encoder and decoder</u>. The <u>latent variable z<sub>i</sub> is sampled from the latent space</u>. The input of encoder in VAE is the data point x<sub>i</sub> and the output of decoder in VAE is its reconstruction x<sub>i</sub>.



#### Encoder of Variational Autoencoder

- The encoder of VAE models the distribution  $q(z_i) = \mathbb{P}(z_i | x_i, \theta_e)$  where the parameters of distribution  $\theta_e$  are the weights of encoder layers in VAE.
- The input and output of encoder are  $x_i \in \mathbb{R}^d$  and  $z_i \in \mathbb{R}^p$ , respectively.
- As the figure depicts, the output neurons of encoder are supposed to determine the parameters of the conditional distribution  $\mathbb{P}(z_i | x_i, \theta_e)$ . If this conditional distribution has <u>*m* number of parameters</u>, we have <u>*m* sets of output neurons from the encoder</u>, denoted by  $\{e_i\}_{j=1}^m$ . 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.,  $z_i \in \mathbb{R}^p$ . If the distribution  $\mathbb{P}(z_i | x_i, \theta_e)$  is a <u>multivariate Gaussian distribution</u>, we have two sets of output neurons for encoder where one set has *p* neurons for the <u>mean</u> of this distribution  $\mu_{z|x} = e_1 \in \mathbb{R}^p$  and the other set has  $(p \times p)$  neurons for the <u>covariance</u> of this distribution  $\Sigma_{z|x} = \text{matrix}$  form of  $e_2 \in \mathbb{R}^{p \times p}$ . If the covariance matrix is <u>diagonal</u>, the second set has *p* neurons rather than  $(p \times p)$  neurons. In this case, we have  $\Sigma_{z|x} = \text{diag}(e_2) \in \mathbb{R}^{p \times p}$ .
- Any distribution with any number of parameters can be chosen for P(z<sub>i</sub> | x<sub>i</sub>, θ<sub>e</sub>) but the multivariate Gaussian with diagonal covariance is very well-used:

• Let the network weights for the output sets of encoder,  $\{e_j\}_{j=1}^m$ , be denoted by  $\{\theta_{e,j}\}_{j=1}^m$ . As the input of encoder is  $x_i$ , the *j*-th output set of encoder can be written as  $e_j(x_i, \theta_{e,j})$ . In the case of multivariate Gaussian distribution for the latent space, the parameters are  $\mu_{z|x} = e_1(x_i, \theta_{e,1})$  and  $\Sigma_{z|x} = \text{diag}(e_2(x_i, \theta_{e,2}))$ .

# 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:

$$\mathbf{z}_i \bigotimes \mathbf{q}(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta}_e).$$
(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  $\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta_d)$  where  $\theta_d$  are the weights of decoder layers in VAE.
- The input and output of decoder are z<sub>i</sub> ∈ ℝ<sup>p</sup> and x<sub>i</sub> ∈ ℝ<sup>d</sup>, respectively. The output neurons of decoder are supposed to either generate the reconstructed data point or determine the parameters of the conditional distribution P(x<sub>i</sub> | z<sub>i</sub>, θ<sub>d</sub>).
- The former is more common.
- In the latter case, if this conditional distribution has *I* number of parameters, we have *I* sets of output neurons from the decoder, denoted by {*d*<sub>j</sub>}<sub>j=1</sub>. 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,  $\{d_j\}_{j=1}^l$ , be denoted by  $\{\theta_{d,j}\}_{j=1}^l$ . As the input of decoder is  $z_i$ , the *j*-th output set of decoder can be written as  $d_j(z_i, \theta_{d,j})$ .

# Training Variational Autoencoder with Expectation Maximization

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

$$\begin{array}{ll} \overbrace{\text{K-step:}}^{\text{E-step:}} & q^{(t)} := \arg \max_{q} & \mathcal{L}(q, \theta^{(t-1)}), \\ & & & \\ \hline \text{M-step:} & \theta^{(t)} := \arg \max_{\theta} & \mathcal{L}(q^{(t)}, \theta). \end{array}$$

 Inspired by that, VAE uses EM for training where the ELBO is a function of encoder weights θ<sub>e</sub>, decoder weights θ<sub>d</sub>, and data point x<sub>i</sub>:

$$\mathcal{O}(\mathfrak{p}) \xrightarrow{} \mathsf{E}\text{-step:} \quad \left( \stackrel{\scriptstyle (e)}{\theta_e} \right) := \arg \max_q \quad \mathcal{L}(\theta_e, \theta_d^{(t-1)}, \mathbf{x}_i), \tag{16}$$

$$\begin{array}{c} \begin{array}{c} & \\ & \\ \end{array} \end{array} \overset{}{\longrightarrow} M \text{-step:} \qquad \\ \begin{array}{c} & \\ & \\ \end{array} \overset{}{\longrightarrow} \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ & \\ \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ & \\ \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ \end{array} \overset{}{\longrightarrow} \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ & \\ \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ & \\ \end{array} \overset{}{\longrightarrow} \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ \end{array} \overset{}{\end{array} \end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ \end{array} \overset{}{\longrightarrow} \end{array} \overset{}{\end{array} \overset{}{\end{array}} \overset{}{\end{array} \overset{}{\longrightarrow} \begin{array}{c} \\ \end{array} \overset{}{\end{array} \end{array} \overset{}{\end{array} \end{array} \overset{}{\end{array} \end{array} \overset{}{\longrightarrow} \begin{array}{c} \end{array} \overset{}{\end{array} \end{array} \overset{}{} \end{array} \overset{}{\end{array} \end{array}$$

# Training Variational Autoencoder with Expectation Maximization

We had:

$$\begin{cases} \mathsf{E}\text{-step:} \quad \boldsymbol{\theta}_{e}^{(t)} := \arg (\max_{q} \mathcal{L}(\boldsymbol{\theta}_{e}, \boldsymbol{\theta}_{d}^{(t-1)}, \mathbf{x}_{i}), \\ \mathsf{M}\text{-step:} \quad \boldsymbol{\theta}_{d}^{(t)} := \arg (\max_{q} \mathcal{L}(\boldsymbol{\theta}_{e}^{(t)}, \boldsymbol{\theta}_{d}, \mathbf{x}_{i}). \end{cases}$$

 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:

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:

• The ELBO is simplified as:

$$\sum_{i=1}^{b} \mathcal{L}(q,\theta) \stackrel{(4)}{=} \bigoplus_{i=1}^{b} \operatorname{KL}(q(z_{i}) || \mathbb{P}(\mathbf{x}_{i}, z_{i} | \theta_{d}))$$

$$\stackrel{(12)}{=} -\sum_{i=1}^{b} \operatorname{KL}(\mathbb{P}(z_{i} | \mathbf{x}_{i}, \theta_{e}) || \mathbb{P}(\mathbf{x}_{i}, z_{i} | \theta_{d})).$$
(20)

- Note that the parameter of  $\mathbb{P}(x_i, z_i | \theta_d)$  is  $\theta_d$  because  $z_i$  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.

• We continue the simplification of ELBO:

$$\sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \mathsf{KL}(\mathbb{P}(\mathbf{z}_{i} | \mathbf{x}_{i}, \theta_{e}) || \mathbb{P}(\mathbf{x}_{i}, \mathbf{z}_{i} | \theta_{d}))$$

$$= -\sum_{i=1}^{b} \mathbb{E}_{q(t-1)(\mathbf{z}_{i})} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_{i} | \mathbf{x}_{i}, \theta_{e})}{\mathbb{P}(\mathbf{x}_{i}, \mathbf{z}_{i} | \theta_{d})} \right) \right]$$

$$= -\sum_{i=1}^{b} \mathbb{E}_{q(t-1)(\mathbf{z}_{i})} \left[ \log \left( \frac{\mathbb{P}(\mathbf{z}_{i} | \mathbf{x}_{i}, \theta_{e})}{\mathbb{P}(\mathbf{x}_{i}, \mathbf{z}_{i} | \theta_{d})} \right) \right].$$
(21)

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

$$\underbrace{\mathbf{z}_{i,j} \sim \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta}_e)}_{\mathbf{z}_i \in \{1, \dots, \ell\}}, \quad \forall j \in \{1, \dots, \ell\}.$$
(22)

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

$$\mathbb{E}_{\sim \mathbb{P}(\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \boldsymbol{\theta}_{e})}[f(\boldsymbol{z}_{i})] \approx \frac{1}{\ell} \sum_{j=1}^{\ell} f(\boldsymbol{z}_{i,j}),$$
(23)

where  $f(z_i)$  is a function of  $z_i$ .



• We had:

$$\bigstar \quad \sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \mathbb{E}_{\sim \mathbb{P}(z_i \mid x_i, \theta_e)} \Big[ \log \big( \frac{\mathbb{P}(z_i \mid x_i, \theta_e)}{\mathbb{P}(x_i, z_i \mid \theta_d)} \big) \Big]$$

• Here, the approximation is:

$$\sum_{i=1}^{b} \mathcal{L}(q, \theta) \approx \sum_{i=1}^{b} \widetilde{\mathcal{L}}(q, \theta)$$
$$= -\sum_{i=1}^{b} \left\{ \frac{1}{\ell} \sum_{j=1}^{\ell} \log\left(\frac{\mathbb{P}(\boldsymbol{z}_{i,j} \mid \boldsymbol{x}_{i}, \theta_{e})}{\mathbb{P}(\boldsymbol{x}_{i}, \boldsymbol{z}_{i,j} \mid \theta_{d})}\right) \right\}$$
$$= \sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \left[ \log\left(\mathbb{P}(\boldsymbol{x}_{i}, \boldsymbol{z}_{i,j} \mid \theta_{d})\right) - \log\left(\mathbb{P}(\boldsymbol{z}_{i,j} \mid \boldsymbol{x}_{i}, \theta_{e})\right) \right].$$
(24)

• We can simplify the ELBO using another approach:

$$\sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \operatorname{KL}\left(\mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \mid \mathbb{P}(x_{i}, z_{i} \mid \theta_{d})\right)$$

$$= -\sum_{i=1}^{b} \int \mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \log\left(\frac{\mathbb{P}(z_{i} \mid x_{i}, \theta_{e})}{\mathbb{P}(x_{i} \mid z_{i}, \theta_{d})}\right) dz_{i}$$

$$= -\sum_{i=1}^{b} \int \mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \log\left(\frac{\mathbb{P}(z_{i} \mid x_{i}, \theta_{e})}{\mathbb{P}(z_{i})}\right) dz_{i}$$

$$= -\sum_{i=1}^{b} \int \mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \log\left(\frac{\mathbb{P}(x_{i} \mid z_{i}, \theta_{d})}{\mathbb{P}(z_{i})}\right) dz_{i}$$

$$= -\sum_{i=1}^{b} \int \mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \log\left(\mathbb{P}(x_{i} \mid z_{i}, \theta_{d})\right) dz_{i}$$

$$= -\sum_{i=1}^{b} \operatorname{KL}(\mathbb{P}(z_{i} \mid x_{i}, \theta_{e}) \mid \mathbb{P}(z_{i}))$$

$$+ \sum_{i=1}^{b} \mathbb{E}_{\sim \mathbb{P}(z_{i} \mid x_{i}, \theta_{e})} \left[\log\left(\mathbb{P}(x_{i} \mid z_{i}, \theta_{d})\right)\right]. \quad (25)$$

• We found:

$$\bigstar \quad \sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \mathsf{KL}\big(\mathbb{P}(\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \theta_{e}) \parallel \mathbb{P}(\boldsymbol{z}_{i})\big) + \sum_{i=1}^{b} \underbrace{\mathbb{E}_{\sim \mathbb{P}(\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \theta_{e})}\Big[\log\big(\mathbb{P}(\boldsymbol{x}_{i} \mid \boldsymbol{z}_{i}, \theta_{d})\big)\Big]}_{i=1}.$$

The second term in the above equation can be estimated using <u>Monte Carlo</u> approximation [7] where we draw *l* samples {*z<sub>i</sub>*,*j*}<sup>*l*</sup><sub>*i*=1</sub> from P(*z<sub>i</sub>* | *x<sub>i</sub>*, *θ<sub>e</sub>*):

$$\sum_{i=1}^{b} \mathcal{L}(q, \theta) \approx \sum_{i=1}^{b} \underbrace{\mathcal{L}}(q, \theta)$$
$$= -\sum_{i=1}^{b} \underbrace{\mathsf{KL}}(\mathbb{P}(z_{i} \mid \mathbf{x}_{i}, \theta_{e}) \parallel \mathbb{P}(z_{i})) + \sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \log\left(\mathbb{P}(\mathbf{x}_{i} \mid z_{i,j}, \theta_{d})\right).$$
(26)

We had:

We had:  

$$\underbrace{} \qquad \underbrace{} \qquad$$

1

• 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  $\{z_{i,j}\}_{j=1}^{\ell}$  from  $\mathbb{P}(z_i | x_i, \theta_e)$ :

$$\sum_{i=1}^{b} \mathcal{L}(q, \theta) \approx \sum_{i=1}^{b} \widetilde{\mathcal{L}}(q, \theta)$$

$$= -\sum_{i=1}^{b} \underbrace{\mathbb{E}_{\sim \mathbb{P}(\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \theta_{e})} \left[ \log \left( \frac{\mathbb{P}(\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \theta_{e})}{\mathbb{P}(\boldsymbol{z}_{i})} \right) \right]}_{\boldsymbol{z}} + \sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \log \left( \mathbb{P}(\boldsymbol{x}_{i} \mid \boldsymbol{z}_{i,j}, \theta_{d}) \right)$$

$$\bigotimes_{i=1}^{b} -\sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \log \left( \mathbb{P}(\boldsymbol{z}_{i,j} \mid \boldsymbol{x}_{i}, \theta_{e}) \right) - \log \left( \mathbb{P}(\boldsymbol{z}_{i,j}) \right) + \sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \log \left( \mathbb{P}(\boldsymbol{x}_{i} \mid \boldsymbol{z}_{i,j}, \theta_{d}) \right).$$
(27)

In case we have some <u>families</u> of distributions, such as <u>Gaussian</u> distributions, for

 *P*(z<sub>i,j</sub> | x<sub>i</sub>, θ<sub>e</sub>) and *P*(z<sub>i,j</sub>), the first term in Eq. (26) can be <u>computed analytically</u>. In the following, we simply Eq. (26) further for Gaussian distributions.

# Simplification Type 2 for Special Case of Gaussian Distributions

 We can compute the <u>KL divergence</u> in the <u>first term</u> of Eq. (26) analytically for <u>univariate</u> or <u>multivariate Gaussian</u> distributions. For this, we need two following <u>lemmas</u> (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:  
 $KL(p_1 \| p_2) = \log(\frac{\sigma_2}{\sigma_1}) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma^2} - \frac{1}{2}.$  (28)

#### Lemma

The KL divergence between two <u>multivariate</u> Gaussian distributions  $\underline{p_1 \sim \mathcal{N}(\mu_1, \mathbf{\Sigma}_1)}$  and  $p_2 \sim \mathcal{N}(\mu_2, \mathbf{\Sigma}_2)$  with dimensionality p is:

$$\mathcal{KL}(p_1 \| p_2) = \frac{1}{2} \Big( \log(\frac{|\mathbf{\Sigma}_2|}{|\mathbf{\Sigma}_1|}) - p + \operatorname{tr}(\mathbf{\Sigma}_2^{-1}\mathbf{\Sigma}_1) + (\mu_2 - \mu_1)^\top \mathbf{\Sigma}_2^{-1}(\mu_2 - \mu_1) \Big).$$
(29)

# Simplification Type 2 for Special Case of Gaussian Distributions

Consider the case in which we have:

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

$$\underbrace{\mathbb{P}(\mathbf{z}_i)}_{\mathbf{0}} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}}), \tag{31}$$

where  $z_i \in \mathbb{R}^{p}$ . Note that the parameters  $\mu_{z|x}$  and  $\Sigma_{z|x}$  are trained in neural network while the parameters  $\mathbb{P}(z_{i,j})$  can be set to  $\mu_z = 0$  and  $\Sigma_z = I$  (inspired by the prior distribution of z in factor analysis).

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

$$\sum_{i=1}^{b} \mathcal{L}(\boldsymbol{q}, \boldsymbol{\theta}) \approx \sum_{i=1}^{b} \widetilde{\mathcal{L}}(\boldsymbol{q}, \boldsymbol{\theta})$$

$$= -\sum_{i=1}^{b} \frac{1}{2} \Big( \log(\frac{|\boldsymbol{\Sigma}_{z}|}{|\boldsymbol{\Sigma}_{z|\times}|}) - \boldsymbol{p} + \operatorname{tr}(\boldsymbol{\Sigma}_{z}^{-1}\boldsymbol{\Sigma}_{z|\times}) + (\boldsymbol{\mu}_{z} - \boldsymbol{\mu}_{z|\times})^{\top} \boldsymbol{\Sigma}_{z}^{-1}(\boldsymbol{\mu}_{z} - \boldsymbol{\mu}_{z|\times}) \Big)$$

$$+ \sum_{i=1}^{b} \frac{1}{\ell} \sum_{j=1}^{\ell} \log \left( \mathbb{P}(\boldsymbol{x}_{i} \mid \boldsymbol{z}_{i,j}, \boldsymbol{\theta}_{d}) \right).$$
(32)

#### Training Variational Autoencoder with Approximations

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

$$\begin{cases} \mathsf{E}\text{-step:} \quad \boldsymbol{\theta}_{e}^{(t)} \coloneqq \boldsymbol{\theta}_{e}^{(t-1)} + \eta_{e} \frac{\partial \sum_{i=1}^{b} \widehat{\mathcal{G}}(\boldsymbol{\theta}_{e}, \boldsymbol{\theta}_{d}^{(t-1)}, \mathbf{x}_{i})}{\partial \boldsymbol{\theta}_{e}}, \\ \mathsf{M}\text{-step:} \quad \boldsymbol{\theta}_{d}^{(t)} \coloneqq \boldsymbol{\theta}_{d}^{(t-1)} + \eta_{d} \frac{\partial \sum_{i=1}^{b} \widehat{\mathcal{G}}(\boldsymbol{\theta}_{e}^{(t)}, \boldsymbol{\theta}_{d}, \mathbf{x}_{i})}{\partial \boldsymbol{\theta}_{d}}, \end{cases}$$

are replaced by the following equations:

$$\begin{aligned} & \mathsf{E}\text{-step:} \quad \boldsymbol{\theta}_{e}^{(t)} := \boldsymbol{\theta}_{e}^{(t-1)} + \eta_{e} \frac{\partial \sum_{i=1}^{b} \left( \widehat{\mathcal{L}} \boldsymbol{\theta}_{e}, \boldsymbol{\theta}_{d}^{(t-1)}, \boldsymbol{x}_{i} \right)}{\partial \boldsymbol{\theta}_{e}}, \end{aligned}$$
(33)

$$\bigwedge \text{M-step:} \quad \boldsymbol{\theta}_{d}^{(t)} := \boldsymbol{\theta}_{d}^{(t-1)} + \eta_{d} \frac{\partial \sum_{i=1}^{b} \left( \widetilde{\mathcal{L}}(\boldsymbol{\theta}_{e}^{(t)}, \boldsymbol{\theta}_{d}, \boldsymbol{x}_{i}) \right)}{\partial \boldsymbol{\theta}_{d}}, \tag{34}$$

where the approximated ELBO was introduced in previous sections.

# The Reparameterization Trick

Sampling the ℓ samples for the latent variables, i.e. Eq. (15):

$$\boldsymbol{z}_i \sim \boldsymbol{q}(\boldsymbol{z}_i) = \mathbb{P}(\boldsymbol{z}_i \,|\, \boldsymbol{x}_i, \boldsymbol{\theta}_e),$$

blocks the gradient flow because computing the derivatives through  $\mathbb{P}(z_i | x_i, \theta_e)$  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  $z_i \sim \mathbb{P}(z_i | x_i, \theta_e)$ , we assume  $z_i$  is a random variable but is a deterministic function of another random variable  $\epsilon_i$  as follows:

$$\mathbf{z}_i = g(\boldsymbol{\epsilon}_i, \mathbf{x}_i, \boldsymbol{\theta}_e), \tag{35}$$

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

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

### The Reparameterization Trick

• The Eqs. (21) and (25):

$$\begin{cases} \sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \mathbb{E}_{\sim \mathbb{P}(z_i \mid \mathbf{x}_i, \theta_e)} \Big[ \log \left( \frac{\mathbb{P}(z_i \mid \mathbf{x}_i, \theta_e)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \theta_d)} \right) \Big], \\ \sum_{i=1}^{b} \mathcal{L}(q, \theta) = -\sum_{i=1}^{b} \mathsf{KL} \big( \mathbb{P}(z_i \mid \mathbf{x}_i, \theta_e) \parallel \mathbb{P}(z_i) \big) + \sum_{i=1}^{b} \mathbb{E}_{\sim \mathbb{P}(z_i \mid \mathbf{x}_i, \theta_e)} \Big[ \log \big( \mathbb{P}(\mathbf{x}_i \mid z_i, \theta_d) \big) \big], \end{cases}$$

both contain an expectation of a function  $f(z_i)$ . Using this technique, this expectation is replaced as:

- Using the reparameterization technique, the <u>encoder</u>, which implemented  $\mathbb{P}(z_i | x_i, \theta_e)$ , is replaced by  $g(\epsilon_i, x_i, \theta_e)$  where in the <u>latent space</u> between encoder and decoder, we have  $\epsilon_i \sim \mathbb{P}(\epsilon)$  and  $z_i = g(\epsilon_i, x_i, \theta_e)$ .
- A simple example for the reparameterization technique is when z<sub>i</sub> and ε<sub>i</sub> are univariate Gaussian variables:

$$egin{aligned} & (z_i \sim \mathcal{N}(\mu, \sigma^2), \ & (\epsilon_i \sim \mathcal{N}(0, 1), \ & (z_i = g(\epsilon_i) = \mu + \sigma \epsilon_i. ) \end{aligned}$$

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**,  $\sum_{i=1}^{b} -\hat{\mathcal{L}}(\theta, \mathbf{x}_i)$ , as the loss function:

$$\bigstar \left[ \boldsymbol{\theta}^{(t)} := \boldsymbol{\theta}^{(t-1)} - \eta \frac{\partial \sum_{i=1}^{b} - \widetilde{\mathcal{L}}(\boldsymbol{\theta}, \mathbf{x}_{i})}{\partial \boldsymbol{\theta}} \right]$$
(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<sub>i</sub> to the encoder to determine the parameters of the conditional distribution of latent space, i.e., P(z<sub>i</sub> | x<sub>i</sub>, θ<sub>e</sub>).
- Then, from this distribution, we <u>sample the latent variable z</u><sub>i</sub> from the latent space and generate the corresponding **reconstructed data point x**<sub>i</sub> 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
  - Iower 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.



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 <u>deep learning course</u> 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

- 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(R) 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. Karray, 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. Karray, 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.