

20170243 김재운

$$\log p_{\theta}(x) = \log \left( \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \right] \right)$$

$$(a) \text{KL} (q_{\phi}(z|x), p_{\theta}(z))$$

Training VAE  $\Leftrightarrow$  Maximization ELBO

$$\arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

$$\log p_{\theta}(x) \geq \mathcal{L}(x, \theta, \phi)$$

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

$$= \mathbb{E}_z \left[ \log p_{\theta}(x|z) \right] - \text{KL} (q_{\phi}(z|x), p_{\theta}(z))$$

$$(b) q_{\phi}(z|x) = N(M_{z|x}, \Sigma_{z|x}^2), P_{\theta}(z) = N(\mu, \sigma^2)$$

$KL(q_{\phi}(z|x), P_{\theta}(z))$  has an analytic solution if both  $q$  and  $p$  follows the normal distribution. And all of integrals are about marginalization over  $Z$ .

$$\left\{ \begin{array}{l} q_{\phi}(z|x) = N(M_{z|x}, \Sigma_{z|x}^2) = \frac{1}{\sqrt{2\pi \Sigma_{z|x}^2}} \exp\left(-\frac{(z - M_{z|x})^2}{2\Sigma_{z|x}^2}\right) \\ P_{\theta}(z) = N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi \sigma^2}} \exp\left(-\frac{(z - \mu)^2}{2\sigma^2}\right) \end{array} \right.$$

$$\begin{aligned} \rightarrow KL(q_{\phi}(z|x), P_{\theta}(z)) &= \int q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{P_{\theta}(z)} dz \\ &= \int q_{\phi}(z|x) \log q_{\phi}(z|x) dz - \int q_{\phi}(z|x) \log P_{\theta}(z) dz \\ &= \underbrace{\int N(M_{z|x}, \Sigma_{z|x}^2) \log N(M_{z|x}, \Sigma_{z|x}^2) dz}_{\textcircled{1}} - \underbrace{\int N(M_{z|x}, \Sigma_{z|x}^2) \log N(\mu, \sigma^2) dz}_{\textcircled{2}} \end{aligned}$$

$$\textcircled{1} = \int q_{\phi}(z|x) \log q_{\phi}(z|x) dz$$

$$= \int A_q \exp\left(-\frac{(z - M_{z|x})^2}{2\Sigma_{z|x}^2}\right) \left( \log A_q + -\frac{(z - M_{z|x})^2}{2\Sigma_{z|x}^2} \right) dz \dots$$

$$\frac{1}{\sqrt{2\pi \Sigma_{z|x}^2}} = A_q$$

$$= \int A_q \exp(-B_q^2) (\log A_q - B_q^2) \sqrt{2\Sigma_{z|x}} dB_q \dots$$

$$\frac{z - M_{z|x}}{\sqrt{2\Sigma_{z|x}}} = B_q$$

$$\frac{1}{\sqrt{2\Sigma_{z|x}}} dz = dB_q$$

$$= \int A_q \exp(-B_q^2) \sqrt{2\Sigma_{z|x}} \log A_q dB_q - \int A_q \exp(-B_q^2) \sqrt{2\Sigma_{z|x}} B_q^2 dB_q$$

$$\begin{aligned}
&= \int \frac{\sqrt{2} \sqrt{\tau_{2|x}}}{\sqrt{2\pi} \sqrt{\tau_{2|x}}} \exp(-B_q^2) \log(2\pi \sqrt{\tau_{2|x}})^{-\frac{1}{2}} dB_q - \int \frac{\sqrt{2} \sqrt{\tau_{2|x}}}{\sqrt{2\pi} \sqrt{\tau_{2|x}}} \exp(-B_q^2) B_q^2 dB_q \\
&= -\frac{\log(2\pi \sqrt{\tau_{2|x}})}{2\sqrt{\pi}} \int \exp(-B_q^2) dB_q - \frac{1}{\sqrt{\pi}} \int B_q^2 \exp(-B_q^2) dB_q
\end{aligned}$$

Gaussian integral

$$\left\{
\begin{array}{l}
\int \exp(-x^2) dx = \sqrt{\pi} \\
\int x^n \exp(-ax^2) dx = \frac{2(n-1)!}{2 \frac{(n+2)}{2} a^{\frac{n}{2}}} \sqrt{\frac{\pi}{a}} \quad (n: \text{even number})
\end{array}
\right.$$

$$\begin{aligned}
&= -\frac{\log(2\pi \sqrt{\tau_{2|x}})}{2\sqrt{\pi}} \sqrt{\pi} - \frac{1}{\sqrt{\pi}} \cdot \frac{2}{4} \cdot \sqrt{\pi} \dots \boxed{n=2, a=1} \\
&= -\frac{\log(2\pi \sqrt{\tau_{2|x}})}{2} - \frac{1}{2} \\
&= -\frac{1}{2} (\log(2\pi \sqrt{\tau_{2|x}}) + 1)
\end{aligned}$$

$$\begin{aligned}
② &= - \int q_{f_\phi}(z|x) \log p_\theta(z) dz \\
&= - \int q_{f_\phi}(z|x) \log \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right) \right\} dz \\
&= - \int q_{f_\phi}(z|x) \log \frac{1}{\sqrt{2\pi\sigma^2}} dz + \int q_{f_\phi}(z|x) \left( \frac{(z-\mu)^2}{2\sigma^2} \right) dz \\
&= \frac{1}{2} \log(2\pi\sigma^2) \int q_{f_\phi}(z|x) dz + \frac{1}{2\sigma^2} \left\{ \int q_{f_\phi}(z|x) z^2 dz - \int q_{f_\phi}(z|x) 2\mu z dz \right. \\
&\quad \left. + \int \mu^2 q_{f_\phi}(z|x) dz \right\} \\
&\dots \boxed{\int q_{f_\phi}(z|x) dz = 1}
\end{aligned}$$

$$= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\mathbb{E}(z^2|x) - 2\mu \mathbb{E}(z|x) + \mu^2)$$

$$= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + \mu_{z|x}^2 - 2\mu \mu_{z|x} + \mu^2) \dots$$

$$= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2)$$

$\sigma^2 = \mathbb{E}[x^2] - \mu^2$
$\mathbb{E}[x^2] = \sigma^2 + \mu^2$

$$\therefore KL(q_\theta(z|x), p_\theta(z)) = \textcircled{1} + \textcircled{2}$$

$$= \left( -\frac{1}{2} \log (2\pi\sigma_{z|x}^2) - \frac{1}{2} \right) + \left( \frac{1}{2} \log (2\pi\sigma^2) + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) \right)$$

$$= \frac{1}{2} \log \left( \frac{2\pi\sigma^2}{2\pi\sigma_{z|x}^2} \right) + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) - \frac{1}{2}$$

$$= \frac{1}{2} \log \frac{\sigma^2}{\sigma_{z|x}^2} + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) - \frac{1}{2}$$

$$= \frac{1}{2} \cdot 2 \cdot \log \frac{\sigma}{\sigma_{z|x}} + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) - \frac{1}{2}$$

$$= \log \frac{\sigma}{\sigma_{z|x}} + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) - \frac{1}{2}$$

If  $\mu = \mu_{z|x}$  &  $\sigma = \sigma_{z|x}$ , the KL divergence becomes zero.

$$(c) P_\theta(x|z) = \mathcal{N}(M_{z|x}, \Sigma_{z|x}^2)$$

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \left( \frac{P_\theta(x, z)}{q_\phi(z|x)} \right) \right] \rightarrow \text{intractable!}$$

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \left( \frac{P_\theta(x, z)}{q_\phi(z|x)} \right) \right] \approx \sum_{j=1}^J \nabla_\phi \log \left( \frac{P_\theta(x, z^{(j)})}{q_\phi(z^{(j)}|x)} \right)$$

with  $z^{(j)} \sim q_\phi(z^{(j)}|x)$  → still intractable!

→ reparametrization trick

Represent random variable  $z \sim \mathcal{N}(M_{z|x}, \Sigma_{z|x}^2)$  as a function of  $M_{z|x}, \Sigma_{z|x}^2$ , and  $\epsilon \sim \mathcal{N}(0, 1)$ :

$$z = \Sigma_{z|x} \epsilon + M_{z|x}$$

Sample of  $z$  using sample from  $\mathcal{N}(0, 1)$ :

$$z^{(j)} = \Sigma_{z|x} \epsilon^{(j)} + M_{z|x} \quad \forall j = 1, \dots, J$$

Approximate  $\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \log \left( \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right]$  with sampling and reparametrization trick

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \log \left( \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right]$$

$$= \nabla_{\phi} \mathbb{E}_z \left[ \log \left( \frac{p_{\theta}(x|z) p_{\theta}(z)}{q_{\phi}(z|x)} \right) \right]$$

$$= \nabla_{\phi} \mathbb{E}_z \left[ \log p_{\theta}(x|z) + \log \left( \frac{p_{\theta}(z)}{q_{\phi}(z|x)} \right) \right]$$

$$= \nabla_{\phi} \left( \mathbb{E}_z \left[ \log p_{\theta}(x|z) \right] + \mathbb{E}_z \left[ \log \left( \frac{p_{\theta}(z)}{q_{\phi}(z|x)} \right) \right] \right)$$

$$= \nabla_{\phi} \left( \mathbb{E}_z \left[ \log p_{\theta}(x|z) \right] - \mathbb{E}_z \left[ \log \left( \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \right) \right] \right)$$

$$= \nabla_{\phi} \left( \mathbb{E}_z \left[ \log p_{\theta}(x|z) \right] - KL(q_{\phi}(z|x), p_{\theta}(z)) \right)$$

$$\approx \underbrace{\frac{1}{J} \sum_{j=1}^J \nabla_{\phi} \log p_{\theta}(x|z^{(j)})}_{\textcircled{1}} - \underbrace{\nabla_{\phi} KL(q_{\phi}(z|x), p_{\theta}(z))}_{\textcircled{2}}$$

$$\textcircled{1} = \frac{1}{J} \sum_{j=1}^J \nabla_{\phi} \log p_{\theta}(x|z^{(j)} + \mu_{z|x}) \dots$$

from reparametrization  
&  
sampling

$$\textcircled{2} = - \nabla_{\phi} \left( \log \frac{\sigma}{\sigma_{z|x}} + \frac{1}{2\sigma^2} (\sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2) - \frac{1}{2} \right) \dots$$

from (b) result

∴ approximation result = ①+②

$$= \frac{1}{J} \sum_{j=1}^J \nabla_{\phi} \log P_{\phi}(x | \mathcal{T}_{z|x}, \mathcal{E}^{(j)} + \mu_{z|x}) \\ - \nabla_{\phi} \left( \log \frac{1}{J z|x} + \frac{1}{2J^2} \left( \sigma_{z|x}^2 + (\mu_{z|x} - \mu)^2 \right) - \frac{1}{2} \right)$$

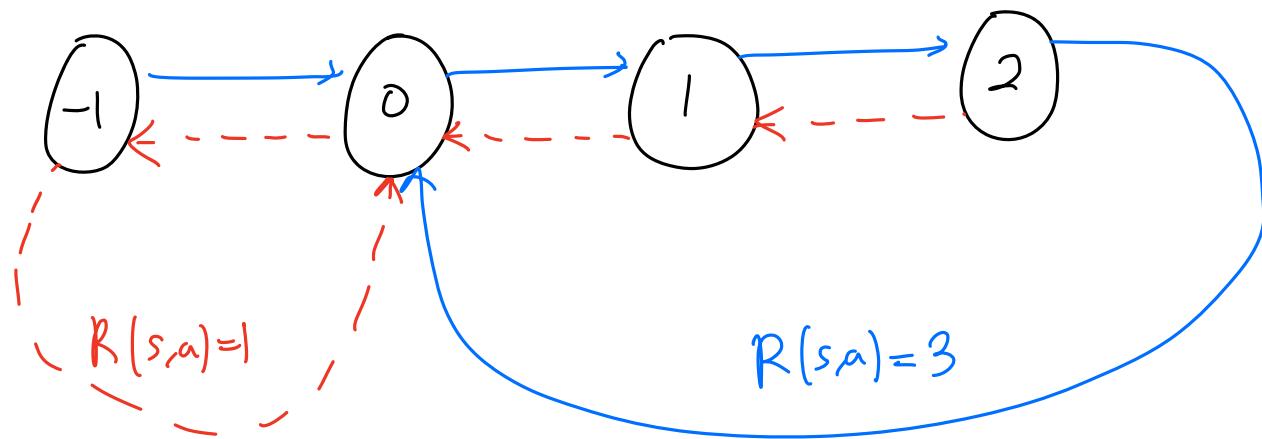
(d)

```
def reparameterize(self, mu, logvar):  
    #####  
    ## implement the reparameterize function  
    std = torch.exp(logvar / 2)  
    eps = torch.randn_like(std)  
    #####  
    return eps.mul(std).add_(mu)
```

PROBLEMS	OUTPUT	DEBUG CONSOLE	TERMINAL
Train Epoch: 26 [25600/60000 (43%)]			Loss: 101.455467
Train Epoch: 26 [38400/60000 (64%)]			Loss: 99.498177
Train Epoch: 26 [51200/60000 (85%)]			Loss: 108.227180
===== Epoch: 26 Average loss: 103.9785 =====			
Train Epoch: 27 [0/60000 (0%)]	Loss: 100.385132		
Train Epoch: 27 [12800/60000 (21%)]		Loss: 101.097992	
Train Epoch: 27 [25600/60000 (43%)]		Loss: 103.895073	
Train Epoch: 27 [38400/60000 (64%)]		Loss: 101.896179	
Train Epoch: 27 [51200/60000 (85%)]		Loss: 103.818451	
===== Epoch: 27 Average loss: 103.8807 =====			
Train Epoch: 28 [0/60000 (0%)]	Loss: 104.533897		
Train Epoch: 28 [12800/60000 (21%)]		Loss: 103.687950	
Train Epoch: 28 [25600/60000 (43%)]		Loss: 103.391647	
Train Epoch: 28 [38400/60000 (64%)]		Loss: 104.651192	
Train Epoch: 28 [51200/60000 (85%)]		Loss: 105.794754	
===== Epoch: 28 Average loss: 103.7341 =====			
Train Epoch: 29 [0/60000 (0%)]	Loss: 104.814400		
Train Epoch: 29 [12800/60000 (21%)]		Loss: 104.002884	
Train Epoch: 29 [25600/60000 (43%)]		Loss: 104.355698	
Train Epoch: 29 [38400/60000 (64%)]		Loss: 103.353668	
Train Epoch: 29 [51200/60000 (85%)]		Loss: 101.666145	
===== Epoch: 29 Average loss: 103.6788 =====			
(assn) simjaeyoon@simjaeyoonui-MacBookPro AIGS515_ASSN5 %			□



2



$$(s, a) \in \{(-1, -), (2, +)\} \rightarrow \text{non-zero}$$

$$r(-1, -) = 1$$

$$r(2, +) = 3$$

(a)  $\gamma = 0.10$

optimal policy  $\pi^* = (\underset{s: -1}{\circ}, \underset{s: 0}{\circ}, \underset{s: 1}{\circ}, \underset{s: 2}{\circ})$

optimal value function  $v^*(s, \gamma) = (\underset{s: -1}{\circ}, \underset{s: 0}{\circ}, \underset{s: 1}{\circ}, \underset{s: 2}{\circ})$

It is possible to generate 16 ( $= 2^4$ ) combinations of policies. We can use some iteration method to compute value function and find optimal policy. To compute the value function, we can use recursive function form. Since discount factor  $\gamma$  is near to 0, it would converge much faster. So it needs shorter iteration to find optimal solution. We need large reward as soon as possible.

① (-, -, +/-, +/-)

$$V_{\pi}(-1) = 1 + 0 + \gamma^2 \cdot 1 + 0 + \gamma^4 \cdot 1 + \dots = \frac{1}{1 - \gamma^2} = 1.01$$

② (-, -, +/-, +/-)

$$V_{\pi}(0) = 0 + \gamma \cdot 1 + 0 + \gamma^3 \cdot 1 + 0 + \gamma^5 \cdot 1 + \dots = \frac{\gamma}{1 - \gamma^2} = 0.10$$

③ (-, -, +, +)

$$V_{\pi}(1) = 0 + \gamma^1 \cdot 3 + 0 + \gamma^3 \cdot 1 + 0 + \gamma^5 \cdot 1 + 0 + \gamma^7 \cdot 1 = 3 + \frac{\gamma^3}{1 - \gamma^2} = 0.301$$

④ (-, -, +/-, +)

$$V_{\pi}(2) = 3 + 0 + \gamma^2 \cdot 1 + 0 + \gamma^4 \cdot 1 + \dots = 3 + \frac{\gamma^2}{1 - \gamma^2} = 3.01$$

$$\pi_{short}^* = (-, -, +, +), V^*(s, r) = (1.01, 0.10, 0.301, 3.01)$$

(b)  $\gamma = 0.99$ .

Method is similar to (a), but the discount factor is close to 1, so it is okay to get reward slower.

-1 (-, +, +, +)

$$V_{\pi}(-1) = 1 + 0 + 0 + \gamma^3 \cdot 3 + 0 + 0 + \gamma^6 \cdot 3 + \dots = 1 + \frac{3 \cdot \gamma^3}{1 - \gamma^3} = 99.001$$

0 (+/-, +, +, +)

$$V_{\pi}(0) = 0 + 0 + \gamma^2 \cdot 3 + 0 + 0 + \gamma^5 \cdot 3 + \dots = \frac{3 \gamma^2}{1 - \gamma^3} = 98.991$$

1 (+/-, +, +, +)

$$V_{\pi}(1) = 0 + \gamma \cdot 3 + 0 + 0 + \gamma^4 \cdot 3 + 0 + 0 + \dots = \frac{3 \gamma}{1 - \gamma^3} = 99.991$$

2 (+/-, +, +, +)

$$V_{\pi}(2) = 3 + 0 + 0 + \gamma^3 \cdot 3 + 0 + 0 + \gamma^6 \cdot 3 = 3 + \frac{3 \gamma^3}{1 - \gamma^3} = 101.001$$

$$\pi_{long}^* = (-, +, +, +), V^*(s, \gamma) = (99.001, 98.991, 99.991, 101.001)$$

$\underbrace{\hspace{20em}}$   
approximation.

(c)  $\gamma=0.99$  with problem 2a.

$$V^{\pi_{\text{short}}^*}(s, \gamma=0.99) = (50.25, 49.75, 51.13, 52.25)$$

$$V^*(s, \gamma=0.10) = (1.01, 0.10, 0.301, 3.01)$$

$\gamma=0.10 \rightarrow$  accumulate  $\downarrow$  step  $\rightarrow$  converge fast  
 $\rightarrow$  low optimal value function.

$\gamma=0.99 \rightarrow$  accumulate  $\uparrow$  step  $\rightarrow$  converge slow  
 $\rightarrow$  high optimal value function.

$\gamma$  controls the accumulation of value function from longer timestep. If  $\gamma$  is close to 1, it can consider far future, however if  $\gamma$  is close to 0, it can consider near future.

When discount factor becomes larger ( $0.10 \rightarrow 0.99$ ), then optimal value function results become larger comparing with  $\pi_{\text{short}}^*$ . Comparing the result between  $\pi_{\text{long}}^*$  and  $\pi_{\text{short}}^*$  with  $\gamma=0.99$ ,  $\pi_{\text{long}}^*$  is optimal policy when  $\gamma$  is large.