

# VAE 정리

VAE는 잠재 변수를 가지는 생성 모델이다.

생성 방식 :  $z$ 가 표준 정규 분포를 따른다는 가정(prior)을 기반으로  $p(z)$ 에서  $z$ 를 추출하고, 추출된  $z$ 를 기반으로  $x$ 를 생성하는 모델이다.

$$\log P_\theta(x) = \log \frac{P_\theta(x, z)}{P_\theta(z|x)} \quad (\text{균형 방정식 } p(x, y) = p(x|y)p(y) \Rightarrow \log P_\theta(x) \text{ 변형})$$

$$\log \frac{P_\theta(x, z)}{P_\theta(z|x)} \text{ 에서 } P_\theta(z|x) = \frac{P_\theta(z, x)}{P_\theta(x)} = \frac{P_\theta(z, x)}{\sum_z P_\theta(x, z)} \text{ 로 연립 등식 제거 ... , log-sum이 생긴다.}$$

**Solution 1:**  $P_\theta(z|x)$ 의 근사값으로 임의의 확률분포  $q(z)$ 를 이용.

$$\log P_\theta(x) = \log \frac{P_\theta(x, z)}{P_\theta(z|x)} \times \frac{q(z)}{1}$$

$$= \log \frac{P_\theta(x, z)}{q(z)} \times \frac{q(z)}{P_\theta(z|x)}$$

$$= \log \frac{P_\theta(x, z)}{q(z)} + \log \frac{q(z)}{P_\theta(z|x)} \quad \text{②}$$

①  $P_\theta(z|x)$ 를  $q(z)$ 로 바꿈. (ok)  $\leftarrow$   $P_\theta(z|x)$ 가 남아있음 (문제)

**Solution 2:** ②번 항을 KL 식으로 바꾸기

$$\log P_\theta(x) = \log P_\theta(x) \sum_z q(z) \quad \text{①}$$

$$= \sum_z q(z) \log P_\theta(x) \quad \text{②}$$

$$= \sum_z q(z) \left( \log \frac{P_\theta(x, z)}{q(z)} + \log \frac{q(z)}{P_\theta(z|x)} \right) \quad \text{KL로 유도됨!}$$

$$= \sum_z q(z) \log \frac{P_\theta(x, z)}{q(z)} + \sum_z q(z) \log \frac{q(z)}{P_\theta(z|x)}$$

$$= \sum_z q(z) \log \frac{P_\theta(x, z)}{q(z)} + D_{KL}(q(z) \parallel P_\theta(z|x))$$

$\Rightarrow$  ①, ② 어디에도 log-sum 형태가 없음!

앞서 다뤘듯이 잠재 변수를 가지는 생성 모델에서 하나의 데이터  $x$ 에 대한 log-likelihood는 공통 확률 분포에서의 잠재 변수 marginalization으로 인해 log-sum의 형태를 지니고 이는 해석적으로 최적화(최대화) 하기 어려운 문제점이 있다.

이 log-sum의 형태를 없애기 위해 원래  $\log P_\theta(x)$ 의 값을 근사하기 위해(더 구체적으로는  $\log P_\theta(x)$ 를 풀어서 구성한  $p_\theta(z|x)$ 항을 근사하기 위해)  $q(z)$ 를 도입한다. (VAE의 경우  $q(z)$ 가 관측 데이터  $x$ 를 기반으로 만들어지기에,  $q(z|x)$ 로 표현함)

이렇게  $q(z)$ 를 도입하여 log-likelihood의 식을 풀어내면 원래의 log-likelihood와 식이 근사한 ELBO항과 KL항으로 풀려진다. 우리의 목표는 원래의 log-likelihood항을 최대화 하는 것인데 이 항이 log-sum의 형태를 지니고 있기에 해석적인 최적화가 어렵고 앞서 sum 형태를 없애기 위해 풀어진 두 항 중 log-likelihood항과 유사한 1) ELBO항을 간접적으로 최대화하면서도 2)이 ELBO항이 우리가 원래 최대화시키고자 했던 log-likelihood와 비슷하도록 만들어야한다. 1)은 ELBO를 구성하는  $\theta$ 를 업데이트하면서 달성되고 2)는  $q(z)$ 를  $p_\theta(z|x)$ 로 업데이트 함으로써 KL 항을 0으로 만들며 달성된다.

결론적으로 도입하는  $q(z)$ 가 도입된 이유는 원래의 log-likelihood항에 있는  $p_\theta(z|x)$ 를 근사하기 위함이다.

1)  $\rightarrow$  M step

2)  $\rightarrow$  E step

그러나 VAE의 경우 전통적인 EM 알고리즘 수행을 위한 E-step이 어렵다.  $q(z)$ 를  $p_{\theta}(z|x)$ 로 업데이트 하는 과정에서 뒤의 항을 구하는 것이 어렵기 때문이다.  $z$ 가 연속적인 값을 가지는 벡터이기에 주변화가 어려움.

$$\begin{aligned} \text{for VAE, } P_{\theta}(z^{(m)}|x^{(m)}) &= \frac{P_{\theta}(x^{(m)}, z^{(m)})}{P_{\theta}(x^{(m)})} & p(x) &= \int p(x, z) dz \\ & & &= \int p(z) p(x|z) dz \\ &= \frac{P_{\theta}(x^{(m)}, z^{(m)})}{\int P_{\theta}(x^{(m)}, z^{(m)}) dz} \rightarrow \text{marginalization} \\ \int P_{\theta}(x, z) dz &\text{ is easily countable in GMM since latent variable } z \text{ is discrete.} \\ \Leftrightarrow \sum_z P_{\theta}(x, z) dz \\ \text{But in VAE } z &\text{ is continuous and } z \text{ is vector} \rightarrow \int P_{\theta}(x, z) dz \text{ is impossible (or very hard)} \end{aligned}$$

$q(z)$ 를  $p_{\theta}(z|x)$ 로 근사하기 위해(업데이트 하기 위해) 다른 방법을 이용하는 것이다.

그래서  $q(z)$ 를 파라미터로 조절 가능한 tractable한 분포로 바꾼다. 다시 말해 파라미터로 조절 가능한 분포인 **가우시안 분포**로 만든다. 이후에 ELBO를  $q(z)$ 를 구성하는 파라미터(pie)와 theta에 대해 최대화되도록 신경망을 training하면  $q(z)$ 는  $p_{\theta}(z|x)$ 에 근사되고(EM에서의 E-step)  $x$ 에 대한 likelihood도 최대화 된다(M-step)

→ 이렇듯 계산이 어려운 복잡한 분포  $p_{\theta}(z|x)$ 를 tractable한 분포로 근사하는 방식을 변분 추론/변분 근사라고 부른다.

→ 신경망 학습을 통해 동시에 두가지 종류의 파라미터를 학습시켜 ELBO가 최대화 되도록 업데이트 하면 ELBO는 원래의 log-likelihood에 가까워 지면서도( $q(z)$ 는  $p_{\theta}(z|x)$ 에 근사되고(EM에서의 E-step)), 원래의 log-likelihood는 커진다.

그런데 매번 하나의 데이터  $x$ 에 대해  $q(z)$ 에 대한 파라미터를 일일이 준비할 수는 없음. 그래서 이 파라미터 만드는 작업을  $x$ 를 입력으로 받는 신경망을 통해 구현함.

신경망은  $x$ 를 입력으로 받아  $q(z)$ 의 평균 벡터와 공분산 행렬의 대각 성분을 구성하는 벡터를 출력함. VAE에서는  $x$ 로부터  $z$ 의 분산이 결정되는 것이기에,  $q_{\pi}(z|x)$ 로 표현되고 이때  $\pi$ 는 VAE에서 인코더 신경망의 파라미터임. 결론적으로 ELBO항의  $q(z|x)$ 를 표현하기 위해 인코더 신경망이 생겨난 것임.

아래는 VAE에서 하나의 데이터  $x$ 에 대한 ELBO항의 식이다.

$$\begin{aligned} \text{Single Sample } x : \text{ELBO}(x; \theta, \phi) &= \int q_{\phi}(z|x) \log \frac{P_{\theta}(x, z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log \frac{P_{\theta}(x|z) P(z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log P_{\theta}(x|z) dz + \int q_{\phi}(z|x) \log \frac{P(z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log P_{\theta}(x|z) dz - \int q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{P(z)} dz \\ &= \underbrace{E_{q_{\phi}(z|x)} [\log P_{\theta}(x|z)]}_{J_1} - \underbrace{D_{KL}(q_{\phi}(z|x) || P(z))}_{J_2} \end{aligned}$$

ELBO항 또 위와 같이 두개의 항으로 쪼개짐.

$J_1$  is  $\int q_\phi(z|x)$ 를 따르는  $z$ 에 대한  $\log p_\theta(x|z)$ 의 Expectation  
 Monte carlo approximation:  $\int q_\phi(z|x)$ 에서  $z$ 를 "1"개만 sampling 해서 근사

$$u, \sigma = \text{Neural Net}(x; \phi)$$

$$z \sim \mathcal{N}(z; u, \sigma^2 I) \quad (\text{sample only 1})$$

$$\hat{x} = \text{Neural Net}(z; \theta)$$

$$J_1 \approx \log p_\theta(x|z)$$

$$\Leftrightarrow J_1 \approx \log \mathcal{N}(x; \hat{x}, I)$$

$$= \log \left( \frac{1}{\sqrt{(2\pi)^D}} \exp \left( -\frac{1}{2} (x - \hat{x})^T I^{-1} (x - \hat{x}) \right) \right)$$

$$= -\frac{1}{2} (x - \hat{x})^T (x - \hat{x}) + \log \frac{1}{\sqrt{(2\pi)^D}} \quad (I^{-1} = I, |I| = 1)$$

$$= -\frac{1}{2} \sum_{d=1}^D (x_d - \hat{x}_d)^2 + \log \frac{1}{\sqrt{(2\pi)^D}}$$

$\downarrow$  defined by  $\theta, \phi$        $\nwarrow$  constant

$$\arg \max_{\theta, \phi} J_1 = \arg \max_{\theta, \phi} -\frac{1}{2} \sum_{d=1}^D (x_d - \hat{x}_d)^2$$

$$= \arg \min_{\theta, \phi} \sum_{d=1}^D (x_d - \hat{x}_d)^2 : \text{reconstruction error}$$

$J_2$   
 - minimize  $J_2$  for maximizing ELBO  $(x; \phi, \theta)$

$$q_\phi(z|x) : \mathcal{N}(z; u, \sigma^2 I)$$

$$p(z) : \mathcal{N}(z; 0, I)$$

why  $\otimes$ ? :  $D_{KL}(q_\phi(z|x) || p(z))$  값도 해줘야 하는 것임.

$$J_2 = D_{KL}(q_\phi(z|x) || p(z))$$

$$= -\frac{1}{2} \sum_{h=1}^H (1 + \log \sigma_h^2 - u_h^2 - \sigma_h^2)$$

$q(z) = \mathcal{N}(z; u, \sigma^2 I), p(z) = \mathcal{N}(z; 0, I)$  (두 정규분포 Normal distribution 문제)  
 $D_{KL}(q||p) = -\frac{1}{2} \sum_{h=1}^H \left( 1 + \log \frac{\sigma_h^2}{\sigma_h^2} - \frac{(u_h - 0)^2}{\sigma_h^2} - \frac{\sigma_h^2}{\sigma_h^2} \right)$  ( $D_{KL}(p||q)$ 도 해줘야 하는 것임) ...  $\otimes$

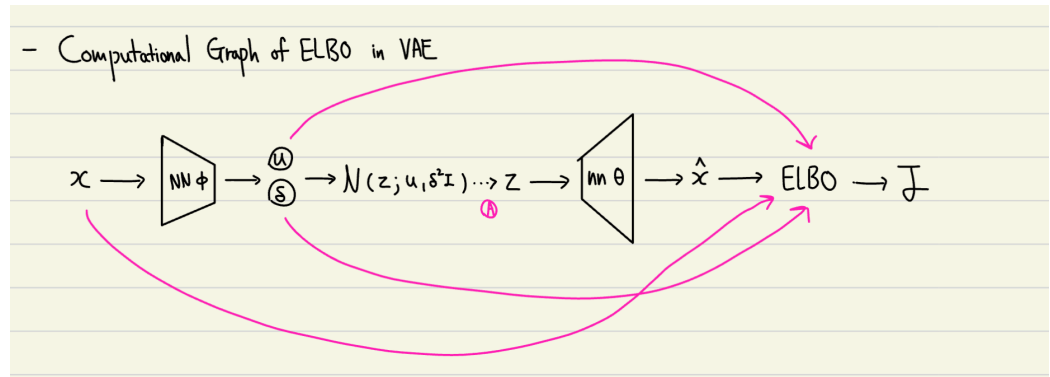
$$\text{minimize } J_2 \Leftrightarrow q_\phi(z|x) = p(z) \quad (\text{consistency / regularization term})$$

$$\therefore \text{ELBO}(x, \theta, \phi) \approx \underbrace{-\frac{1}{2} \sum_{d=1}^D (x_d - \hat{x}_d)^2}_{\text{reconstruction loss}} + \underbrace{\frac{1}{2} \sum_{h=1}^H (1 + \log \sigma_h^2 - u_h^2 - \sigma_h^2)}_{\text{regularization term}} + \text{const}$$

ELBO값을 최대화 시키기 위해서 원본 관측 데이터를 복원하면서도( $J_1$  : reconstruction Term),  $q(z|x)$  분포가  $z$ 의 prior인  $p(z)$ 에 가까워 지도록( $J_2$  : Regularization Term) 인코더와 디코더의 파라미터가 업데이트됨.

$J_1$ 의 값은 Monte Carlo 근사를 이용했기에(주어진  $x$ 에 대해  $z$ 를 한개만 샘플링), 결론적으로는 ELBO의 근사값을 구한 것임

따라서 VAE의 모델의 ELBO값을 구하기 위해서는 인코더와 디코더 구조가 모두 필요한 것임. Training 시에 ELBO값 최적화를 위한 계산 그래프는 다음과 같음.



$z$  sampling과정에서 gradient 연산이 끊기는 문제가 발생하기에 reparameterization trick을 도입함.