

Linear Regression & MLE

* Assumption: x, y in linear relationship (only 1 feature for each data x)

$$y = ax + b + \epsilon \quad (\epsilon \sim N(0, \sigma^2), \sigma \text{ is constant}) \quad \text{(Vector version for single data } x: y = W^T x + \epsilon \text{ where } x = [1 \ x_1 \ x_2 \ \dots \ x_n])$$

* Linear Regression Model

$$\hat{y} = Wx + b \quad \text{where } W, b \text{ are parameters}$$

$$\text{(Vector version: } \hat{y} = W^T x)$$

* MLE in Linear Regression? (For a single data y given x)

$$\epsilon = y - \hat{y} \sim N(0, \sigma^2) \quad (\text{where } \sigma \text{ is constant})$$

$$y \sim N(y; \hat{y}, \sigma) \Leftrightarrow y \sim p(y|x; W, b) \quad (\sigma \text{ is constant})$$

σ is constant and parameters W, b are to be estimated.

$$\log P_\theta(y) = \log p(y|x; W, b) = \log N(y; \hat{y}, \sigma)$$

$$= \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\hat{y})^2}{2\sigma^2}} \right)$$

$$= \log \frac{1}{\sqrt{2\pi}\sigma} + \log e^{-\frac{(y-\hat{y})^2}{2\sigma^2}}$$

$$= -\frac{1}{2\sigma^2}(y-\hat{y})^2 + \log \frac{1}{\sqrt{2\pi}\sigma}$$

$$= -\frac{1}{2\sigma^2}(y - Wx - b)^2 + \log \frac{1}{\sqrt{2\pi}\sigma}$$

does not change (since parameters are only W, b)

$$\begin{aligned} & \underset{w, b}{\operatorname{argmax}} \log N(y; \hat{y}, s) \\ \Leftrightarrow & \underset{w, b}{\operatorname{argmax}} \log P(y|x; w, b) \end{aligned}$$

$$\Leftrightarrow \underset{w, b}{\operatorname{argmax}} -\frac{1}{2s^2} (y - wx - b)^2$$

$$\Leftrightarrow \underset{w, b}{\operatorname{argmax}} -(y - wx - b)^2$$

$$\Leftrightarrow \underset{w, b}{\operatorname{argmin}} (y - wx - b)^2$$

$$\Leftrightarrow \underset{w, b}{\operatorname{argmin}} (y - (wx + b))^2 \text{ (square of residual!)} \leadsto \text{This is why objective function for Linear Regression leads to MSE Loss.}$$

For n datas?

$$L(w, b) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - (wx^{(n)} + b))^2 \text{ (MSE)}$$

object function

VAE : Variational Autoencoder

→ 주어진 data에 따라 유연하게 그 형태가 결정되는 확률분포 ($P_\theta(x)$ 를 신경망으로 구현)

* Review

purpose of generative models

- Single normal distribution (Multivariate)

$$P_\theta(x) = N(x; \theta) \quad \text{where } \theta = \{\mu, \Sigma\}$$

MLE for D

$$D = \{x^{(1)}, x^{(2)} \dots x^{(n)}\}$$

$$L(D; \theta)$$

$$= \log P_\theta(D)$$

$$= \log \left(\prod_{n=1}^N P_\theta(x^{(n)}) \right)$$

$$= \sum_{n=1}^N \log P_\theta(x^{(n)})$$

$$\text{Solve : } \frac{\partial}{\partial \theta} L(D; \theta) = 0$$

- GMM : Gaussian Mixture Model

$$\frac{\partial}{\partial \theta} \log_\theta(D) = 0 \rightarrow \text{해석적으로 풀기 힘들다 (} p(x, z) \text{에서 } p(x) \text{를 도출하기 어려움 } \sum_z \text{ marginalization이 필요하고 이는 log-sum은 } \frac{\partial}{\partial \theta} \log_\theta(D) = 0 \rightarrow \text{해석적으로 풀기 힘들다 (} p(x, z) \text{에서 } p(x) \text{를 도출하기 어려움 } \sum_z \text{ marginalization이 필요하고 이는 log-sum은 use ELBO, 유효함)}$$

$$\text{ELBO}(D; \theta, q) \leq \log P_\theta(D)$$

$$\Leftrightarrow \sum_{n=1}^N \sum_{z^{(n)}} q^{(n)}(z^{(n)}) \log \frac{P_\theta(x^{(n)}, z^{(n)})}{q^{(n)}(z^{(n)})} \leq \log P_\theta(D)$$

E-M with ELBO (θ, q update stepwisely)

E: $q^{(n)}(z^{(n)}) = P_\theta(z^{(n)} | x^{(n)})$ (for all n datas & for all K latent variables), ELBO log-likelihood approx.

M: $\frac{\partial}{\partial \theta} \text{ELBO} = 0 \rightarrow \text{maximize}$

latent variable (z) generation model

* Sampling (or generating) data with VAE (Decoder)

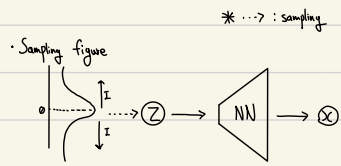
Parameter is fixed!

- ① Sample latent variable z from "fixed" normal distribution ($P(z) = N(z; 0, I)$)
- ② Transform latent variable z to observed data x using neural network (Decoder)

zero vector Identity Matrix

Where $z \in D^H$, fixed params are $\{0, I\}$

① $\therefore p(z) = N(z; 0, I)$



• In GMM z follows categorical distribution where z is discrete variable. In VAE, z is sampled from gaussian distribution ($N(z; 0, I)$) where z is continuous variable which makes representation broader.

→ 더 다양한 표현을 다양 표현 가능

population

• VAE is also a generative model which estimates the observed variable x 's distribution $P(x)$ based on sample.

Since VAE transforms z to x , this models probability $P(x|z)$

Since VAE decoder outputs \hat{x} from sampled z , define the x 's distribution as..

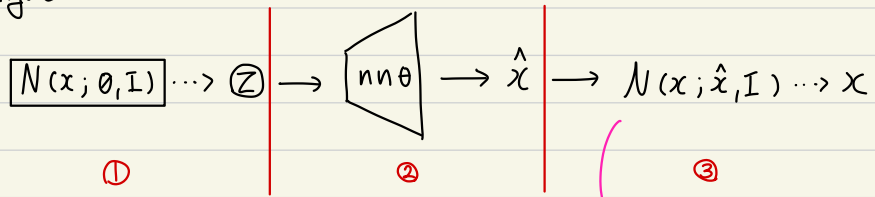
vector

② $\hat{x} = \text{NeuralNet}(z; \theta)$

Identity matrix

③ $P_{\theta}(x|z) = N(x; \hat{x}, I)$ ($x \sim N(x; \hat{x}, I)$)

* figure



→

If the observed data is binary $\{0, 1\}$
 $P_{\theta}(x|z)$ can be modeled as Bernoulli dist.

* Limitation of EM in VAE (Difficulty in E-step)

• ELBO EM Estep for n datas.

$$\text{renew } q^{(n)}(z^{(n)}) = p_{\theta}(z^{(n)} | x^{(n)})$$

$$\text{for VAE, } p_{\theta}(z^{(n)} | x^{(n)}) = \frac{p_{\theta}(x^{(n)}, z^{(n)})}{p_{\theta}(x^{(n)})}$$

$$= \frac{p_{\theta}(x^{(n)}, z^{(n)})}{\int p_{\theta}(x^{(n)}, z^{(n)}) dz} \rightarrow \text{marginalization}$$

Review) latent variable model $p(x, z)$ 에서 x 에 대한

$$\text{if } p(x) = \int p(x, z) dz \text{ 이다.}$$

$$p(x) = \int p(x, z) dz$$

$$= \int p(z) p(x|z) dz$$

$\int p_{\theta}(x, z) dz$ is easily countable in GMM since latent variable z is discrete.

$$\Leftrightarrow \sum_z p_{\theta}(x, z) dz$$

But in VAE z is continuous and z is vector $\rightarrow \int p_{\theta}(x, z) dz$ is impossible (or very hard)

\therefore VAE에서 EM의 E-step은 direct하게 적용할 수 X. \rightarrow Needs improvement in EM algorithm.

"How To Solve?"

* VAE training by improving EM algorithm (Encoder)

• Review : Derive ELBO from $\log P_\theta(x)$ using $q(z)$!

$\log P_\theta(x)$ (log likelihood for single data x)

$$= \int q(z) dz \log P_\theta(x) \quad (\int q(z) dz = 1)$$

$$= \int q(z) \log P_\theta(x) dz$$

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

$$= \int q(z) \log \frac{P_\theta(x, z)}{P_\theta(z|x)} \frac{q(z)}{q(z)} dz \quad \left(\frac{q(z)}{q(z)} = 1 \right)$$

$$= \int q(z) \left(\log \frac{P_\theta(x, z)}{q(z)} + \log \frac{q(z)}{P_\theta(z|x)} \right) dz$$

$$= \underbrace{\int q(z) \log \frac{P_\theta(x, z)}{q(z)} dz}_{\text{ELBO}(x; q, \theta)} + \underbrace{\int q(z) \log \frac{q(z)}{P_\theta(z|x)} dz}_{D_{KL}(q(z) \parallel P_\theta(z|x))}$$

$$\therefore \log P_\theta(x) = \int q(z) \log \frac{P_\theta(x, z)}{q(z)} dz + D_{KL}(q(z) \parallel P_\theta(z|x)) \quad (D_{KL} \geq 0)$$

$$\geq \int q(z) \log \frac{P_\theta(x, z)}{q(z)} dz \quad (\text{ELBO})$$

According to EM algorithm, E fixes θ and renews $q(z) = P_\theta(z|x)$, but according to previous page it is hard to get $P_\theta(z|x)$ since z is continual vector.

* $P(z)$ VS $q(z) (= q(z|x))$

for generative model $p(x, z)$ (with latent variable)

$$p(x) = \int p(x, z) dz = \int p(z) p(x|z) dz \dots \textcircled{A}$$

Since \textcircled{A} is not easily calculable and optimize (for D)

We use $q(z|x)$ (or $q(z)$) to approximate $p(z|x)$

• Comparison

| $p(z)$ | $p(z x)$ |
|--------------------------|--------------------------------|
| 사전 분포 (Prior) | 근사 후분포 (Variational Posterior) |
| 모델이 가정하는 z 에 대한 분포 | 데이터 x 에 대해 조건부인 정해진 분포 |
| z 에 대한 가설 / generation | |

* $P(z)$ VS $q(z) (= q(z|x))$

for generative model $p(x, z)$ (with latent variable)

$$p(x) = \int p(x, z) dz = \int p(z) p(x|z) dz \dots \textcircled{A}$$

Since \textcircled{A} is not easily calculable and optimize (for D)

We use $q(z|x)$ (or $q(z)$) to approximate $p(z|x)$

Comparison?

$p(z)$

- 사전 분포 (Prior)
- 모델이 가정하는 z 에 대한 분포
- z 에 대한 가정 / generation 시 이용하는 z 분포
 - 모델이 정한
- $\log p(x) = \log \int p(z) p(x|z) dz$ 가능

$p(z|x)$

- 근사 후분포 (Variational Posterior)
- 데이터 x 에 대해 조건부인 정해진 분포
- ELBO 개념은 이란 개념을 분포, $p(z|x)$ 를 근사하기 위해
- 사용하게 됨
- $ELBO = E_{q(z|x)}[\log p(x|z)] - D_{KL}[q(z|x) || p(z)]$ 를 가능
- VAE에서는 Encoder에 의해 가능

* Solution

- ① $q(z)$ 를 정해볼로 제안, 제1인 $q(z)$ 의 매개변수 $\psi = \{u, \Sigma\} \Leftrightarrow q_{\psi}(z) = \mathcal{N}(z; u, \Sigma)$
- ② $q(z)$ 를 정해볼로 제1인 상태에서 ELBO를 최대화한다.

$$\log P_{\theta}(x) = \underbrace{\int q_{\psi}(z) \log \frac{P_{\theta}(x, z)}{q_{\psi}(z)} dz}_{\text{ELBO}} + D_{KL}(q_{\psi}(z) \parallel P_{\theta}(z|x)) \dots \textcircled{1}$$

$$\theta, \psi = \underset{\theta, \psi}{\operatorname{argmax}} \text{ELBO}(x; \theta, \psi) = \underset{\theta, \psi}{\operatorname{argmax}} \int q_{\psi}(z) \log \frac{P_{\theta}(x, z)}{q_{\psi}(z)} dz$$

In VAE, cannot find $P_{\theta}(z|x)$ for $q_{\psi}(z) = P_{\theta}(z|x)$ with ψ

But can fit $q_{\psi}(z) = P_{\theta}(z|x)$ by maximizing ELBO (명세적으로 E step을 수행하지 않아도, ELBO를 maximize 하는 과정에서 $q_{\psi}(z)$ 가 $P_{\theta}(z|x)$ 에 fit된다)

Why? ①에서 좌항 $\log P_{\theta}(x)$ 는 ψ 를 가지고 있지 않음

$\therefore \psi$ 가 변해도 ELBO와 KL의 합인 $\log P_{\theta}(x)$ 는 변하지 X

$\Leftrightarrow \psi$ 에 대해 ELBO가 최대가 되면, KL항은 최소가 됨.

$\Leftrightarrow q_{\psi}(z) \approx P_{\theta}(z|x)$ 에 가까워짐. (E step 수행가능)

(간단한 확률분포 $q_{\psi}(z)$ 를 이용해 계산이 불가능한 $P_{\theta}(z|x)$ 를 근사시킴. 이를 variational approximation 또는 variational bayes 라고함.)

변분 근사

* For n datas?

$$\sum_{n=1}^N \text{ELBO}(x^{(n)}; \theta, \mathcal{H}^{(n)}) = \sum_{n=1}^N \int q_{\mathcal{H}^{(n)}}(z) \log \frac{p_{\theta}(x^{(n)}, z)}{q_{\mathcal{H}^{(n)}}(z)} dz$$

• prepare $q_{\mathcal{H}^{(n)}}(z)$ for each $x^{(n)}$ where $\mathcal{H}^{(n)} = \{u^{(n)}, \Sigma^{(n)}\}$... ①

• $\arg \max_{\mathcal{H}^{(n)}} \text{ELBO}(x^{(n)}; \theta, \mathcal{H}^{(n)}) \Leftrightarrow q_{\mathcal{H}^{(n)}}(z) \approx p_{\theta}(z | x^{(n)})$ (각 $\mathcal{H}^{(n)}$ 에 대해 ELBO 최대화)

"Prepare $\mathcal{H}^{(n)}$ for each data $x^{(n)}$? what if $n=1$ billion?"

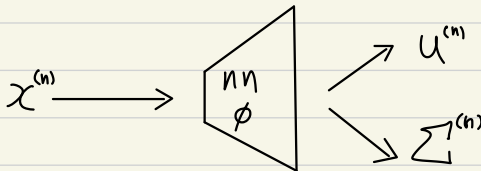


"Use Neural Network for extracting $\mathcal{H}^{(n)} = \{u^{(n)}, \Sigma^{(n)}\}$ for each $x^{(n)}$ "



"The role of Encoder in VAE"

* Encoder in VAE



↗ 분산 추정

• 근사치 분포 $q(z)$ 의 매개변수를 nn으로 출력하는 기법은 amortized inference라 함.

$$z^{(n)} \in D^H \rightarrow u^{(n)} \in D^H, \Sigma^{(n)} \in D^{H \times H}$$

Refine Σ to diagonal matrix

$$\Leftrightarrow \Sigma = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_H^2 \end{bmatrix}$$

elements are standard deviation of each feature

$$\Leftrightarrow \Sigma \text{ is expressible as } \sigma^2 I \text{ where } \sigma \in D^H, I = I_H$$

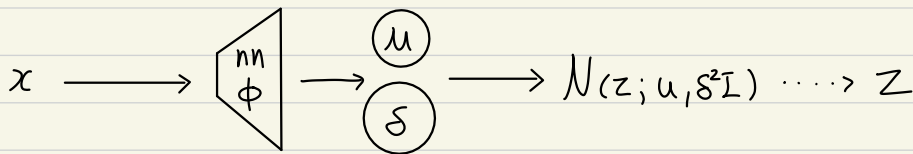
$$\therefore u, \sigma = \text{NeuralNetEncoder}(x; \phi)$$

$$q_\phi(z|x) = N(z; u, \sigma^2 I)$$

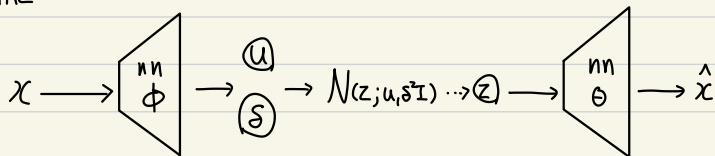
파라미터 ϕ 를 입력

* EM에서는 $x^{(n)}$ 에 대한 $q^{(n)}(z)$ 이고, VAE에서는 $q^{(n)}(z)$ 의 parameter를 x 로 부터 ∇ 를 통해 얻어내기, $q_\phi(z|x)$ 로 표현

• Encoder



• VAE



* ELBO Optimization in VAE

• Decoder

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

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

$$p_{\theta}(x|z) = \mathcal{N}(x; \hat{x}, I) \quad \text{where } \theta \text{ is decoder parameter}$$

• Encoder

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

$$q_{\phi}(z|x) = \mathcal{N}(z; u, \sigma^2 I) \quad \text{where } \phi \text{ is encoder parameter}$$

• ELBO

$$\begin{aligned} \text{Single Sample } x : \text{ELBO}(x; \theta, \phi) &= \int q_{\phi}(z|x) \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log p_{\theta}(x|z) dz + \int q_{\phi}(z|x) \log \frac{p(z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log p_{\theta}(x|z) dz - \int q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{p(z)} dz \\ &= \underbrace{E_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\mathcal{J}_1} - \underbrace{D_{KL}(q_{\phi}(z|x) \parallel p(z))}_{\mathcal{J}_2} \end{aligned}$$

$$D = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\} : \text{ELBO}(D; \theta, \phi) = \sum_{n=1}^N \text{ELBO}(x^{(n)}; \theta, \phi) = \sum_{n=1}^N \int q_{\phi}(z|x) \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} dz$$

$$\text{ELBO}(x; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)]}_{J_1} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{J_2} \rightarrow \begin{matrix} J_1 \text{는 } \phi, \theta \text{에 대해} \\ J_2 \text{는 } \phi \text{와 } p(z) \text{에 대해} \end{matrix}$$

• J_1

J_1 은 $q_{\phi}(z|x)$ 를 따르는 z 에 대한 $\log p_{\theta}(x|z)$ 의 Expectation

Monte carlo approximation - $q_{\phi}(z|x)$ 에서 z 를 "1"개만 sampling 해서 근사

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

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

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

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

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

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

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

$$= -\frac{1}{2} \sum_{d=1}^p (x_d - \hat{x}_d)^2 + \underbrace{\log \frac{1}{\sqrt{(2\pi)^p}}}_{\text{constant}}$$

↓
defined by θ, ϕ

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

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

• J_2

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

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

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

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

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

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

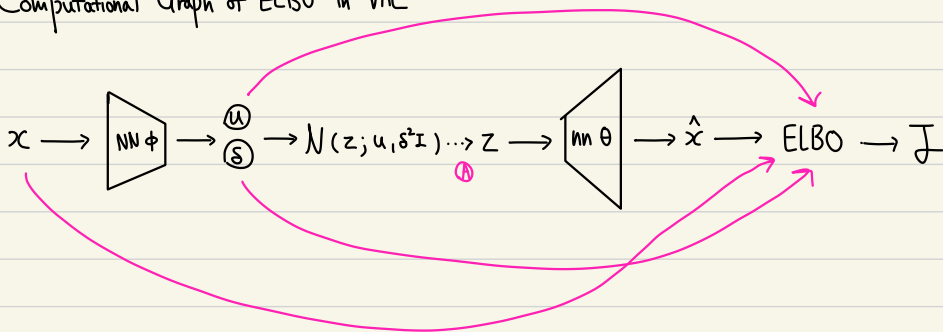
why ②? : $D_{KL}(q_\phi(z|x) \| p(z))$ 값을 해석적으로 구해 줌.

$$q(z) = \mathcal{N}(z; u_1, s_1^2 I), p(z) = \mathcal{N}(z; u_2, s_2^2 I) \text{ (두 정렬표가 Normal distribution일 때)}$$

$$D_{KL}(q \| p) = -\frac{1}{2} \sum_{h=1}^H \left(1 + \log \frac{s_{1,h}^2}{s_{2,h}^2} - \frac{(u_{1,h} - u_{2,h})^2}{s_{2,h}^2} - \frac{s_{1,h}^2}{s_{2,h}^2} \right) \text{ (} D_{KL}(p \| q) \text{도 해석적으로 구할 수 있음)} \dots \textcircled{1}$$

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

- Computational Graph of ELBO in VAE



- update θ, ϕ simultaneously with G.D for maximizing ELBO

→ θ 와 ϕ 를 각각 나눠서 θ, ϕ 를 각각 갱신함 (EM algorithm)

- issue : Gradient flow is unavailable in ①

* Reparameterization trick for solving issue

- Make sampling to enable gradient flow

$$Z \sim \mathcal{N}(Z; u, \sigma^2 I) \text{ sampling } \Leftrightarrow \epsilon \sim \mathcal{N}(\epsilon; 0, I)$$

$$Z = u + \sigma \odot \epsilon \quad \text{where } \odot \text{ is elementwise multiplication}$$

$$Z = u + \sigma \odot \epsilon$$

$$= \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_H \end{bmatrix} + \begin{bmatrix} \sigma_1 \\ \sigma_2 \\ \vdots \\ \sigma_H \end{bmatrix} \odot \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_H \end{bmatrix}$$

$$= \begin{bmatrix} u_1 + \sigma_1 \epsilon_1 \\ u_2 + \sigma_2 \epsilon_2 \\ \vdots \\ u_H + \sigma_H \epsilon_H \end{bmatrix}$$

• Without reparameterization trick

$$\begin{array}{l} \rightarrow \textcircled{u} \rightarrow \mathcal{N}(z; u, \sigma^2 I) \cdots \rightarrow Z \\ \rightarrow \textcircled{\sigma} \rightarrow \mathcal{N}(z; u, \sigma^2 I) \cdots \rightarrow Z \end{array}$$

• With reparameterization trick

$$\begin{array}{l} u \xrightarrow{\quad} \oplus \rightarrow Z \\ \sigma \xrightarrow{\quad} \odot \rightarrow \uparrow \oplus \\ \mathcal{N}(\epsilon; 0, I) \cdots \rightarrow \epsilon \xrightarrow{\quad} \odot \end{array}$$

"Gradient flow available for Encoder"