

# Understanding Variational Autoencoders

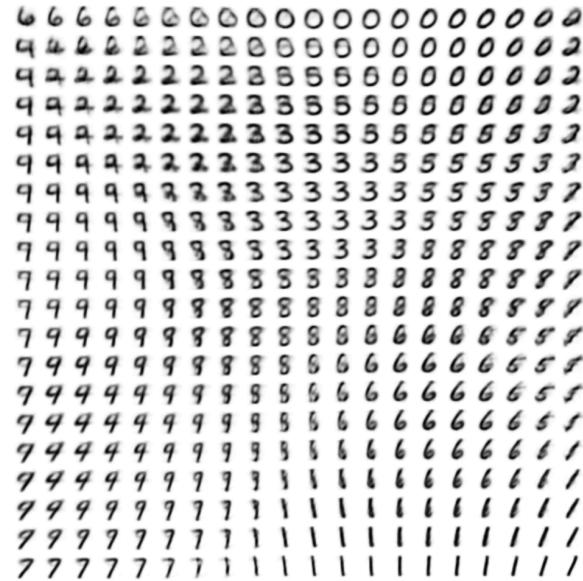
- ICLR 2014

- ## - Unsupervised generative model

- # - Generating samples



(a) Learned Frey Face manifold



(b) Learned MNIST manifold

Auto-Encoding Variational Bayes

**Diederik P. Kingma**  
Machine Learning Group  
Universiteit van Amsterdam  
[dpkingma@gmail.com](mailto:dpkingma@gmail.com)



**Max Welling**  
Machine Learning Group  
Universiteit van Amsterdam  
[welling.max@gmail.com](mailto:welling.max@gmail.com)



## Warm - up

---

- Supervised learning and unsupervised learning
  - Notations of Deep learning
- 

$\mathbf{x}$  input (image or tabular data)

---

$y(\mathbf{x})$  or  $\hat{y}(\mathbf{x})$  output of model

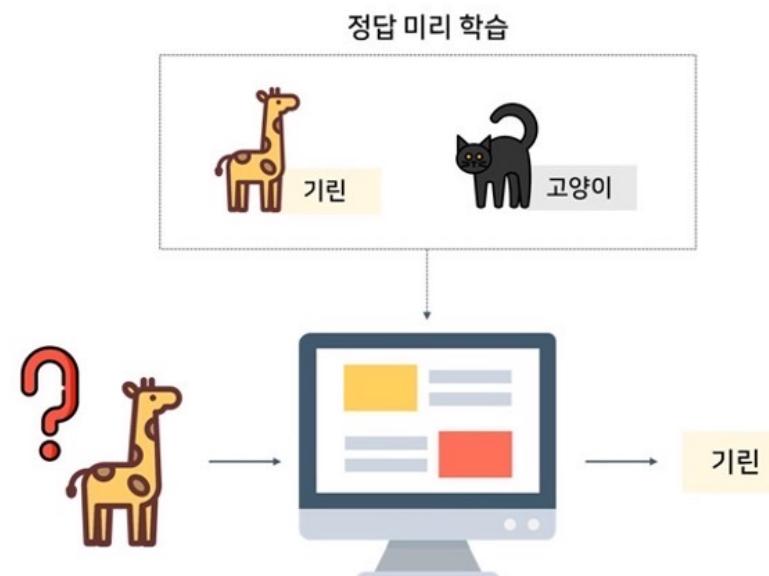
---

$t$  or  $\hat{t}$  target value

---

## Supervised learning

- 학습 데이터에 정답 or 라벨 **가** 주어진 경우
- \* 와 # 설명하는 함수 찾기
- goal :  $y(x) \sim \#$  인  $y(x)$  찾기
- e.g.: regression, classification



# Unsupervised learning

- 데이터의 특징이나 구성 찾는 문제
- 학습 데이터에 x 만 주어진 경우
- e.g: clustering, feature extraction, dimensionality reduction

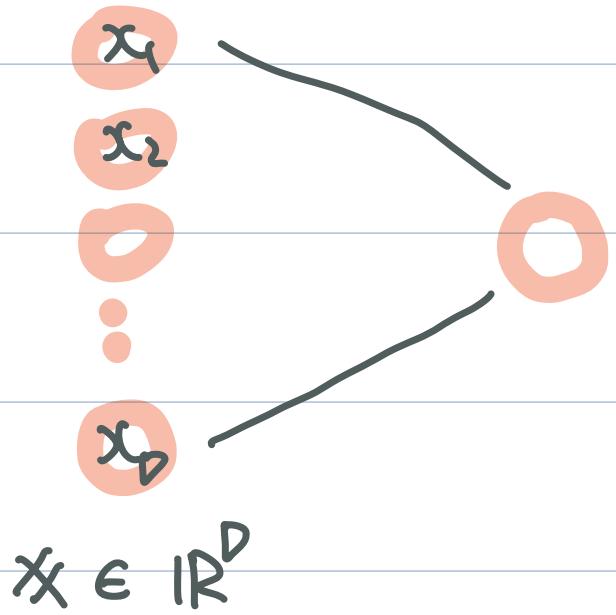


# Notations of Deep learning

- $\mathbf{x}$  input (image or tabular data)
- $y(\mathbf{x})$  or  $\hat{y}(\mathbf{x})$  output of model
- $t$  or  $\hat{t}$  target value
-  Perceptron, activation function, weight, bias
- Input, hidden, output layer, rule of activation functions
- cost (loss) function, optimization algorithm (method)

backpropagation algorithm





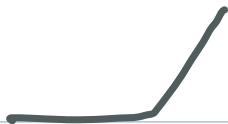
$w \in \mathbb{R}^D$ ,  $b \in \mathbb{R}$   
weight bias

$w^T x + b$  perceptron

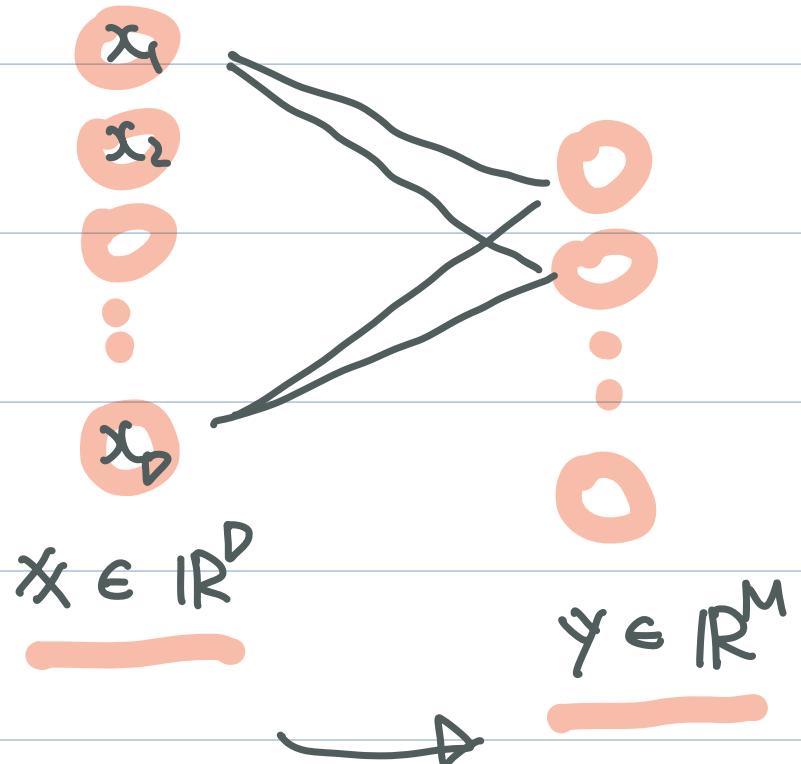
$f$ : non-linear function

is called activation function

Sigmoid, ReLU, tanh, ....



$$y(x) = \max(0, x)$$



$w \in \mathbb{R}^D$ ,  $b \in \mathbb{R}$

weight

bias

$$w^T x + b$$

$W : D \times M$

$b \in \mathbb{R}^M$

weight  
matrix

$$f(x^T w + b)$$

$$\frac{1}{n} \sum (y(x) - t)^2$$

Q1

$x$  input

$t$  target

$x$   
 $\mathbb{R}^5$

$\mathbb{R}^4$

$\mathbb{R}^4$

$\mathbb{R}^3$

$\mathbb{R}$

$$y(x) \sim t$$

# Autoencoder

- $y(x) \sim x$  가 되게끔 cost 정의

e.g.  $\sum_n \{y(x_n) - x_n\}^2$

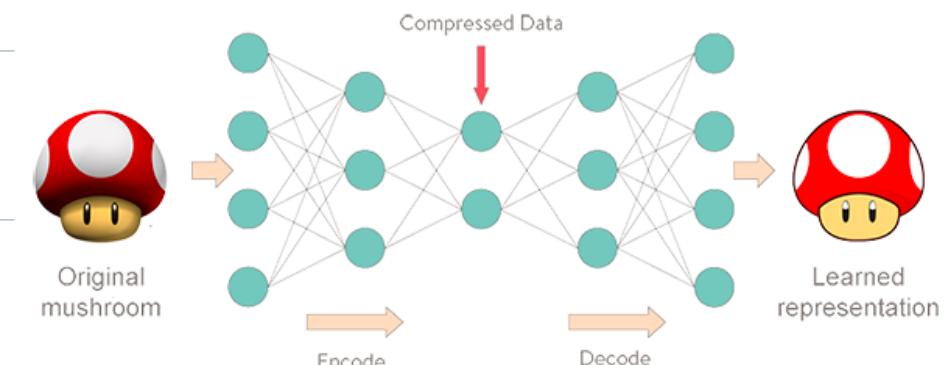
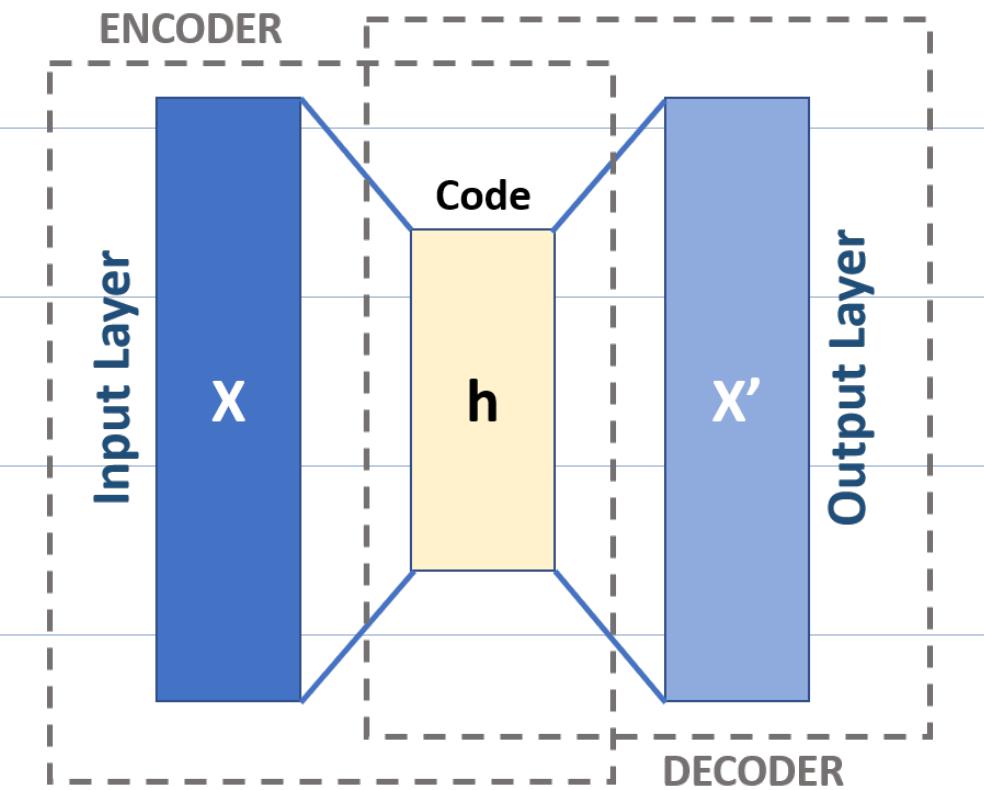
- 일반적으로 encoder와 decoder 모양은

대칭으로 구성

- 구 : latent variable, bottleneck layer

feature or hidden representation.

- Nonlinear identity function



## Applications of Autoencoder

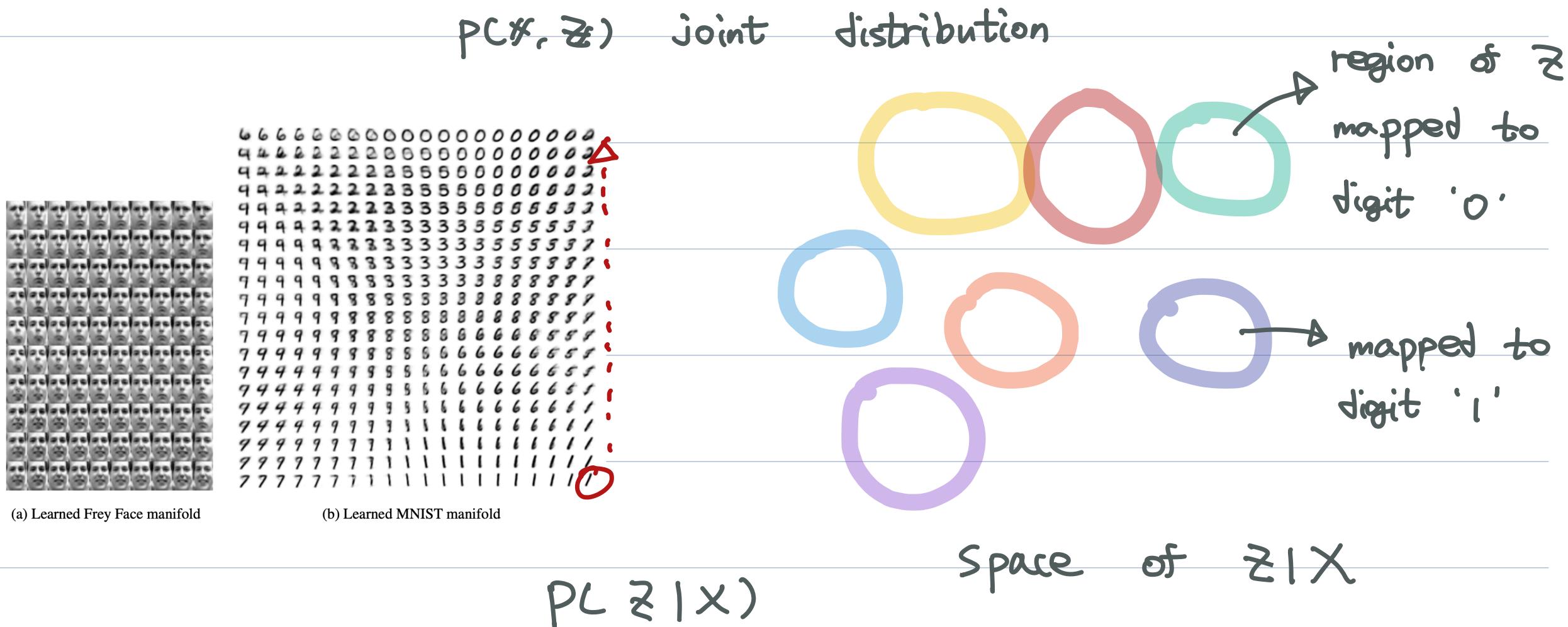
---

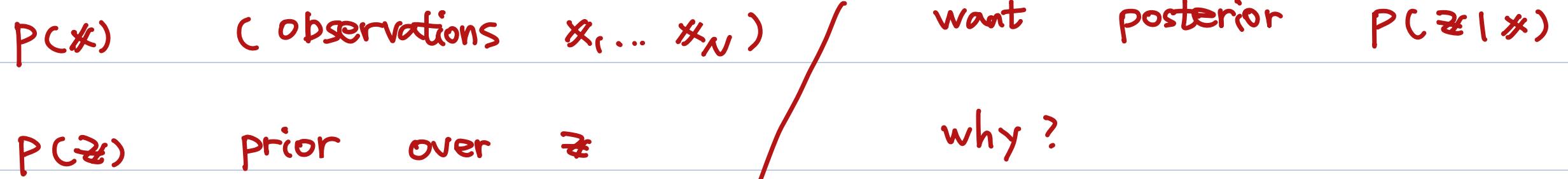
- Feature extract , dimension reduction
  - Anomaly detection (reconstruction based)
  - Fine tuning, denoising
- 
- 
- 
-

# Variational

# Autoencoder (VAE)

- Generative modeling: deal with models of distribution  $p(x)$ .
- Latent variable  $z$ ,  $p(z)$  is defined over  $\mathcal{Z}$





We want to know the posterior  $p(z|x)$  so we  
 (never)  
 will approximate it with variational distribution  $q_\theta(z|x)$

Let  $q_\theta(z|x)$  be a Gaussian distribution as

$$q_\theta(z|x) = N(\underbrace{z| \mu_\theta(x)}, \underbrace{\sigma_\theta(x)^2 I}_{(0 \dots 0)})$$

D-dim

$\mu_\theta$  and  $\sigma_\theta$  is some neural network function of  $x$   
 parametrized by  $\theta$ .

Goal : Maximize Log Marginal likelihood of given data

i.e. maximize  $\frac{1}{N} \sum_n \log P(x_n)$   $(x_n)_{n=1, 2, \dots, N}$

Step 1.

$$\log P(x) = \log \frac{P(z) P(x|z)}{P(z|x)} = \log P(z) + \log P(x|z) - \log P(z|x)$$

Step 2. (approximating posterior  $q(z|x)$  over  $z$ .)

$$\log p(x) = \underbrace{\log p(x)}_{=1} \int q(z|x) dz$$

$$= \int q(z|x) \{ \log P(z) + \underbrace{\log P(x|z) - \log p(z|x)}_{\text{arrow}} \} dz$$

$$= E_{q(z|x)} [\log p(x|z)] + \int q(z|x) \log P(z) dz$$

$$- \int q(z|x) \log p(z|x) dz$$

$$\int f(x) p(x) dx$$

$$= E_{p(x)} [f(x)]$$

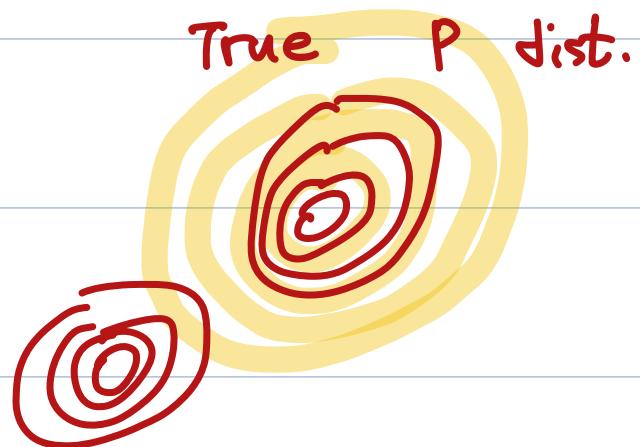
### Step 3

$$\log P(x) \pm \int q(z|x) \log q(z|x) dz$$

???

$$= E_{q(z|x)} [\log P(x|z)] - KL(q(z|x) || P(z)) + KL(q(z|x) || P(z|x))$$

$$(KL(P||Q) := - \int P(z) \log \left\{ \frac{Q(z)}{P(z)} \right\} dz)$$



reverse  $KL - D$  :  $KL(Q || P)$   
true  
approximation

$KL(P || Q)$   $\geq$   $KL(Q || P)$  " unimodal Gaussian

## Clean - up

$$\log P(x)$$

$$= \mathbb{E}_{q(z|x)} [\log P(x|z)] - KL(q(z|x) || P(z)) + \underbrace{KL(q(z|x) || P(z|x))}_{\geq 0}$$

$$\geq \mathbb{E}_{q(z|x)} [\log P(x|z)] - KL(q(z|x) || P(z))$$

$$\approx \frac{1}{J} \sum_{j=1}^J \log P(x|z_j) \quad z_j \sim q(z|x)$$

$x$

$p(z|x)$

$P(x|z)$

1. Maximize  $p(x)$

2. Introduce prior  $p(z)$  and joint distribution  $p(x, z)$

3. Introduce approximating posterior  $q(z|x) = N(z|M(x), \sigma(x)I)$

$p(z|x)$

## Actual calculation

- Encoder (approximating posterior) :  $q(z|x) := N(z| \mu(x), \sigma(x) I)$

where  $\mu, \sigma$  are neural net. ( $D$ -dim random vector)

- Prior  $p(z) = N(z| 0, I)$

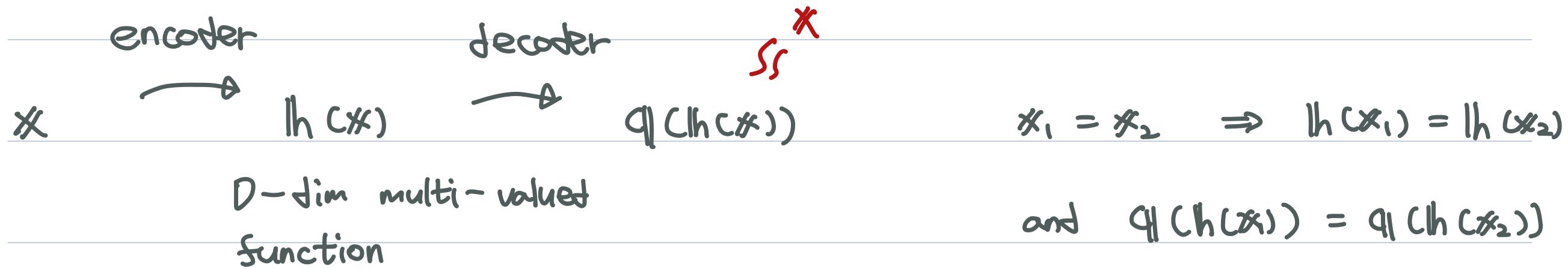
- Likelihood / Decoder :  $p(x|z)$  depending on input data

- Since  $p(z)$  and  $q(z|x)$  are Gaussian,

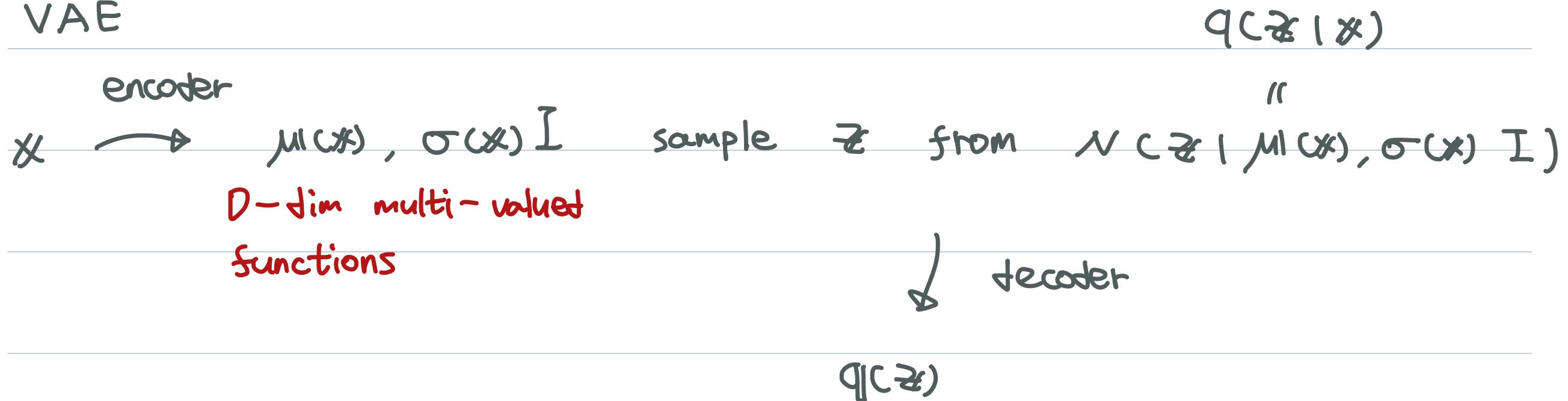
$$KL(q(z|x) || p(z)) = \frac{1}{2} \sum_{d=1}^D \left\{ 1 + \log(\sigma_d^2(x)) - \mu_d^2(x) - \sigma_d^2(x) \right\}$$

where  $\mu(x) = (\mu_1(x), \dots, \mu_D(x))^T$ ,  $\sigma(x) = (\sigma_1(x), \dots, \sigma_D(x))^T$

# Autoencoder



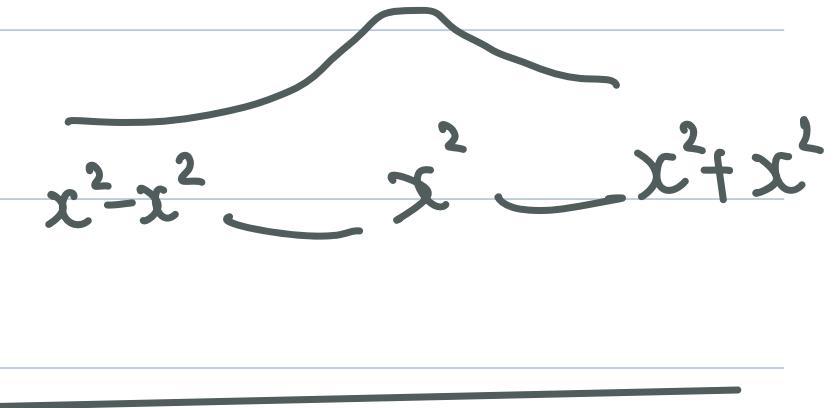
# VAE



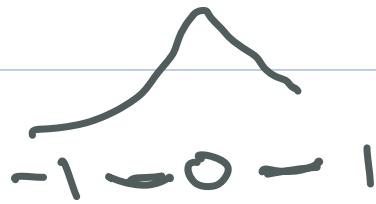
$x$   $\xrightarrow{\quad}$   $N(\mu(x^2), \sigma^2(x^4))$   $\xrightarrow{\quad}$   $z$  sample  
 $z \sim N(\mu(x^2 + x^4 t), \sigma^2(x^2 + x^4 t))$

$$\begin{aligned} \mu(x) &= x^2 \\ \sigma^2(x) &= x^4 \end{aligned}$$

$$x \in D(\mu) \cup \{0\}$$



$$t \sim N(t|0, 1)$$



$$z := x^2 + x^4 t$$

$$w_i^{t+1} = w_i^t - \eta \nabla_{w_i} \text{cost}$$

## Reparametrization trick

If we sample  $\tilde{z}$  from  $N(\tilde{z} | \mu(x), \sigma(x) I)$ , then we lose the information about weights and bias of  $\mu, \sigma$ .

I.e. Sampling operation is non-differentiable.

So we sample  $t$  from  $N(t | 0, I)$  and let

$$\tilde{z} := \mu(x) + \sigma(x) t I$$

$$z \sim N(z | x^2, x^4)$$

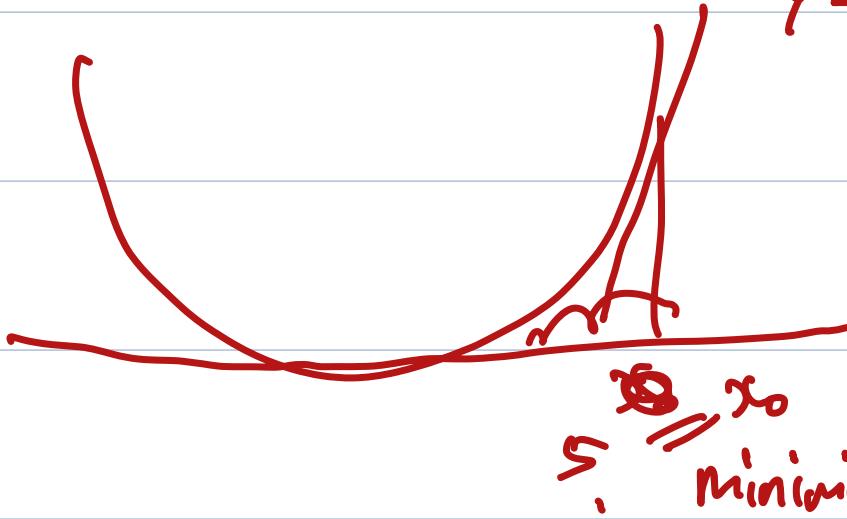
$$\frac{z - x^2}{x^2} \sim N(0, 1)$$

$$t \sim N(t | 0, 1)$$

$$z_i = x^2 + t x^2$$

$$t = \frac{z - x^2}{x^2}$$

$$y = f(x) = x^2$$



Want  $f(x) \approx s(x)$ .

s. minimize  $x^2$

$$x^{t+1} = x^t - \eta \cdot f'(x^t)^{0.1}$$

$2x$

$$x_1 = 5 - 0.1 \cdot 10 = 4$$

$$x_2 = 4 - 0.1 \cdot 8 = 3.2$$

