

$x^1, x^2, \dots, x^m$

$$P_\theta(x) = \int P_\theta(x, z) dz$$
$$= \int P_\theta(x|z) p(z) dz$$

log-likelihood

$$\ell(\theta) = \sum_{i=1}^m \log P_\theta(x^i) = \sum_{i=1}^m \log \int P_\theta(x^i, z) dz$$

Hard!

ELBO

~~$L(q_i, x^i, \theta)$~~

$$L(q_i, x^i, \theta) = \sum q_i(z|x^i) \log \frac{P_\theta(x^i, z)}{q_i(z|x^i)}$$
$$= \log P_\theta(x^i) - KL(q_i(z|x^i) || P_\theta(z|x^i))$$

$$\max_{\theta} \ell(\theta)$$
$$= \max_{\theta} \sum_{i=1}^m \log P_\theta(x^i)$$



$$\max_{\theta, q_1, q_2, \dots, q_m} \sum_{i=1}^m L(q_i, x^i, \theta)$$

start from some  $\theta, q_1, q_2, \dots, q_m$

Loop

$$\nabla_{\theta} = \frac{\partial}{\partial \theta} \sum_{i=1}^m L(q_i, x^i, \theta)$$

$$\theta \leftarrow \theta + \lambda \nabla_{\theta}$$

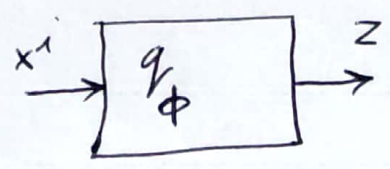
for  $i = 1 \dots m$ :

$$q_i \leftarrow \operatorname{argmax}_q L(q, x^i, \theta)$$

what if we have millions of data?

solution 1  $q_i(z|x^i) \rightarrow q(z)$  X

solution 2:  $q_i(z|x^i) \rightarrow q_{\phi}(z|x^i)$



Amortized Variational Inference

# Amortized VI

start from some  $\theta, \phi$

Loop

$$\nabla_{\theta} = \frac{\partial}{\partial \theta} \sum_{i=1}^m \mathcal{L}(x^i; \phi, \theta)$$

$$\mathcal{L}(x^i; \phi, \theta)$$

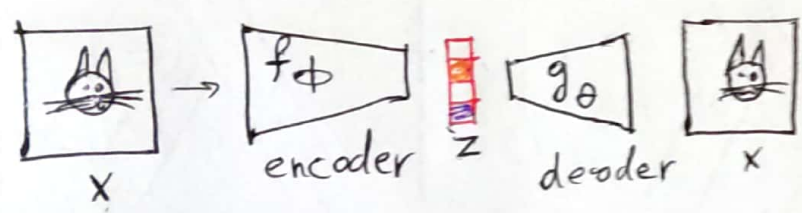
$$\nabla_{\phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^m \mathcal{L}(x^i; \phi, \theta)$$

$$= \sum_z q_{\phi}(z|x^i) \log \frac{p_{\theta}(z, x^i)}{q_{\phi}(z|x^i)}$$

$$\theta \leftarrow \theta + \lambda \nabla_{\theta}$$

$$\phi \leftarrow \phi + \lambda \nabla_{\phi}$$

## Auto-encoders



$$z = f_{\phi}(x)$$

$$y = g_{\theta}(z)$$

find  $\phi, \theta$  such that  $\|y - x\|$  is small

$$\min_{\theta, \phi} \sum_{i=1}^m \|g_{\theta}(f_{\phi}(x^i)) - x^i\|^2$$

- 1- data compression
- ↳ 2- dimensionality reduction
- ↳ 3- feature extraction (unsupervised), representation learning
- ↳ 4- semi-supervised learning

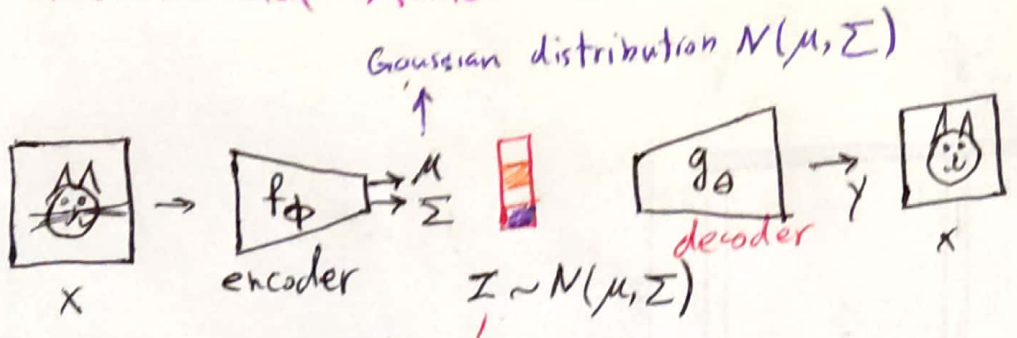
## generating data

choose a random  $z$  (e.g.  $z \sim \mathcal{N}(0, I)$ )

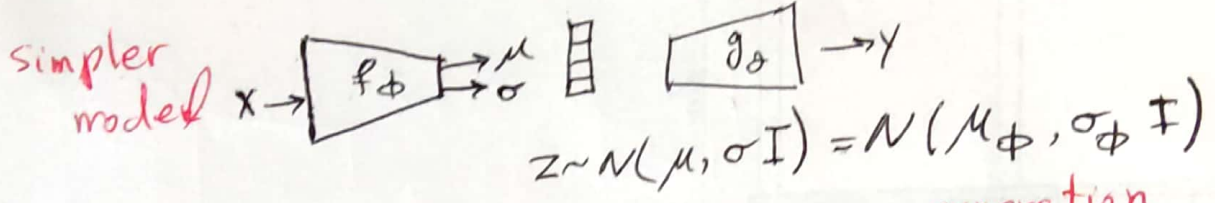
$$\hat{x} = g_{\theta}(z)$$

doesn't work because the distribution of  $z$  is unknown  $\neq \mathcal{N}(0, I)$

# Variational Auto-Encoders



a sample from  $N(\mu, \Sigma)$



generation

$$z \sim N(0, I)$$

$$y = g_\theta(z)$$

$$\mu, \sigma = f_\phi(x)$$

$$z \sim N(\mu, \sigma I)$$

$$\hat{x} = g_\theta(z)$$

$x^1, x^2, \dots, x^m$

$$C(\theta, \phi) = \sum_{i=1}^m \alpha \| g_\theta(z \sim N(\mu_\phi(x^i), \sigma_\phi(x^i) I)) - x^i \|^2 + \sum_{i=1}^m KL(N(\mu_\phi(x^i), \sigma_\phi(x^i) I) \| N(0, I))$$

How to compute  $\frac{\partial}{\partial \theta} C(\theta, \phi)$ ,  $\frac{\partial}{\partial \phi} C(\theta, \phi)$ ?

$$\mu, \sigma = f_\phi(x) \quad \checkmark$$

$$z \sim N(\mu, \sigma I) \quad ?$$

$$y = g_\theta(z) \quad \checkmark$$

How to ~~sample~~ differentiate from the sampling operation?



# Reparameterization trick

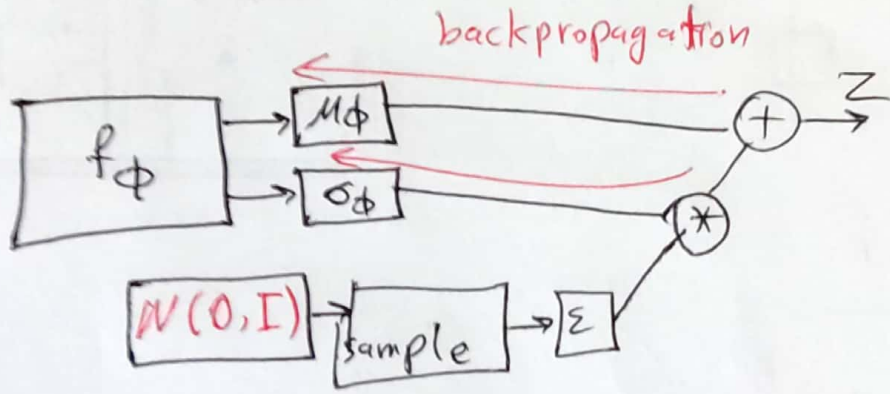
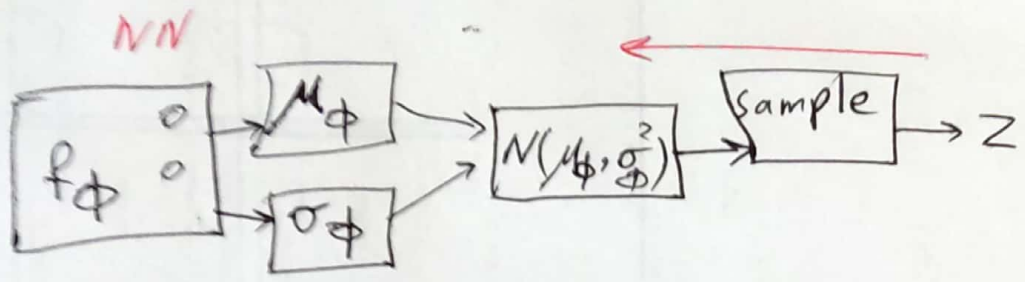
$$z \sim \mathcal{N}(\mu_\phi, \sigma_\phi^2 \mathbf{I})$$



$$z = \mu_\phi + \epsilon \sigma_\phi$$

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

$\epsilon \sim \mathcal{N}(0, \mathbf{I})$   
 $\sigma \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$   
 $\mu + \sigma \epsilon \sim \mathcal{N}(\mu, \sigma^2 \mathbf{I})$



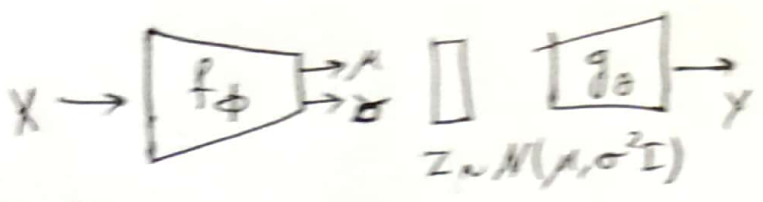
$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

$$z = \mu_\phi + \epsilon \sigma_\phi$$

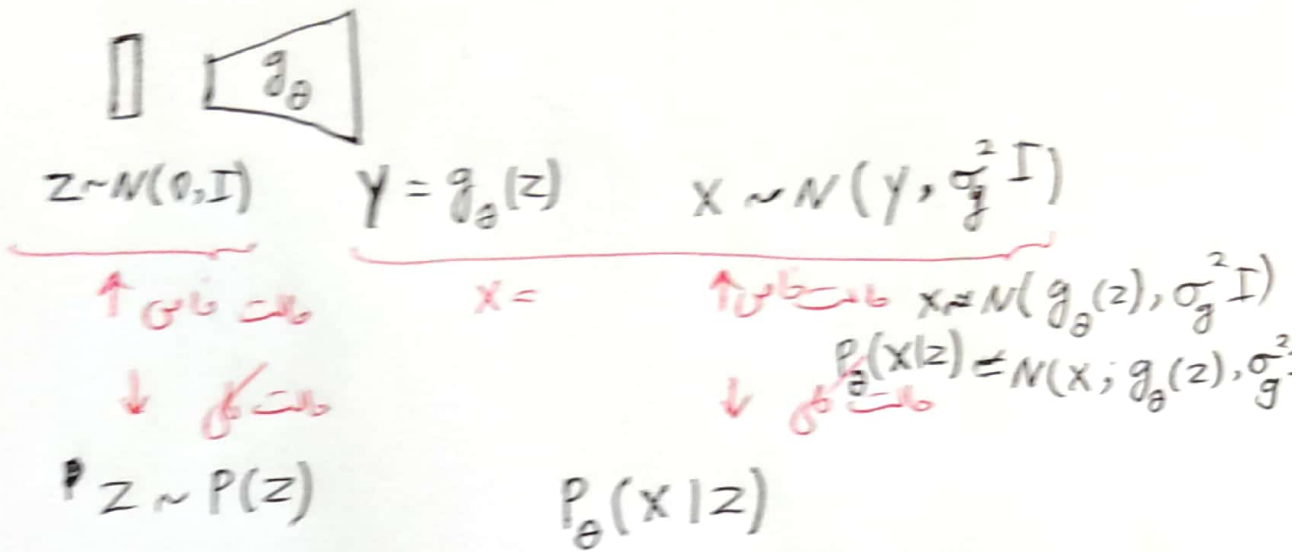
$$\frac{\partial z}{\partial \phi} = \frac{\partial \mu_\phi}{\partial \phi} + \epsilon \frac{\partial \sigma_\phi}{\partial \phi} \quad \checkmark$$

$\epsilon^i \sim \mathcal{N}(0, \mathbf{I})$

$$C(\theta, \phi) = \sum_{i=1}^m \alpha \left\| g_\theta(\mu_\phi(x^i) + \epsilon^i \sigma_\phi(x^i) \mathbf{I}) \right\|^2 + \text{KL} \left( \mathcal{N}(\mu_\phi(x^i), \sigma_\phi^2(x^i) \mathbf{I}) \parallel \mathcal{N}(0, \mathbf{I}) \right)$$



generation



↑  $z \sim N(0, I)$   $x =$   $\uparrow$   $x \sim N(g_\theta(z), \sigma_g^2 I)$   
 ↓  $P(z)$   $\downarrow$   $P_\theta(x|z) \sim N(x; g_\theta(z), \sigma_g^2)$

$P(z) \sim P(z)$        $P_\theta(x|z)$

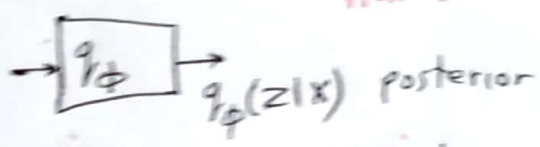
$P_\theta(x, z) = P_\theta(x|z) P(z)$

data  $x^1, x^2, \dots, x^m$

$\max_{\theta} \ell(\theta) = \sum_{i=1}^m \log P_\theta(x^i) = \sum_{i=1}^m \log \int P_\theta(x, z) dz$   
 log-likelihood =  $\sum_{i=1}^m \log \int P_\theta(x^i|z) P(z) dz$

Hard to maximize

ELBO (Approximated)

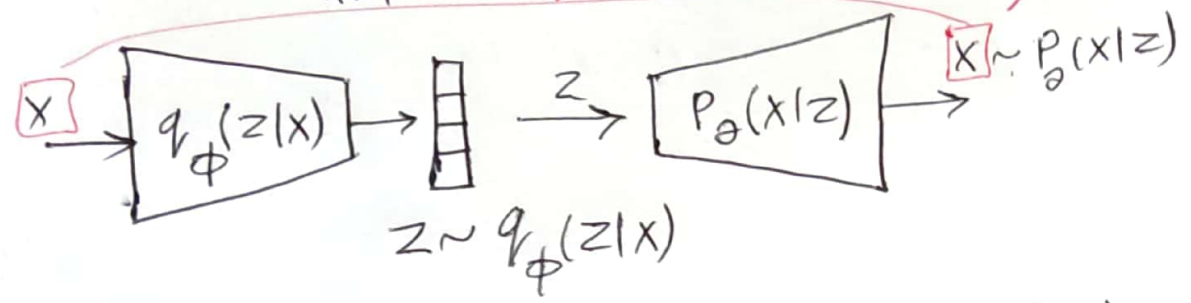


$\max_{\phi, \theta} \mathcal{L}(\phi, \theta) = \sum_{i=1}^m \int q_\phi(z|x^i) \log \frac{P_\theta(x^i, z)}{q_\phi(z|x^i)}$   
 =  $\sum_{i=1}^m \int q_\phi(z|x^i) \log \frac{P_\theta(x^i|z) P(z)}{q_\phi(z|x^i)}$

=  $\sum_{i=1}^m \int q_\phi(z|x^i) \log P_\theta(x^i|z) - \int q_\phi(z|x^i) \frac{q_\phi(z|x^i)}{P(z)}$   
 $\sum_{i=1}^m E_{q_\phi(z|x^i)} \left\{ \log P_\theta(x^i|z) \right\} - KL(q_\phi(z|x^i) \parallel P(z))$   
 $x \rightarrow$  enc  $q(x|z)$       dec  $p(z|x)$

$$\mathcal{L}(\phi, \theta) = \sum_{i=1}^m E_{q_{\phi}(z|x^i)} \left\{ \log p_{\theta}(x^i|z) \right\}$$

$$- \sum_{i=1}^m \text{KL} \left( q_{\phi}(z|x^i) \parallel p(z) \right)$$



صوب صوب

$$q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x^i), \sigma_{\phi}(x^i)^2 I) \quad \epsilon \sim \mathcal{N}(0, I)$$

$$z \sim q_{\phi}(z|x) \Rightarrow z = \mu_{\phi}(x^i) + \epsilon \sigma_{\phi}(x^i)$$

$$p_{\theta}(x|z) = \mathcal{N}(x; g_{\theta}(z), \sigma_g I)$$

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

$$\mathcal{L}(\phi, \theta) = \sum_{i=1}^m E_{q_{\phi}(z|x^i)} \left\{ \log \mathcal{N}(x; g_{\theta}(z), \sigma_g I) \right\}$$

$$- \sum_{i=1}^m \text{KL} \left( \mathcal{N}(\mu_{\phi}(x^i), \sigma_{\phi}(x^i)^2 I) \parallel \mathcal{N}(0, I) \right)$$

$$= \sum_{i=1}^m E_{q_{\phi}(z|x^i)} \left\{ \log \frac{1}{\sqrt{2\pi} \sigma_g^{d/2}} e^{-\frac{\|x - g_{\theta}(z)\|^2}{\sigma_g^2}} \right\} - \text{KL}$$

$$= \sum_{i=1}^m E_{q_{\phi}(z|x^i)} \left\{ -\frac{d}{2} \log \sigma_g + \frac{\|x - g_{\theta}(z)\|^2}{\sigma_g^2} \right\} - \text{KL}$$

one sample approximation  $z^i \sim q_{\phi}(z|x^i)$

$$\approx \sum \frac{\|x^i - g_{\theta}(z^i)\|^2}{\sigma_g^2} + \sum \text{KL}$$