

# VAE using Expectation Maximization strategy.

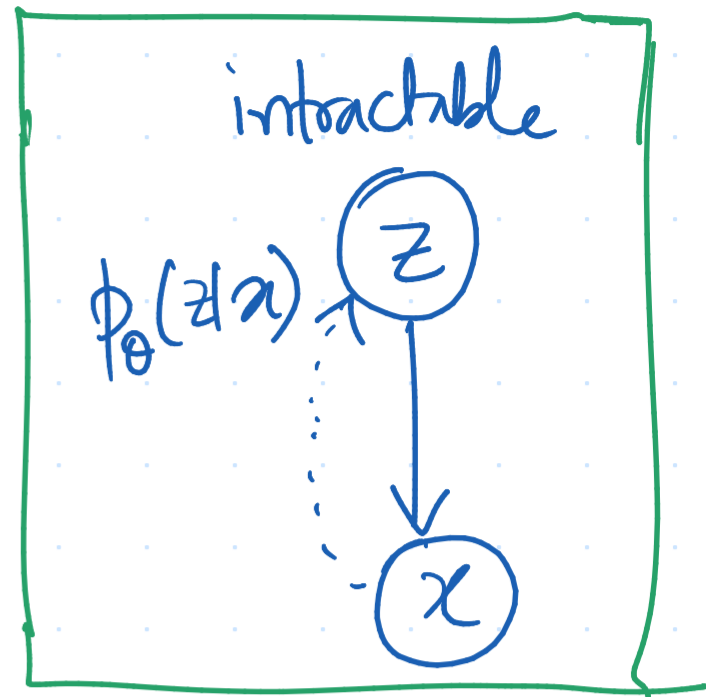
16/3/25

## M.C.M.C. Expectation Maximization.

Assumption :-

MCMC sampling

We can calculate  $p(z|x^{(i)}) \rightarrow$  this may not hold in practice.



E-step : For all  $i$

$$\text{set } Q_i^{(t)}(z) = p(z|x^{(i)}; \theta^{(t)})$$

M-step

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{i=1}^n \sum_z Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z)}$$

← doesn't depend on  $\theta$ .

$$= \arg \max_{\theta} \sum_{i=1}^n \sum_z Q_i^{(t)}(z^{(i)}) \log p(x^{(i)}, z^{(i)}; \theta)$$

$$z \sim N(0, I_{k \times k})$$

$$x|z \sim \text{NeuralNet}_{\theta}(z^{(i)}) \sim \text{Complex}(z^{(i)}, \theta)$$

(same as maximizing this quantity).

$p(z|x) \rightarrow$  impossible, cannot be calculated, very complex

→ intractable.

$$\log\left(\frac{A}{B}\right) = \log A - \log B \xrightarrow{\text{constant}} \text{w.r.t. } \theta$$

hence, we get rid of this form.

$$\Rightarrow \arg \max_{\theta} \sum_{i=1}^n \sum_{z} Q_i^{(t)}(z^{(i)}) \log p(x^{(i)}, z^{(i)}; \theta).$$

$$\Rightarrow \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q_i^{(t)}} [\log p(x^{(i)}, z^{(i)}; \theta)]$$

$Q_i^{(t)}(z)$  in E-step is only used to calculate this expectation.

We are not interested in the density of  $Q$ . We are only interested in the density value of  $Q$  to perform this expectation.

$$\approx \arg \max_{\theta} \sum_{i=1}^n \frac{1}{M} \sum_{m=1}^M \log p(x^{(i)}, z^{(m)}; \theta)$$

(replace with Monte-Carlo estimation of the estimation)

$$z_i^{(m)} \sim Q_i^{(t)}$$

$z_i^{(m)}$  = sampled from  $Q_i^{(t)}$ .

## Q-posterior

Sample from posterior of a complex probability distribution even though we don't know how to evaluate it.

Law of large numbers tells us that as the  $M$  goes to  $\infty$  the Monte-Carlo estimate will tend to converge towards the true expectation.

Gibbs Sampling  
Metropolis Hastings.

We don't know the posterior, but we can approximate the posterior using the Monte-Carlo Technique / Sampling Technique.

Here the guarantee of the likelihood increase after every step doesn't hold anymore, because this is an approximation of the lower bound and not the exact lower bound.

So, there are three ways to do this :-

\* Exact posterior  $\rightarrow$  maths / calculation etc.  $\rightarrow$  (diff)

\* Gibbs / Sampling methods  $\rightarrow$   $\checkmark$

\* Variational Inference  $\rightarrow$  optimization  $\rightarrow$  (previous class).

Jensen's Inequality

$$\log p(x) \geq \underbrace{\text{ELBO}(x; \theta)}$$

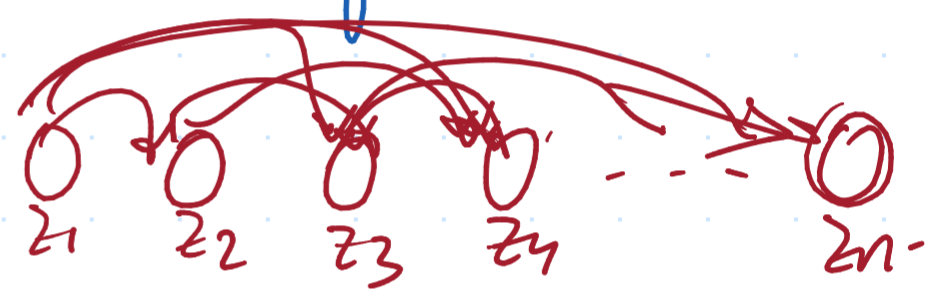
$\rightarrow$  lower bound of the evidence.

$$\log p(x) = \text{ELBO}(x; \mathcal{Q}) + \underbrace{\quad}_{?}$$

$$D_{KL}(\mathcal{Q} | p(z|x)) = \log p(x) - \text{ELBO}(x; \mathcal{Q})$$

$$\log p(x) = \underbrace{\text{ELBO}(x; \mathcal{Q})}_{\text{const. w.r.t. } \mathcal{Q}} + \underbrace{D_{KL}(\mathcal{Q} || p_{z|x})}_{\substack{\text{maximize} \\ \text{w.r.t. } \mathcal{Q}.}} \geq 0.$$

$$p_{z|x} = \underset{q \in \mathcal{Q}}{\text{argmax}} \text{ELBO}(x; q) \leftarrow \text{Variational Inference.}$$



Mean field Assumption:-

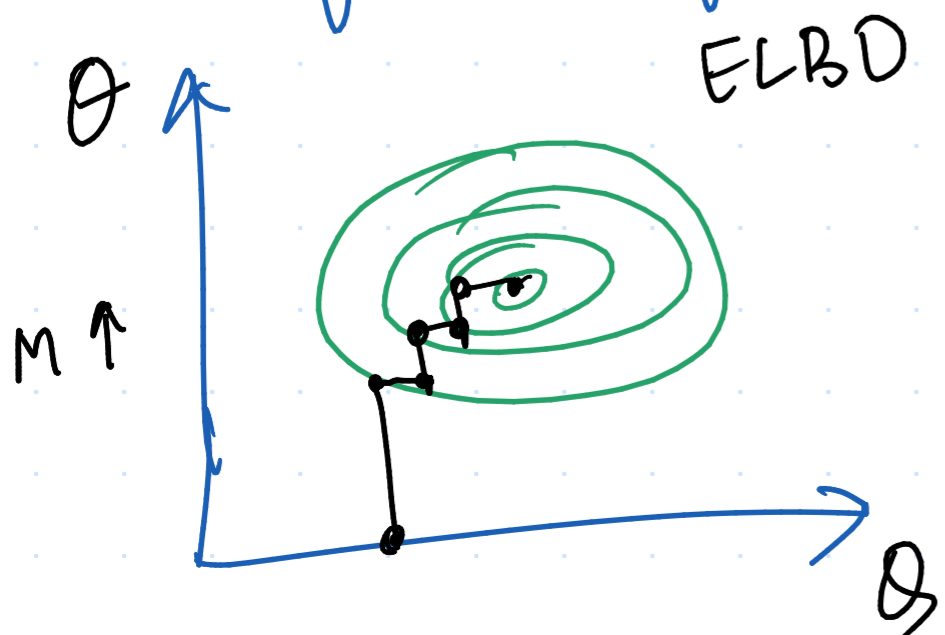
$$\mathcal{Q}(z) = \mathcal{Q}_1(z_1) \mathcal{Q}_2(z_2) \mathcal{Q}_3(z_3) \dots \mathcal{Q}_k(z_k) \quad \text{O O O O}$$

$$z \in \mathbb{R}^k.$$

If we assume that the component of  $z$  vector can be factored into  $k$ -independent scalar, probability distribution, then, this is called mean field assumption.



Mean field inference makes computation easier -



co-ordinate ascent - multiple variables • Start with random initialization, all except some fixed, optimize the ones that are not held fixed.

Worked well with classic EM.

Once updated estimates are got, hold them fixed and update the other ones & so on.

Can we do gradient ascent?  $\rightarrow$  Variational A.E.

$z \sim N(0, I_{k \times k}) \rightarrow$  Latent variable / prior  $\rightarrow$  sampled from some normal distribution.

Important.

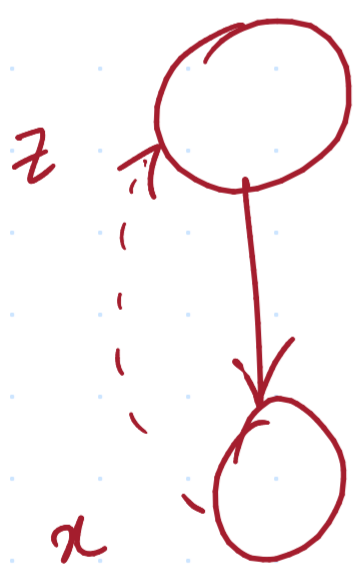
$$p(z|x) = \frac{p(z|x) \phi(x)}{p(z)}$$

posterior  $\downarrow$  evidence  $\swarrow$

likelihood  $\swarrow$  prior  $\uparrow$

$x|z \sim N(g(z; \theta), \sigma^2 I)$   
 $\downarrow$   
Likelihood.

$p(z|x)$ ;  $g \rightarrow$  Neural Network with param,  $\theta$ , we are never going to obtain.



Approximate using variational inference  $p(z|x)$

$\mathcal{Q} \rightarrow$  family for Variational Inference (V.I.)

$$\mathcal{Q}_i(z) = N(\underbrace{g(z^{(i)}; \phi)}_{\in \mathbb{R}^k}, \text{diag}(\underbrace{V(x^{(i)}; \psi)}_{\in \mathbb{R}^{k \times k}})^2)$$

We will get different  $Q \rightarrow$  dist per example. In E.M. E-step is performed separately per example.

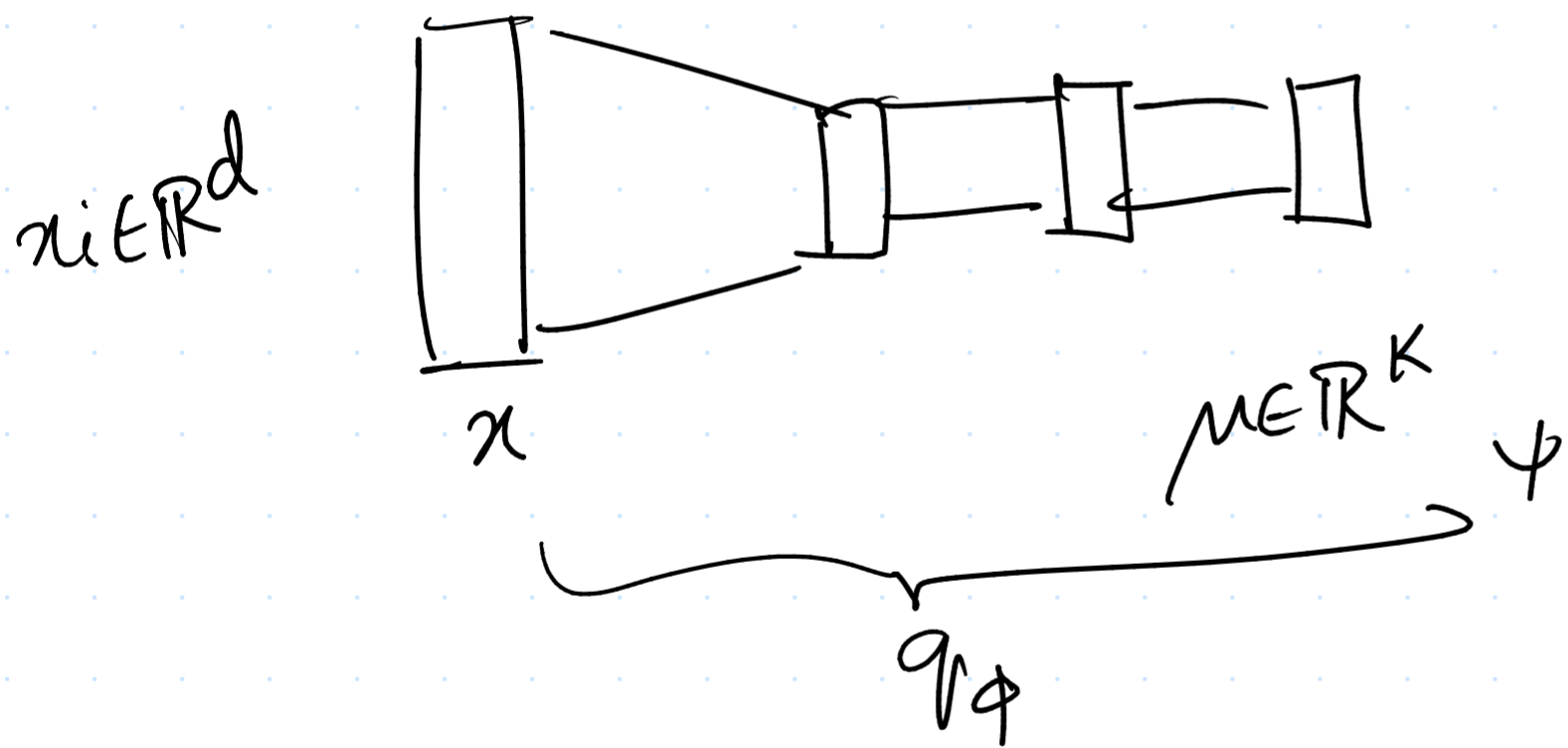
$Q_i \rightarrow$  depends on  $x$   $Q_i = p(z|x)$  N.N.  $\rightarrow$  takes  $x$  as input

"Amortized Influence"

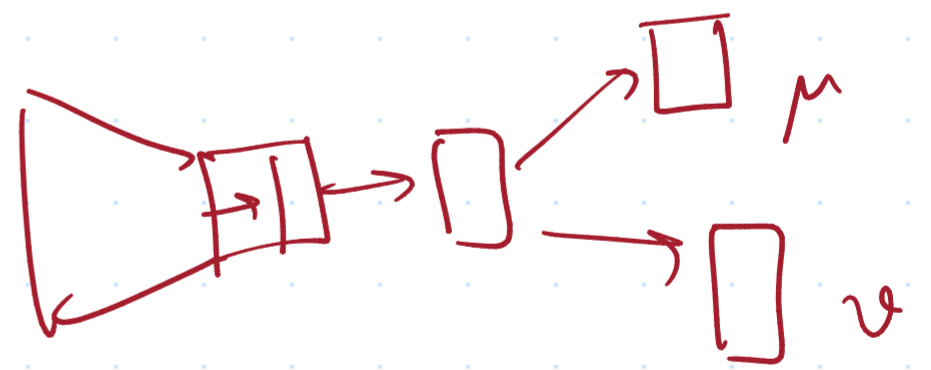
function of  $x$ , and those function is the N.N., so we don't have to optimize for each example independently as the E-M step.

$$q_\phi(x^{(i)}) = \mu^{(i)}$$

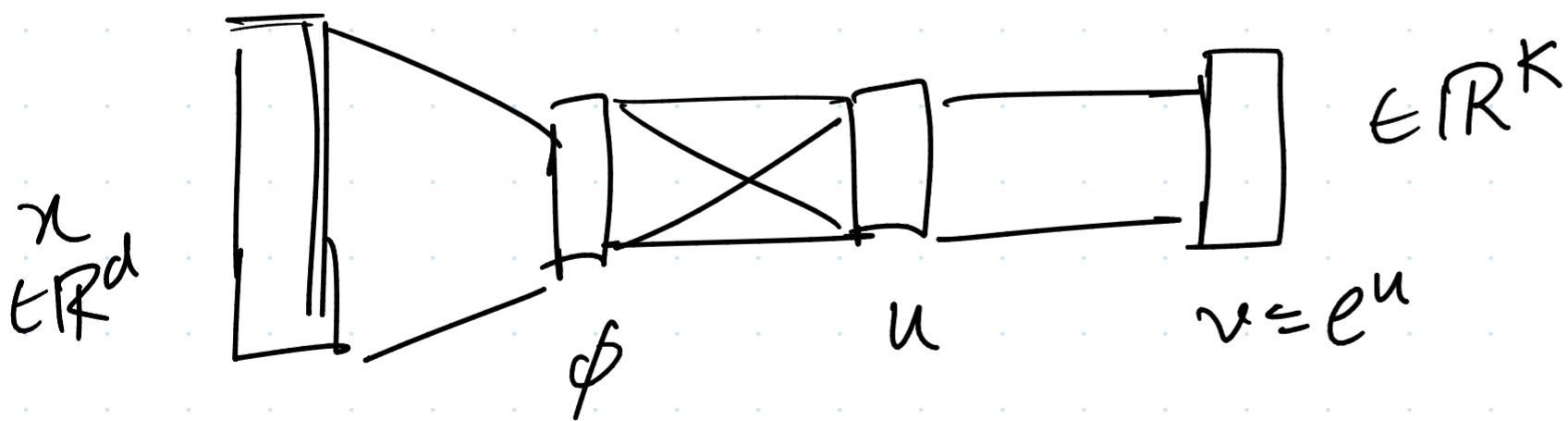
why is  $k < d$ ,  $\mu \rightarrow$  smaller dim.



$Q \rightarrow$  dist over  $z$ .  
and  $z = k$  dim,  $\mu = k$ -dim.



Covariance



to get +ve s.d.

positive = square them, exponentiate each element.

$Q_i \rightarrow$  mean & covariance.

Normal Dist  $\rightarrow$  vector & a matrix  
in high dimensional space.

$$k \begin{bmatrix} v_1 & & & \\ & v_2 & & \\ & & \ddots & \\ & & & v_k \end{bmatrix}$$

Gaussian dist :- & diagonal covariance matrix. If there is  
no correlation b/w two components of a joint Gaussian,  
then, they necessarily must be independent.

Diagonal covariance matrix  $\rightarrow$  making the mean field  
assumption. Each of  $z_i$ 's  $\rightarrow$  independent.

$\phi$  &  $\psi \rightarrow$  shared across all the examples.

Decoder  
 $p(x|z) \cdot p(z)$

$$ELBO(\phi, \psi, \theta) = \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q^{(i)}} \left[ \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

Where,  $Q_i = N(\mu(x^{(i)}; \phi), \text{diag}(v(x^{(i)}; \psi)))$ .

$$ELBO(Q, \theta) = \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim Q^{(i)}} \left[ \log \frac{p(x, z; \theta)}{Q_i(z^{(i)})} \right]$$

ELBO in EM, where each  $\theta \rightarrow$  separately calculated for each example in parallel.

Optimize  $\phi, \psi, \theta$  with gradient descent:-

---

$$\theta := \theta + \eta \nabla_{\theta} \text{ELBO}(\phi, \psi, \theta)$$

$$\phi := \phi + \eta \nabla_{\phi} \text{ELBO}(\phi, \psi, \theta)$$

$$\psi := \psi + \eta \nabla_{\psi} \text{ELBO}(\phi, \psi, \theta)$$



$$\nabla_{\theta} \text{ELBO}(\phi, \psi, \theta)$$

$$= \nabla_{\theta} \sum_{i=1}^N \mathbb{E}_{z^{(i)} \sim q_i} \left[ \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{q_i(z^{(i)})} \right]$$

$$= \sum_{i=1}^N \mathbb{E}_{z^{(i)} \sim q_i} \left[ \nabla_{\theta} \log p(x^{(i)}, z^{(i)}; \theta) \right]$$

$$= \sum_{i=1}^N \mathbb{E}_{z^{(i)} \sim q_i} \left[ \underbrace{\nabla_{\theta} \log p(x^{(i)} | z^{(i)}; \theta)}_{\text{Generator Decoder}} + \nabla_{\theta} \log p(z^{(i)}) \right]$$

Monte Carlo (M.C.)



$$\nabla_{\phi} \text{ELBO}(\phi, \psi, \theta) = \nabla_{\phi} \sum_{i=1}^N \mathbb{E}_{z^{(i)} \sim q_i} \left[ \log \frac{p(x, z; \theta)}{q(z^{(i)})} \right]$$

The distribution  $z^{(i)} \sim q_i$  depends on  $\phi$ , hence we cannot swap the expectations.

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

$$\boxed{z = \epsilon \cdot \sigma + \mu}$$

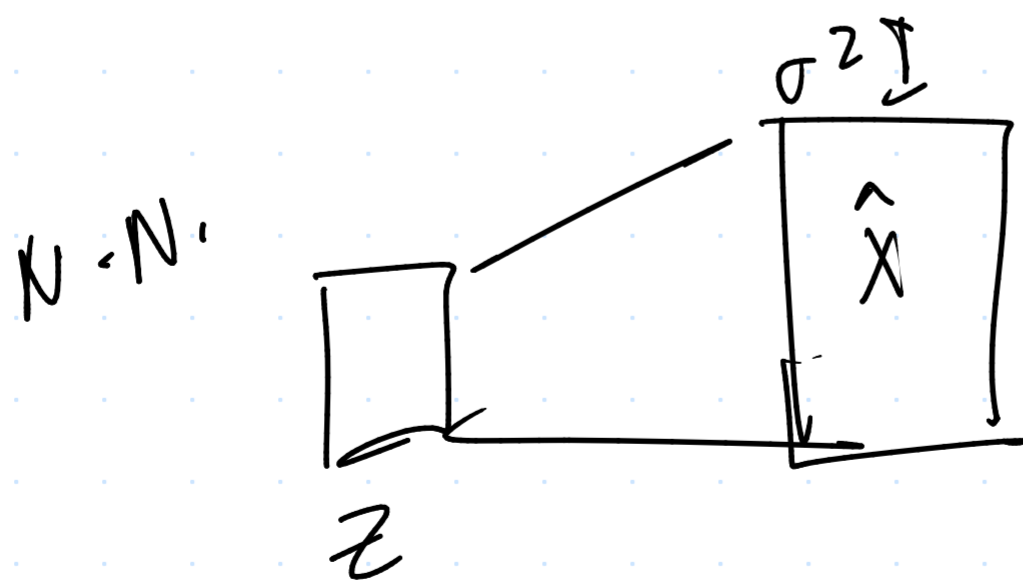
decoupled into this.  
 $\epsilon \sim N(0, 1)$

$$\nabla_{\phi} \sum_{i=1}^N \mathbb{E}_{\epsilon^{(i)} \sim N(0, 1)} \left[ \log \frac{p(x, \boxed{\epsilon^{(i)} \odot \Sigma^{(i)} + \mu^{(i)}}; \theta)}{q(\boxed{\epsilon^{(i)} \odot \Sigma^{(i)} + \mu^{(i)}})} \right]$$

$\uparrow z^{(i)}$   
 $\downarrow z^{(i)}$

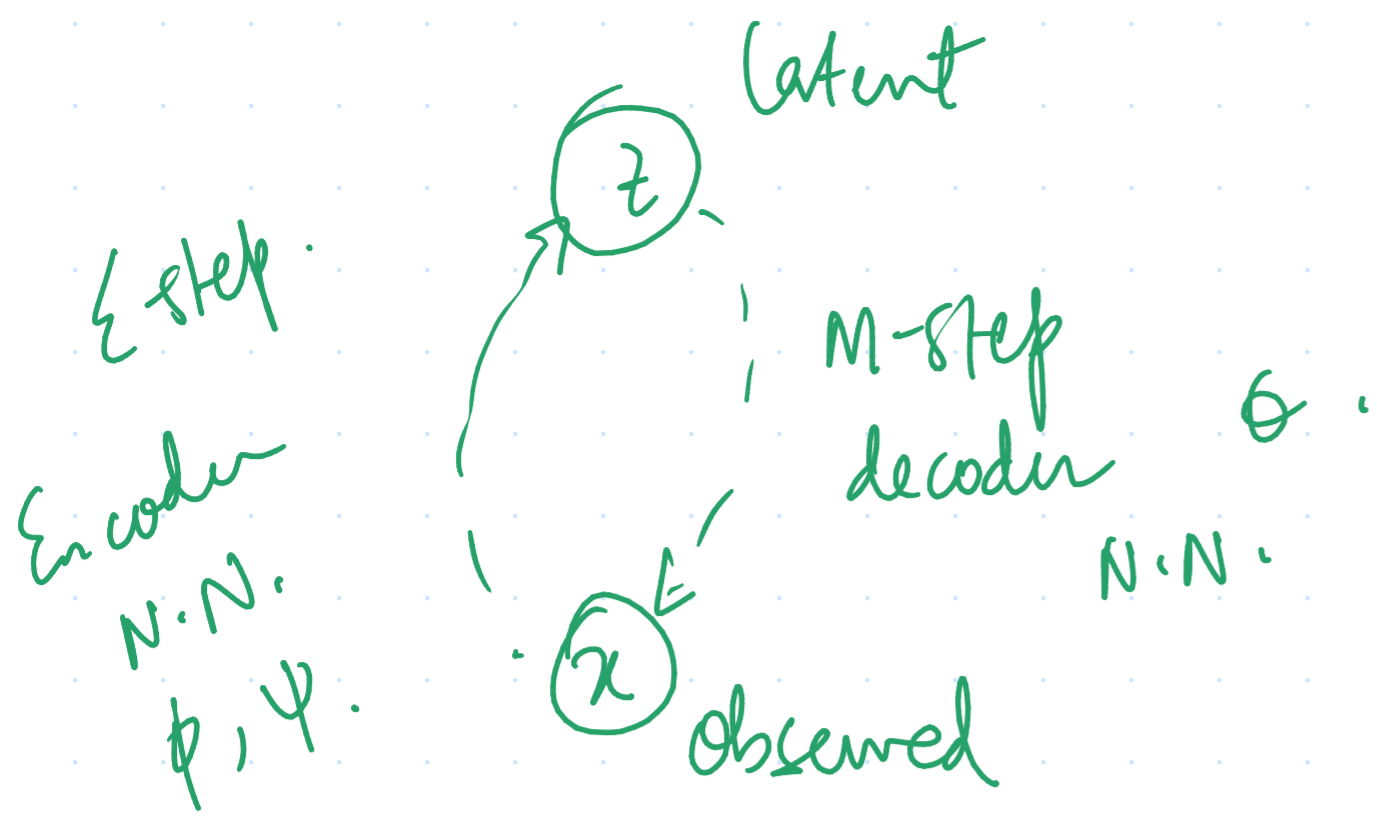
where,  $\mu^{(i)} = q(x^{(i)}; \phi)$

$\Sigma^{(i)} = \text{diag}(v(x^{(i)}); \psi)$ .



$p(x^{(i)} | z^{(i)}; \theta)$ .

$\log p(x^{(i)}; g(z^{(i)}), \sigma^2 I)$



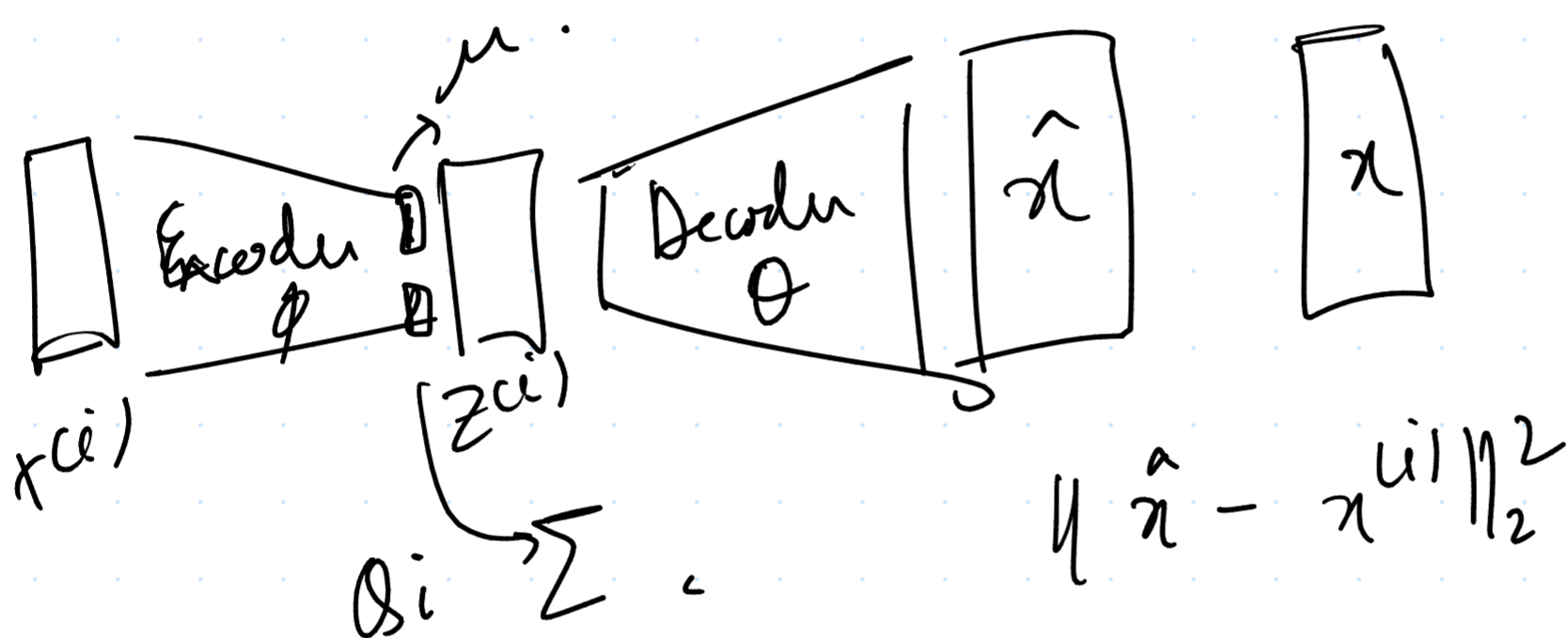
Maximizing the ELBO,

$$\text{ELBO}(\phi, \psi, \theta)$$

Easy gradient.

↓  
Reparameterization trick.

Expectation of the gradients are approximated using Monte-Carlo estimate.



Maximizing the ELBO - calculating the gradients and then gradient ascent steps. Alternate to EM  $\rightarrow$  fit a latent variable model.

