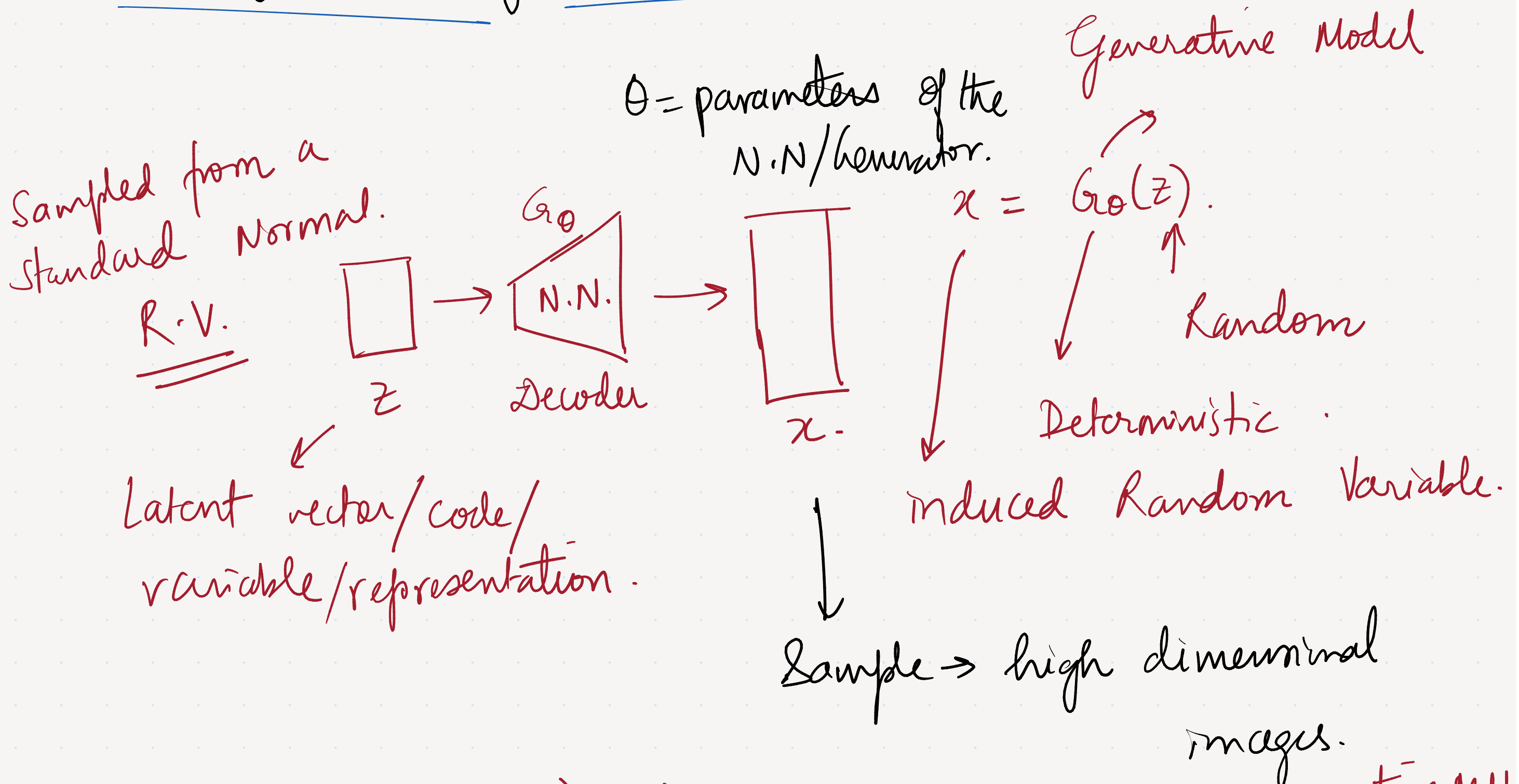


# Continuation of VAE-Theory

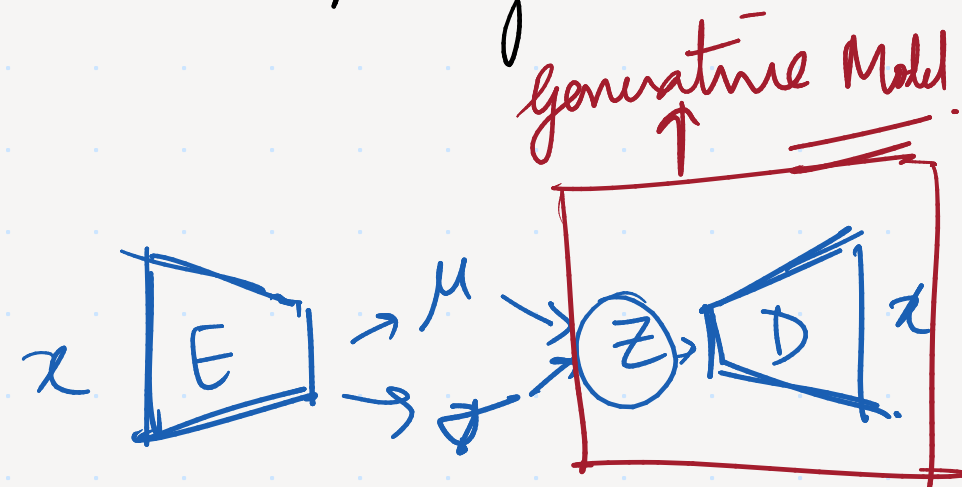
02/03/25

(A more rigorous treatment of the topic)

Sample an image from a distribution } need for  
 Training data for Neural Networks. } generation  
 of images.



Typically,  $z \sim \mathcal{N}(0, I) \rightarrow$  low dimensional.



Autoencoders :-

$x \rightarrow E_\phi \rightarrow h \rightarrow D_\theta \rightarrow x'$

$x' = D_\theta(E_\phi(x)) \approx x$

$\min_{\theta, \phi} \|x_i - D_\theta(E_\phi(x))\|^2$

L2-reconstruction error.

$\mu, \sigma$

$z \sim \mathcal{U}(0, I)$

V.A.E.

Scale with mean & add variance for reparam trick.

Autoencoder is not a generative model, since it is not sampling from a data distribution. Generative model  $\rightarrow$  Define a distribution.

Training a low latent-dimensional generative model by likelihood.

Given data  $\{x_i\}_{i=1,2,\dots,n}$  train a generative model to maximize the likelihood of the observed data under our model.

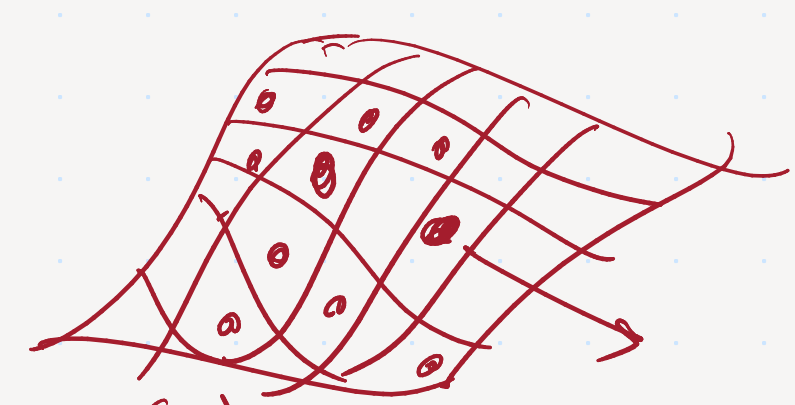
If Generative Model :-

$$G_\theta : \mathbb{R}^k \rightarrow \mathbb{R}^d, \text{ w/ } k < d.$$

(low dim)  $z \rightarrow x$  (high dim).

then  $p(x) = 0$  almost everywhere

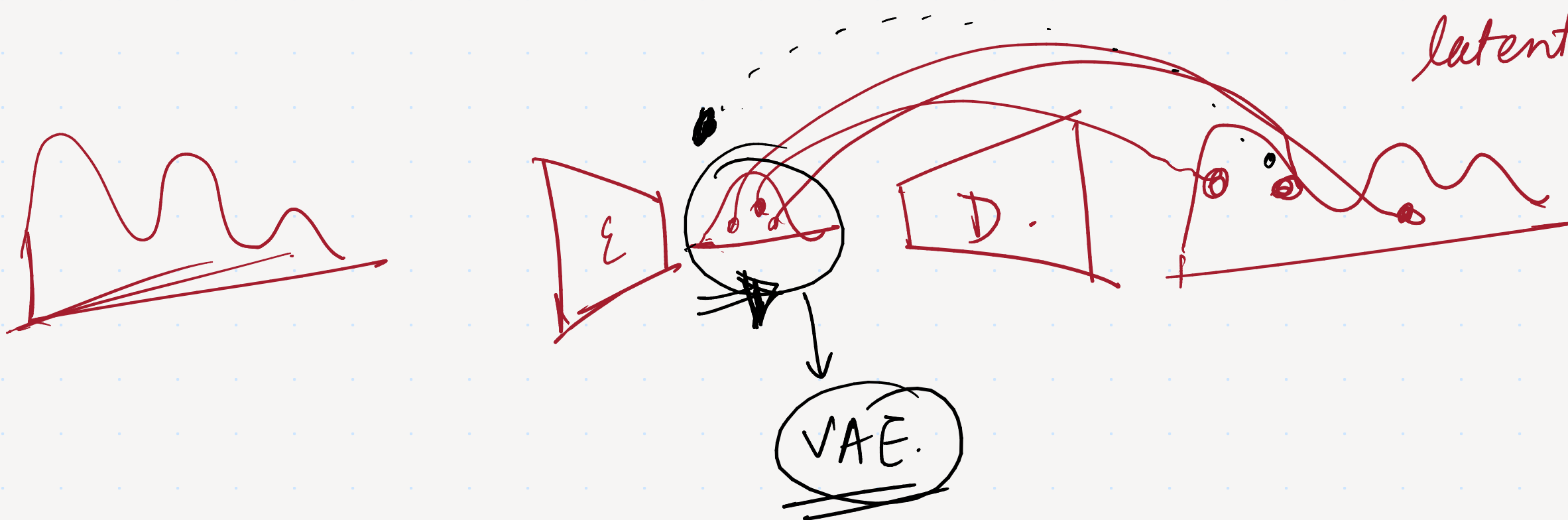
$\odot x$



Range( $G$ ). non-zero likelihood.

$\downarrow$

Manifold of the latent  $z$ .



$\therefore$  we have a non-zero likelihood only on the lower dimensional subset of space (subspace) (i.e.,  $\text{Range}(A)$ )

If we pick a randomly generated point  $x$ , off the manifold of  $A$ , then this point won't have any probability, i.e., on the higher dimension.

So, we can't directly optimize likelihood.

To have a non-zero likelihood everywhere, define noisy observation model. / Noisy inference model.

$$p_{\theta}(x|z) = \mathcal{N}(x; g_{\theta}(z), \eta I) \rightarrow \text{likelihood.}$$

$\mu$                        $\sigma$

$z \rightarrow$  induce a distribution in image space governed by  $\theta$ . Under a simple prior  $p(z)$ , this induces a joint distribution  $p_{\theta}(x, z)$ .

$z \in \mathcal{D} \quad x$   
 $p(x|z)$

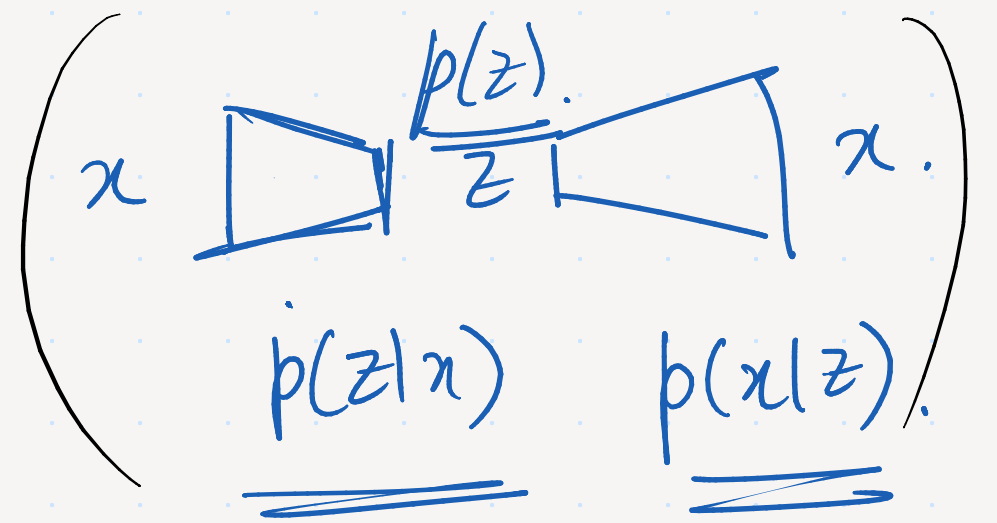
• intractable to evaluate at each iteration, hence we optimize a lower bound instead.



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

→ likelihood of  $z$

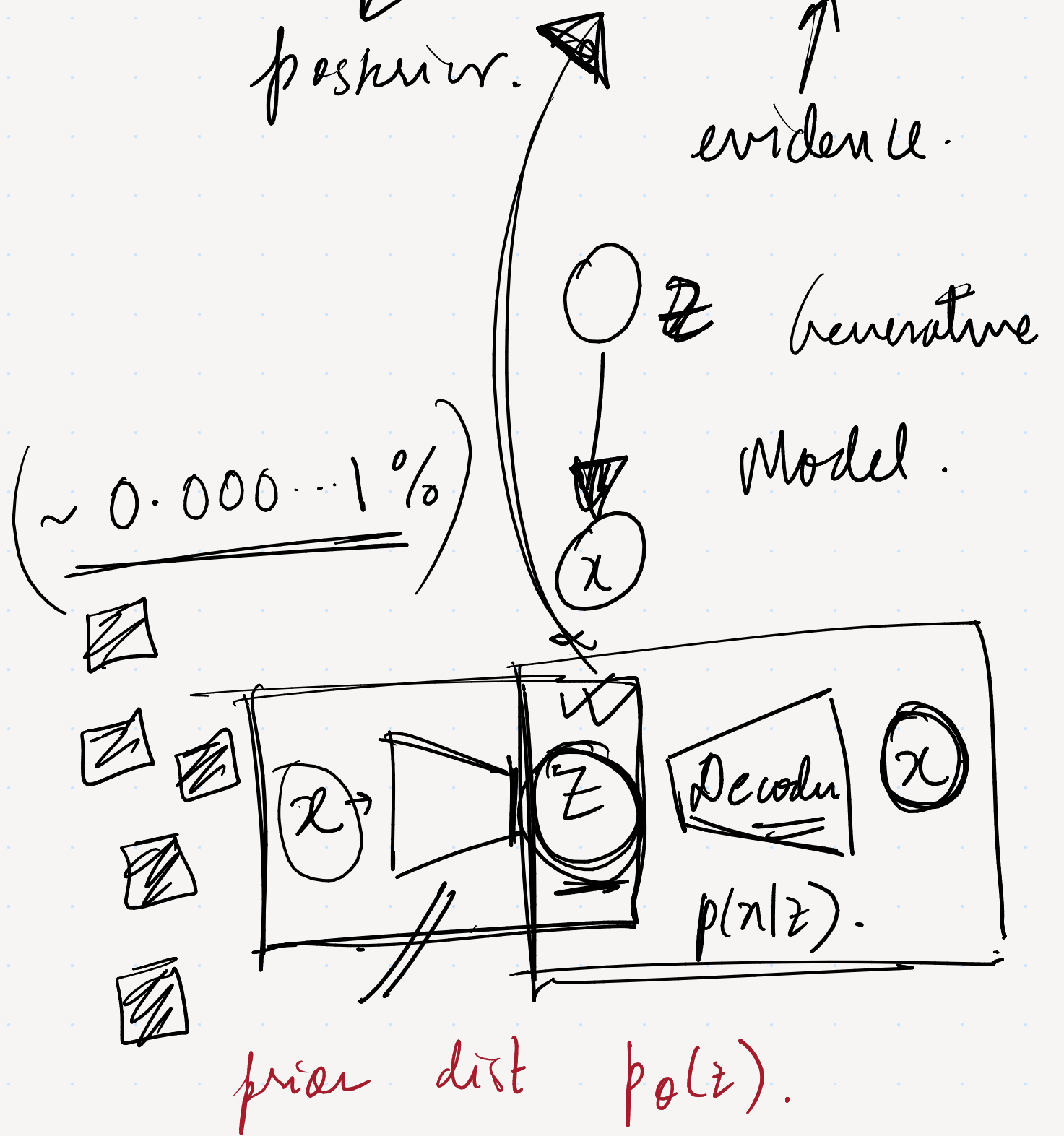
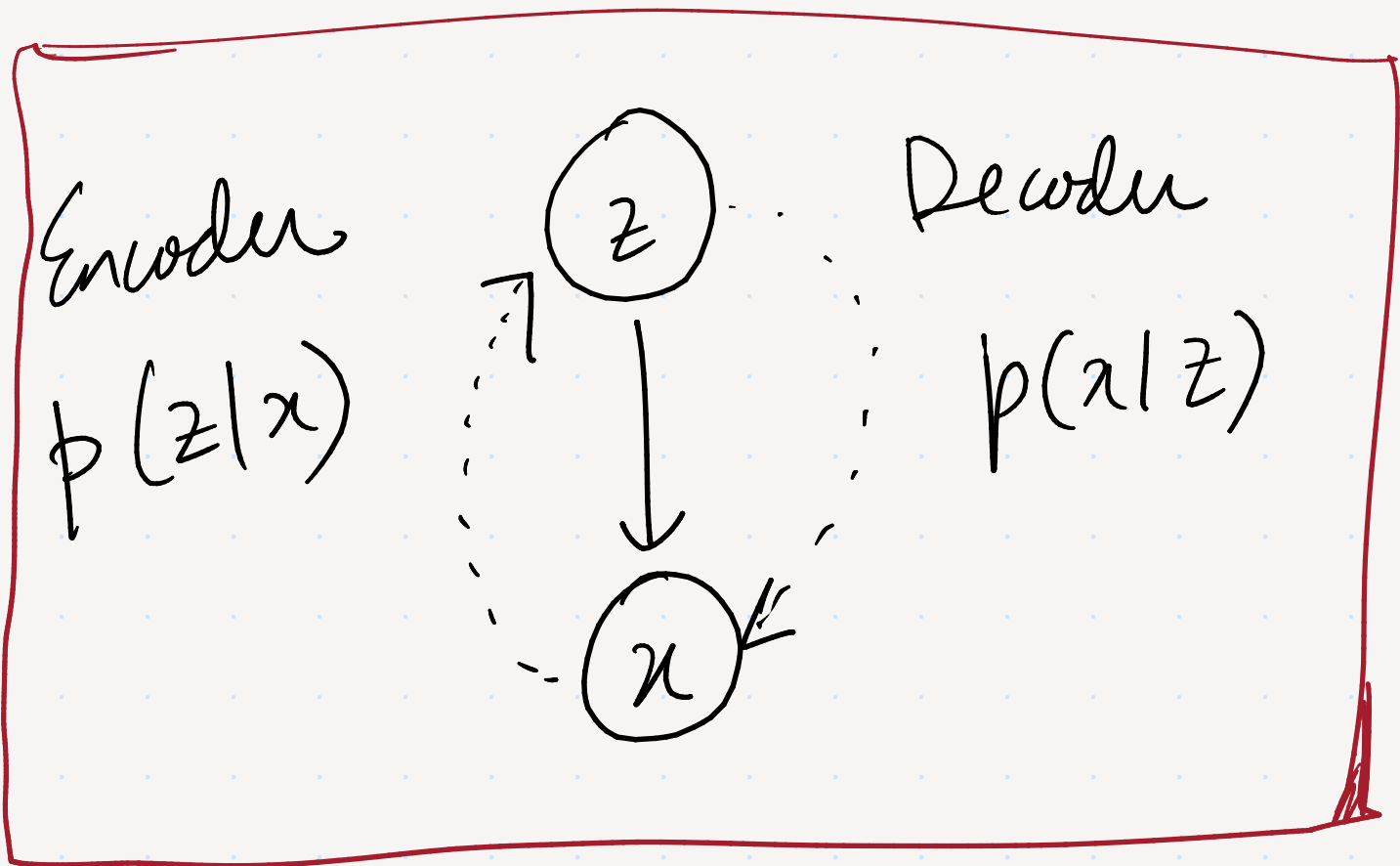
↳ likelihood of  $x|z$



Maximize the log-likelihood of the data.

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

↓ posterior.      ↑ evidence.



So, we optimize  $p(x|z)$   
 → the likelihood.

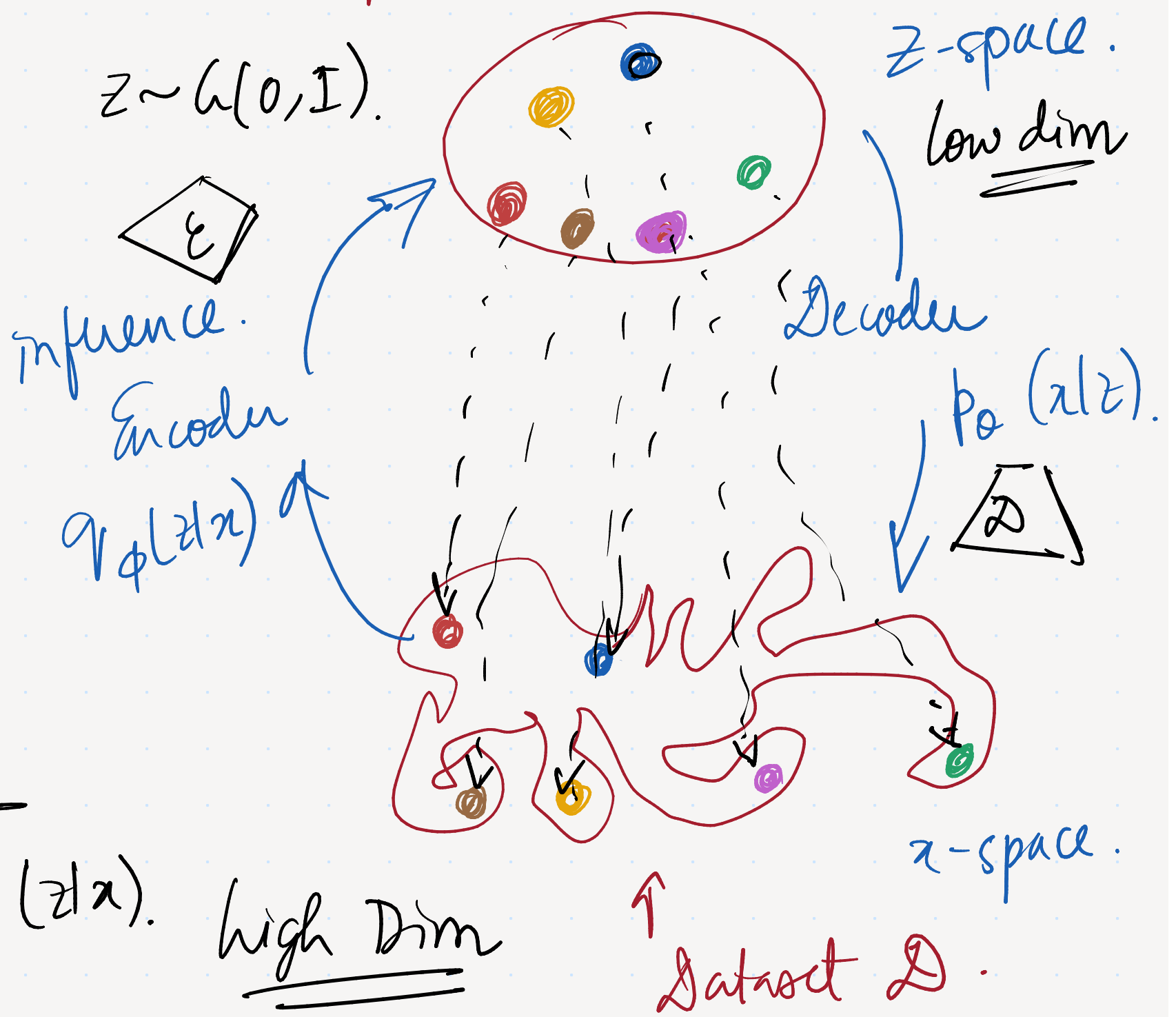
Setup:-

$$z \sim p(z) \rightarrow \text{prior.}$$

$$x \sim p_\theta(x|z) \rightarrow \text{Decoder.}$$

Use  $q_\phi(z|x)$  as a proxy for the encoder.

Intractable to compute  $p_\theta(z|x) \rightarrow$  intractable  $\sim q_\phi(z|x)$ .



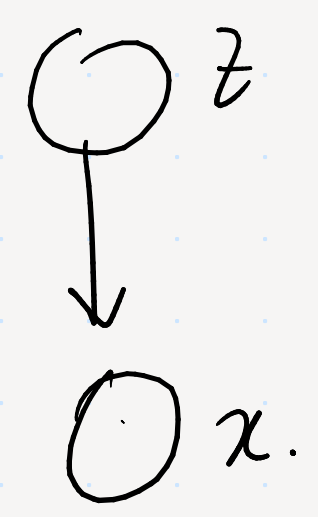


We will find the lower bound to  $p_\theta(x)$

$$p_\theta(z|x) = \frac{p_\theta(z,x)}{p_\theta(x)}$$

$$\log p_\theta(x) = \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x)$$

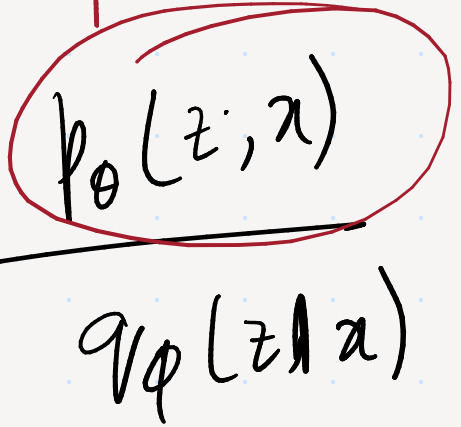
$$= \mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(z|x)}{p_\theta(z|x)}$$



$$= \mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(z|x)}{q_\phi(z|x)} - \frac{q_\phi(z|x)}{p_\theta(z|x)}$$

$$= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(z|x)}{q_\phi(z|x)}}_{L_{\theta, \phi}(x)} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{q_\phi(z|x)}{p_\theta(z|x)}}_{D_{KL}(q_\phi(z|x) \parallel p_\theta(z|x))}$$

*intractable*



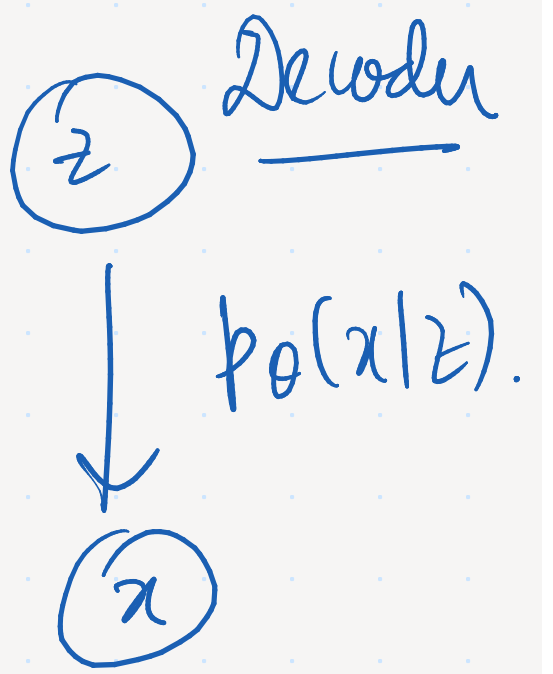
$D_{KL}(q_\phi(z|x) \parallel p_\theta(z|x))$

$\downarrow$   
*intractable.*

Variational Lower Bound  
Evidence Lower Bound.

We will optimize this instead

$$\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(z, x)}{q_\phi(z|x)} \rightarrow \text{intractable.}$$



$\mathcal{L}_{\theta, \phi}(x)$

Decoder  
 $p_\theta(x|z) \cdot p_\theta(z) \sim \text{Simple Normal.}$

$$p_\theta(x|z) = \frac{p_\theta(z, x)}{p_\theta(z)}$$

$$\Rightarrow \mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(x|z) \cdot p_\theta(z)}{q_\phi(z|x)} \text{ surrogate.}$$

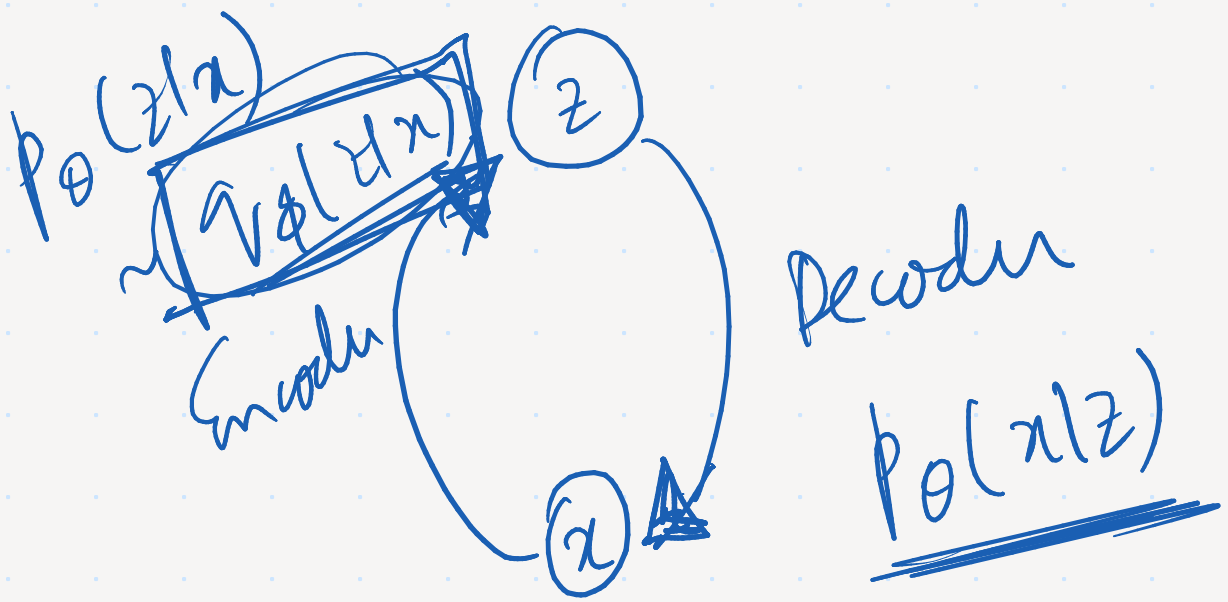
$$p_\theta(z, x) = p_\theta(x|z) \cdot p_\theta(z)$$

$$\Rightarrow \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z)}_{\text{Reconstruction error}} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(z)}{q_\phi(z|x)}}_{\text{DKL}(q_\phi(z|x) \parallel p_\theta(z))}$$

(Reconstruction error)

- DKL ( $q_\phi(z|x) \parallel p_\theta(z)$ )

Regularization



AE



$x \rightarrow \text{Dataset}$   
 $x \rightarrow \text{generate via } z$

$$\underline{\underline{\min \|x - x'\|_2}}$$





$$D_{KL}(q \parallel p) = \mathbb{E}_{z \sim q} \log \frac{q(z)}{p(z)}$$

•  $D_{KL}(q \parallel p) \neq D_{KL}(p \parallel q)$

• Measure of how far  $p$  is from  $q$ , hence if  $p \sim q$  then  $D_{KL}(p \parallel q) = 0$

•  $D_{KL}(q \parallel p) \geq 0$  & is  $= 0$  if  $p = q$ .

$$\mathcal{L}_{\phi, \theta}(x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} \log p_{\theta}(x|z) + \underbrace{\mathbb{E}_{z \sim q_{\phi}(z|x)} \log \frac{p_{\theta}(z)}{q_{\phi}(z|x)}}_{-D_{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z))}$$

max -  $\mathcal{L} \rightarrow$  max  $D_{KL}$

hence  $D_{KL} = 0$  when  $p = q$

So maximizing the VLB,  $\mathcal{L}_{\phi, \theta}$  pushes  $q_{\phi}(z|x)$  towards  $p(z)$ , prevents  $q_{\phi}$  from being a point mass.

$$\log p_{\theta}(x) = \underbrace{\mathbb{E}_{z \sim q_{\phi}(z|x)} \log p_{\theta}(x|z)}_{\mathcal{L}_{\theta, \phi}(x)} - D_{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z)) +$$

$\mathcal{L}_{\theta, \phi}(x)$

$$D_{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z)) \left\{ \mathbb{E}_{z \sim q_{\phi}(z|x)} \log \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \right.$$

So, Maximizing VLB,  $L_{\theta, \phi}$  :-

- Maximizes  $p(x)$   $\rightarrow$  roughly

- Minimizes KL divergence b/w  $q_{\phi}(z|x)$  &  $p_{\theta}(z|x)$  making  $q_{\phi}$  better.

Instead of optimizing  $\sum_{i=1}^n \log p_{\theta}(x_i)$ , we optimize.

We optimize the lower bound to be max.

$$\sum_{i=1}^n L_{\theta, \phi}(x_i)$$

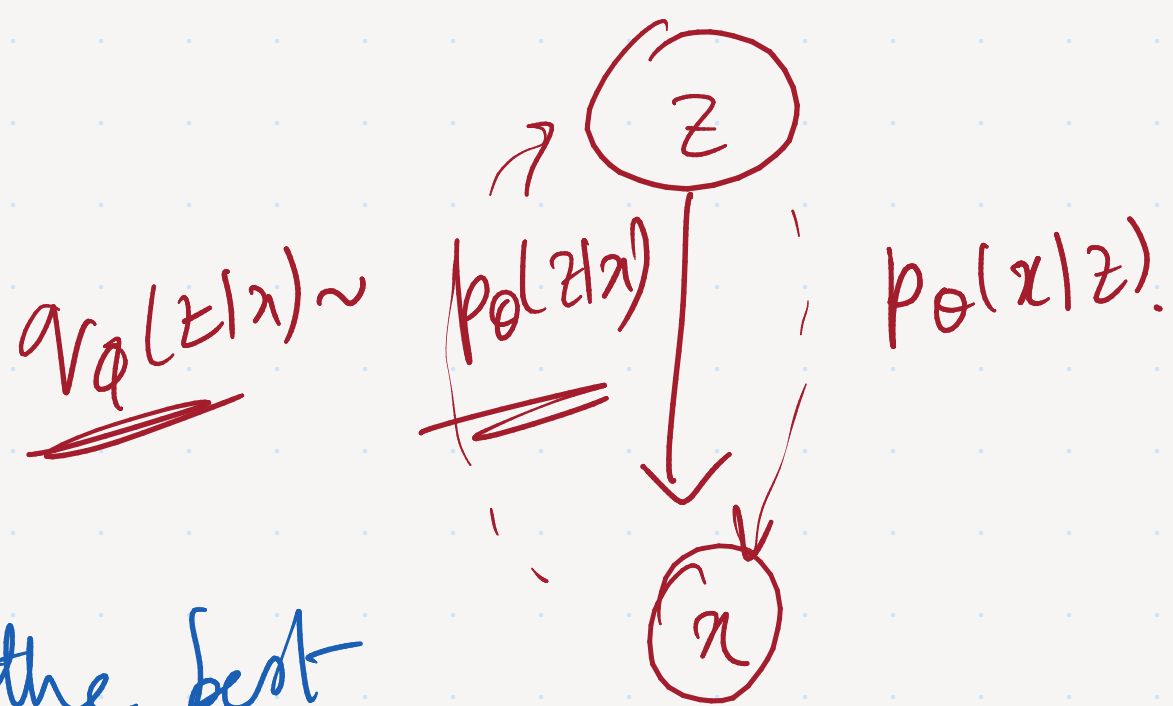
w/  $L_{\theta, \phi}(x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)}$

intractable.

Lower bound.

Optimizing VLB

$$\max_{\theta, \phi} \sum_{i=1}^n L_{\theta, \phi}(x_i)$$



One possibility  $\rightarrow$  for each  $x_i$  find the best  $q_{\phi}(z|x_i)$  by multiple gradient steps in  $\phi$ . Then gradient ascent in  $\theta$ .

This results in expensive inference steps / updates.

Instead

Amortize the inference costs by learning an inference  $n/w$ .

$$x \rightarrow (\mu, \Sigma)$$

$$w/ \quad q_{\phi}(z|x) = N(z; \mu(x), \Sigma(x))$$

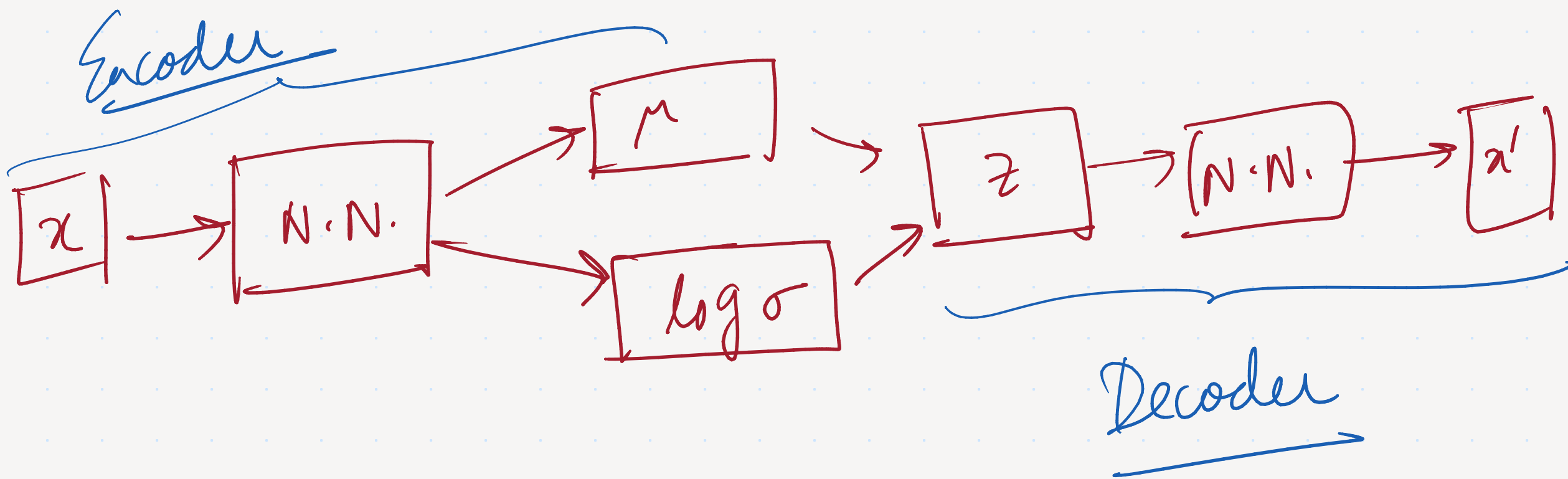
or  $\sigma(x)I$

Parameters of the inference models are shared b/w the data-points

Reparameterization Trick.

VAE - Architecture

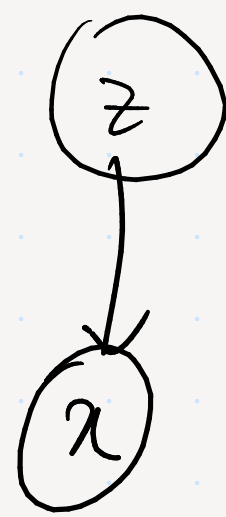
Separate random sources from differentiable quantities.





# Stochastic Gradient Optimization of the VLB.

Dataset  $\mathcal{D} = \{x_i\}_{i=1,2,3,\dots,n}$



Solve,  $\max_{\theta, \phi} \sum_{x_i \in \mathcal{D}} L_{\theta, \phi}(x_i)$

intractable.

$$\text{w/ } \underline{L_{\theta, \phi}(x)} = \mathbb{E}_{z \sim q_{\phi}(z|x)} \left( \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right)$$

Computing  $\nabla_{\theta, \phi} L_{\theta, \phi}(x_i)$  is intractable, but there are unbiased estimators. We need to sample all of  $z$  that are

possible & then differentiate b/w  $\theta$  &  $\phi$ , and hence

this is intractable.

Easy to get unbiased  $\nabla_{\theta} L_{\theta, \phi}$  :-

$$\nabla_{\theta} L_{\theta, \phi}(x) = \nabla_{\theta} \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \log p_{\theta}(x, z) - \log q_{\phi}(z|x) \right]$$

$$= \mathbb{E}_{z \sim q_{\phi}(z|x)} \nabla_{\theta} \left[ \log p_{\theta}(x, z) - \log q_{\phi}(z|x) \right]$$

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \nabla_\theta \log p_\theta(x, z).$$

Sample from  $\mathcal{D}$   
and evaluate at  
randomly chosen  $z$ .

$$\approx \nabla_\theta \log p_\theta(x, z) \rightarrow w / z \sim q_\phi(z|x)$$

↓  
unbiased estimate.

But it is not easy to get an unbiased estimate of  $\nabla_\phi L_{\theta, \phi}(x)$  since  $\nabla$  doesn't commute.

$$\nabla_\phi L_{\theta, \phi}(x) = \nabla_\phi \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

$$\neq \mathbb{E}_{z \sim q_\phi(z|x)} \nabla_\phi (\log p_\theta(x, z) - \log q_\phi(z|x))$$

Since,  $q_\phi(z|x)$  will keep on changing, when we differentiate w.r.t.  $\phi$

$$\text{Recall, } q_\phi(z|x) = \mathcal{N}(z; \mu(x), \sigma(x) \mathbf{I})$$

$$= \mu(x) + \sigma(x) \cdot \epsilon, \quad w / \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$

$$\mathcal{L}_{\theta, \phi}(x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} \left( \log p_{\theta}(x, z) - \log q_{\phi}(z|x) \right)$$

$$= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left( \log p_{\theta}(x, z) - \log q_{\phi}(z|x) \right)$$

form an estimator of  $\mathcal{L}_{\theta, \phi}(x)$  as  $\tilde{\mathcal{L}}_{\theta, \phi}(x)$  by:-

$$\epsilon \sim p(\epsilon)$$

$$z = \mu_{\phi}(x) + \sigma_{\phi}(x) \cdot \epsilon = g(\phi, x, \epsilon)$$

$$\hat{\mathcal{L}}_{\theta, \phi}(x) = \log p_{\theta}(x, z) - \log q_{\phi}(z|x)$$

Unbiased estimate of  $\nabla_{\phi} \mathcal{L}_{\theta, \phi}(x)$  :-

$$\nabla_{\phi} \hat{\mathcal{L}}_{\theta, \phi}(x)$$

Note :-  $\mathbb{E}_{\epsilon \sim p(\epsilon)} \hat{\mathcal{L}}_{\theta, \phi}(x) = \mathcal{L}_{\theta, \phi}(x)$

So,  $\mathbb{E}_{\epsilon \sim p(\epsilon)} \nabla_{\phi} \tilde{\mathcal{L}}_{\theta, \phi}(x) = \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} \hat{\mathcal{L}}_{\theta, \phi} = \nabla_{\phi} \mathcal{L}_{\theta, \phi}$

optimize VAE params. wr.t. stochastic gradient.