

# Variational Autoencoders (VAEs)

Jayin Khanna  
CSD662: ADL

# Latent Variable Models & The Learning Goal

## 1. What is a Latent Variable Model?

We assume our observed data  $\mathbf{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  is generated from some true but unknown distribution  $p_{\mathbf{x}}(\mathbf{x})$ .

A latent variable model explains this data by introducing an unobserved (latent) variable  $\mathbf{z}$ .

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x}, \mathbf{z}) \\ &\text{or equivalently} \\ p_{\theta}(\mathbf{x}) &= \int p(\mathbf{z}) p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \end{aligned}$$

Intuition:

- **Discrete  $\mathbf{z}$  (e.g.,  $\mathbf{z} \in \{1, \dots, M\}$ ):**  
The latent variable  $\mathbf{z}_i$  acts as a cluster assignment for each data point  $\mathbf{x}_i$ .  
→ Example: Gaussian Mixture Models (GMMs)
- **Continuous  $\mathbf{z}$  (e.g.,  $\mathbf{z} \in \mathbb{R}^k$ , with  $k \ll d$ ):**  
The latent variable  $\mathbf{z}_i$  is a low-dimensional feature vector or compressed representation of  $\mathbf{x}_i$ .  
→ Foundation of Autoencoders and Variational Autoencoders (VAEs)

## Example — Mixture of Gaussians (Discrete Latent Variable) 1D Mixture of Gaussians

In a **1D mixture of Gaussians**, the latent variable  $\mathbf{z}$  is **discrete**, and the prior  $\mathbf{Pr}(\mathbf{z})$  is a **categorical distribution** with one probability  $\lambda_n$  for each possible value of  $\mathbf{z}$ .

$$\mathbf{Pr}(\mathbf{z} = n) = \lambda_n$$

The **likelihood** of the data  $\mathbf{x}$  given  $\mathbf{z} = n$  is normally distributed with mean  $\mu_n$  and variance  $\sigma_n^2$ :

$$\mathbf{Pr}(\mathbf{x} \mid \mathbf{z} = n) = \mathcal{N}(\mu_n, \sigma_n^2)$$

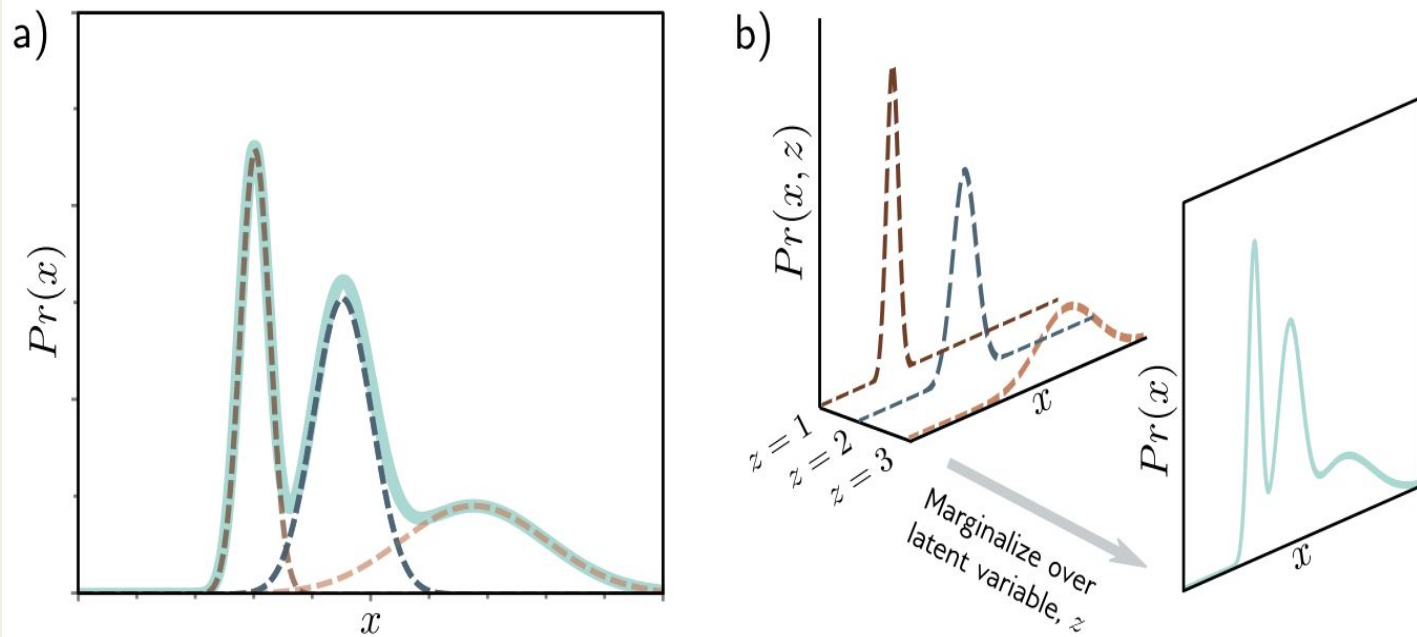
By marginalizing over all possible latent variable values, we get the overall data distribution:

$$\begin{aligned} Pr(x) &= \sum_{n=1}^N Pr(x, z = n) \\ &= \sum_{n=1}^N Pr(x|z = n) \cdot Pr(z = n) \\ &= \sum_{n=1}^N \lambda_n \cdot \text{Norm}_x[\mu_n, \sigma_n^2]. \end{aligned}$$

**Interpretation:**

From **simple components** (a prior over  $\mathbf{z}$  and Gaussian likelihoods), we can construct a **complex multi-modal distribution**.

Each mode corresponds to a **cluster**, and the mixture weights  $\lambda_n$  control their relative importance.



**Figure 17.1** Mixture of Gaussians (MoG). a) The MoG describes a complex probability distribution (cyan curve) as a weighted sum of Gaussian components (dashed curves). b) This sum is the marginalization of the joint density  $Pr(x, z)$  between the continuous observed data  $x$  and a discrete latent variable  $z$ .

# Nonlinear Latent Variable Models

## Continuous Multivariate Latent Variable Model

In the **nonlinear latent variable model**, both the data  $\mathbf{x}$  and the latent variable  $\mathbf{z}$  are **continuous and multivariate**.

**Prior:**

$$\Pr(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$$

**Likelihood:**

The likelihood of  $\mathbf{x}$  given  $\mathbf{z}$  and parameters  $\boldsymbol{\phi}$  is also Gaussian,

but its **mean** is a **nonlinear function** of  $\mathbf{z}$ , modeled by a neural network  $\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}]$ , and its covariance is  $\sigma^2 \mathbf{I}$  (spherical):

$$\Pr(\mathbf{x} | \mathbf{z}, \boldsymbol{\phi}) = \mathcal{N}(\mathbf{x} | \mathbf{f}[\mathbf{z}, \boldsymbol{\phi}], \sigma^2 \mathbf{I})$$

**Marginalizing over  $\mathbf{z}$ :**

$$\begin{aligned} \Pr(\mathbf{x} | \boldsymbol{\phi}) &= \int \Pr(\mathbf{x}, \mathbf{z} | \boldsymbol{\phi}) d\mathbf{z} \\ &= \int \Pr(\mathbf{x} | \mathbf{z}, \boldsymbol{\phi}) \cdot \Pr(\mathbf{z}) d\mathbf{z} \\ &= \int \text{Norm}_{\mathbf{x}} [\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}], \sigma^2 \mathbf{I}] \cdot \text{Norm}_{\mathbf{z}} [\mathbf{0}, \mathbf{I}] d\mathbf{z}. \end{aligned}$$

**Interpretation:**

This can be viewed as an **infinite mixture of Gaussians**,

where each  $\mathbf{z}$  contributes a Gaussian centered at  $\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}]$ ,

weighted by  $\Pr(\mathbf{z})$ . The network  $\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}]$  captures the key structure of the data, and  $\sigma^2 \mathbf{I}$  represents residual noise.

# General Principle for Learning Latent Variable Models

Suppose we are given data

$$\mathcal{D} = \{x_i\}_{i=1}^n \quad \text{i.i.d.} \quad \sim p_X.$$

Let  $p_\theta(x) = \int_z p_\theta(x, z) dz$  be a latent variable model. The goal is to estimate the model parameters  $\theta$  given data  $\mathcal{D}$ .

$$\theta^* = \arg \min_{\theta} \text{KL}(p_X \| p_\theta).$$

Expanding the KL divergence:

$$\theta^* = \arg \min_{\theta} \left[ \int_x p_X(x) \log \frac{p_X(x)}{p_\theta(x)} dx \right]. \quad \begin{array}{l} \text{Dropping the constant entropy term } H(p_X) = - \int_x p_X(x) \log p_X(x) dx \text{ (independent of } \\ \theta), \text{ this becomes:} \end{array} \quad \theta^* = \arg \min_{\theta} \left[ - \int_x p_X(x) \log p_\theta(x) dx \right].$$

- The KL minimization objective compares the true data distribution  $p_X$  with the model distribution  $p_\theta$ . Minimizing it makes the model approximate the data as closely as possible.

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{p_X} [\log p_\theta(x)].$$

Denoting the log-likelihood as  $\ell(\theta) = \log p_\theta(x)$ , this is precisely the **Maximum Likelihood Estimation (MLE)** principle:

$$\theta^* = \arg \max_{\theta} \ell(\theta) = \arg \max_{\theta} \mathbb{E}_{p_X} [\log p_\theta(x)].$$

In a latent variable model:

$$p_\theta(x) = \int_z p_\theta(x, z) dz,$$

so the log-likelihood involves integrating out the latent variable  $z$ .

# Evidence Lower Bound (ELBO)

Consider  $\ell(\theta) = \log p_{\theta}(x)$

$$= \log \int_{\mathbf{z}} p_{\theta}(x, \mathbf{z}) d\mathbf{z}$$

Sps  $q(\mathbf{z}|x)$  denote a distrib<sup>n</sup> over latent variable,  $\mathbf{z}$

$$\ell(\theta) = \log \int_{\mathbf{z}} p_{\theta}(x, \mathbf{z}) \frac{q(\mathbf{z}|x)}{q(\mathbf{z}|x)} d\mathbf{z}$$

$$= \log \int_{\mathbf{z}} q(\mathbf{z}|x) \left[ \frac{p_{\theta}(x, \mathbf{z})}{q(\mathbf{z}|x)} \right] d\mathbf{z}$$

Can't calculate!  $\left. \begin{array}{l} \text{here, this} \end{array} \right\} = \log \left[ \mathbb{E}_{q(\mathbf{z}|x)} \left( \frac{p_{\theta}(x, \mathbf{z})}{q(\mathbf{z}|x)} \right) \right]$  we can not estimate the log of an expectation using the statistical estimates, hence we need to pull out the expectation outside the log.

## Jensen's Inequality

Jensen's Inequality:

$$\log \mathbb{E}(\cdot) \geq \mathbb{E} \log(\cdot)$$

$$\ell(\theta) = \log \mathbb{E}_{q(\mathbf{z}|x)} \left( \frac{p_{\theta}(x, \mathbf{z})}{q(\mathbf{z}|x)} \right) \geq \mathbb{E}_{q(\mathbf{z}|x)} \log \left( \frac{p_{\theta}(x, \mathbf{z})}{q(\mathbf{z}|x)} \right)$$

$\ell(\theta) \geq \mathcal{J}_{\theta}(q)$ , called the Evidence lower Bound [ELBO] denote with  $\mathcal{J}_{\theta}(q)$

$\mathcal{J}_{\theta}(q)$  is a fn of both the model parameters  $\theta$  & the dens<sup>y</sup> on  $\mathbf{z}$ ,  $q(\mathbf{z}|x)$

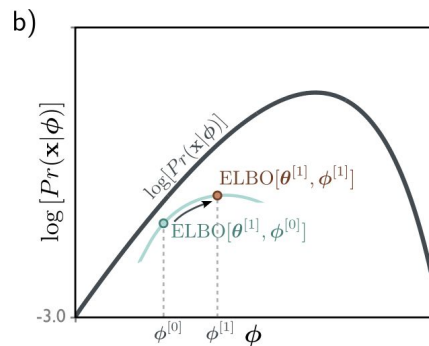
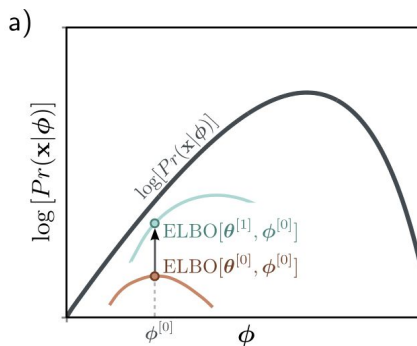
$\theta^* = \arg \max_{\theta} \ell(\theta)$ : original max likelihood problem

$\theta^*, q^* = \arg \max_{\theta, q} \mathcal{J}_{\theta}(q)$ : Approx with ELBO

The fundamental problem solved in any latent var generative model is given below

ELBO maximization, find  $\theta$  &  $q$  that maximizes the ELBO,  $\mathcal{J}_{\theta}(q)$

$$\theta^*, q^* = \arg \max_{\theta, q} \mathbb{E}_{q(\mathbf{z}|x)} \frac{p_{\theta}(x, \mathbf{z})}{q(\mathbf{z}|x)}$$



# ENCODER AND DECODER IN VAEs

\* Goal of a Neural lat. var. model:

- Learn a lat. var. model with unknown  $p_\theta(z|x)$
- Enable sampling/generation from  $p_\theta(z|x)$
- Enable posterior inference, compute/estimate  $q(z|x)$

$$\begin{aligned} \text{Consider } J_\theta(q) &= \mathbb{E}_{q(z|x)} \log \frac{p_\theta(x, z)}{q(z|x)} \\ &= \mathbb{E}_{q(z|x)} \log p_\theta(x|z) - \mathbb{E}_{q(z|x)} \log \frac{q(z|x)}{p_\theta(z)} \end{aligned}$$

$$J_\theta(q) = \mathbb{E}_{q(z|x)} \log p_\theta(x|z) - D_{KL}(q(z|x) || p_\theta(z))$$

\* Represent  $p_\theta(x|z)$  &  $q(z|x)$  via Neural network

$p_\theta(x|z)$ : conditional data likelihood

$q(z|x)$ : variational latent posterior density

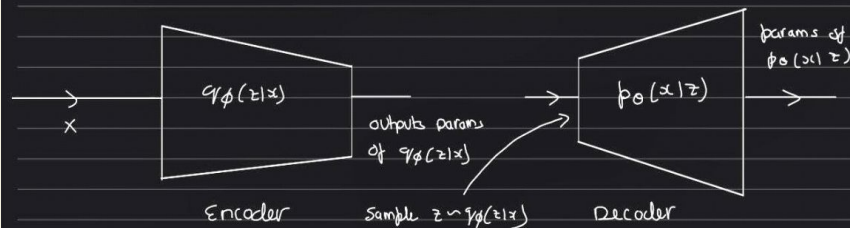
Given  $D = \{x_i\}_{i=1}^n \sim \text{iid } p_x$

$$p_\theta(z) = \int_z p_\theta(x, z) dz : \text{model}$$

The follow<sup>n</sup> objective is optimized

$$J_\theta(q) = \mathbb{E}_{q(z|x)} \log p_\theta(x|z) - D_{KL}(q(z|x) || p_\theta(z))$$

$q(z|x)$  &  $p_\theta(x|z)$  are rep using prob NNs



There is no direct connection b/w

the encoder & Dec

$\phi$  &  $\theta$ : the weights of Enc & Dec resp.

$$\phi^*, \theta^* = \underset{\theta, \phi}{\text{argmax}} J_\theta(q_\phi)$$

$$\phi^{t+1} \leftarrow \phi^t + \alpha \nabla_\phi J_\theta(q_\phi)$$

$$\theta^{t+1} \leftarrow \theta^t + \alpha \nabla_\theta J_\theta(q_\phi)$$

# The Reparameterization trick

\* To compute the gradients of  $J_\phi(q_\phi)$  w.r.t  $\phi$ :

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

$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(z|x)} \log p_\theta(x|z) &= \nabla_\phi \int_z q_\phi(z|x) \cdot \log p_\theta(x|z) \\ &= \int \nabla_\phi q_\phi(z|x) \cdot \log p_\theta(x|z) \end{aligned} \quad \left[ \begin{array}{l} \text{since gradient is a} \\ \text{linear operator} \end{array} \right]$$

Explain!!

$$= \int_z q_\phi(z|x) \left( \nabla_\phi \log p_\theta(x|z) \right) + \underbrace{\int_z \left( \nabla_\phi q_\phi(z|x) \right) \log p_\theta(x|z)}_{\substack{\text{is not an expectation} \\ \therefore \text{can't be computed}}}$$

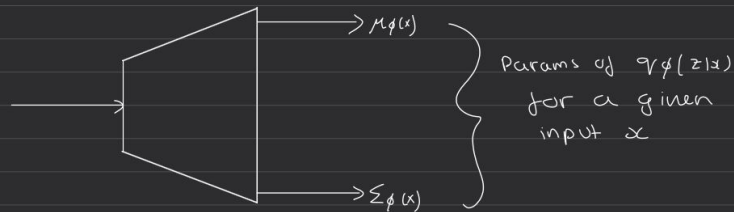
There is one more complication; the network involves a sampling step, and it is difficult to differentiate through this stochastic component. However, differentiating past this step is necessary to update the parameters  $\theta$  that precede it in the network.

Fortunately, there is a simple solution; we can move the stochastic part into a branch of the network that draws a sample  $\epsilon^*$  from  $\text{Norm}_\epsilon[0, \mathbf{I}]$  and then use the relation:

$$\mathbf{z}^* = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \boldsymbol{\epsilon}^*, \quad (17.25)$$



$$\text{let } q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x), \Sigma_{\phi}(x))$$



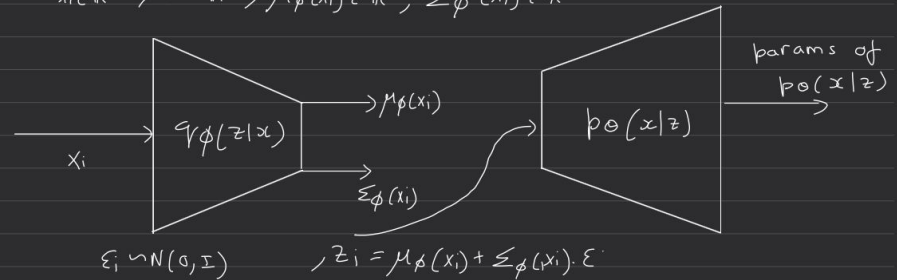
let  $\epsilon \sim \mathcal{N}(0, I)$ , then

$$z = \mu_{\phi}(x) + \Sigma_{\phi}(x)\epsilon = g(\epsilon)$$

$$z \sim \mathcal{N}(z; \mu_{\phi}(x), \Sigma_{\phi}(x))$$

then we have

$$x_i \in \mathbb{R}^d; z_i \in \mathbb{R}^k; \mu_{\phi}(x_i) \in \mathbb{R}^k; \Sigma_{\phi}(x_i) \in \mathbb{R}^{k \times k}$$



hence,

$$\begin{aligned} \nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)} \log p_{\theta}(x|z) &= \nabla_{\phi} \mathbb{E}_{p_{\epsilon}} \log p_{\theta}(x|g(\epsilon)) \\ &= \nabla_{\phi} \left( \frac{1}{M} \sum_{j=1}^M \log p_{\theta}(x|g(\epsilon_j)) \right) \end{aligned}$$

$$\epsilon_1, \epsilon_2, \dots, \epsilon_M \sim p_{\epsilon}$$

To compute  $\log p_\theta(x | g(\varepsilon_j))$ :

parameterize  $p_\theta(x|z)$  via some known distribution & use the decoder to output its params

Eg.  $p_\theta(x|z) = \mathcal{N}(x; \hat{x}_\theta(z), \mathbb{I})$

$$\begin{aligned} \Rightarrow \log p_\theta(x|z) &= \log \left[ \frac{1}{(2\pi)^{d/2}} \exp[-\|x - \hat{x}_\theta(z)\|_2^2] \right] \\ &= -\|x - \hat{x}_\theta(z)\|_2^2 \end{aligned}$$

hence,

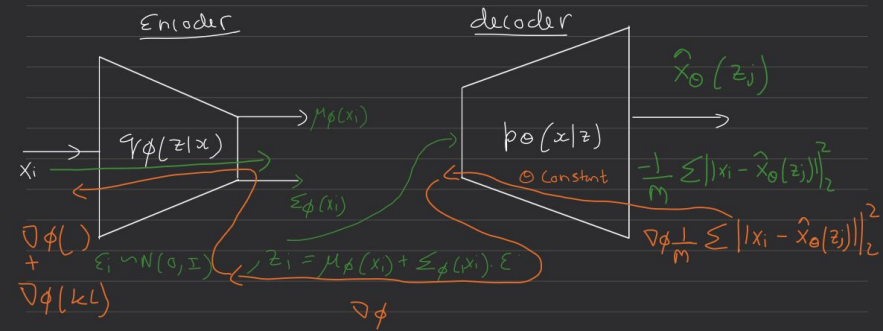
$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(z|x)} \log p_\theta(x|z) &= \nabla_\phi \mathbb{E}_{p_\theta} \log p_\theta(x | g(\varepsilon)) \\ &= \nabla_\phi \left[ -\frac{1}{n} \sum_{j=1}^n \|x_i - \hat{x}_\theta(z_j)\|_2^2 \right] \end{aligned}$$

where  $z_j = \mu_\phi(x_i) + \varepsilon_j \leq \phi(x_i)$ ,

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_j \sim \mathcal{N}(0, \mathbb{I})$

# Training the Network:

## 1. The Encoder



Forward pass: i) Given an  $x_i \in D$ , pass it through the Encoder, to obtain  $\mu_\phi(x_i)$  &  $\epsilon_\phi(x_i)$

ii) Sample  $\epsilon_1, \dots, \epsilon_M \sim \mathcal{N}(0, I)$  outside NNs

iii) Compute  $z_1, \dots, z_M$  via reparametrization

iv) Pass all of  $z_1, \dots, z_M$  through the Decoder to compute  $\hat{x}_0(z_j)$

v) Compute  $-\frac{1}{M} \sum_{j=1}^M \|x_i - \hat{x}_0(z_j)\|_2^2$

training for Encoder (one training step)

$$\phi^{t+1} \leftarrow \phi^t + \alpha \nabla_\phi J_\phi(q_\phi)$$

Backward pass:

i) Compute the grad  $\nabla_\phi \left( -\frac{1}{M} \sum_{j=1}^M \|x_i - \hat{x}_0(z_j)\|_2^2 \right)$

ii) Backprop Through

→ Decoder, with fixed param  $\theta$

→ Reparam, Id of  $\phi$

→ Encoder

Thus we compute  $\nabla_\phi \mathbb{E}_{q_\phi(z|x)} \log p_\theta(x|z)$

To compute  $\nabla_\phi D_{KL}[q_\phi(z|x) || p_\theta(z)]$ :

$p_\theta(z)$ : latent prior, typically assumed to be  $\mathcal{N}(0, I)$

$$\therefore D_{KL}[q_\phi(z|x) || p_\theta(z)]$$

$$= D_{KL}\left[\mathcal{N}(\mu_\phi(x), \epsilon_\phi(x)) || \mathcal{N}(0, I)\right]$$

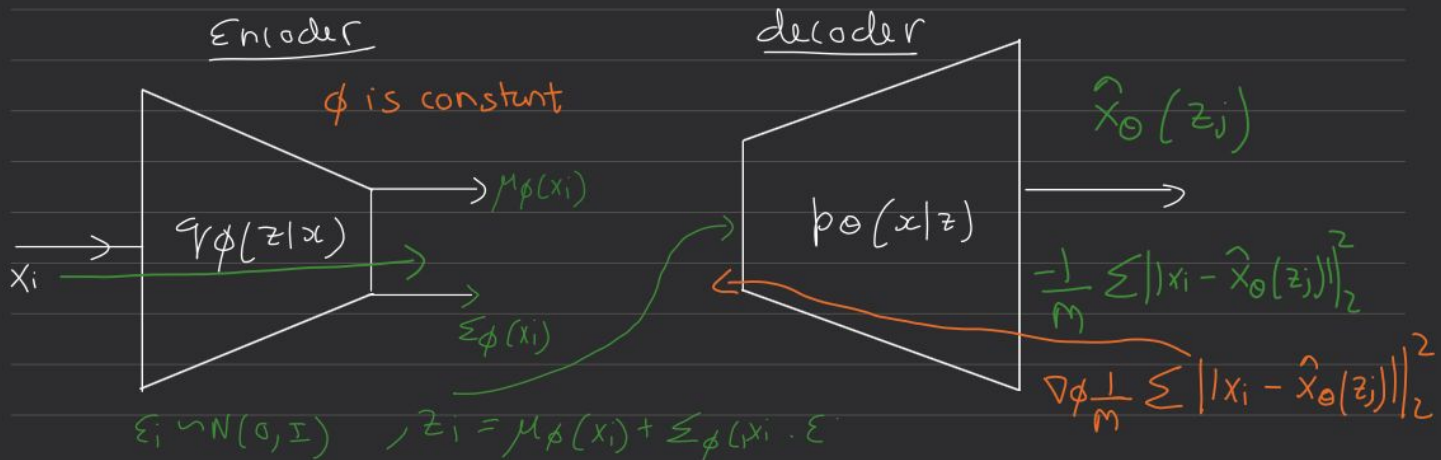
$$= \underbrace{-\frac{1}{2} \log |\epsilon_\phi(x)| - d + \text{tr}[\epsilon_\phi(x)] + \|\mu_\phi(x)\|_2^2}_{\text{Can be Estimated}}$$

hence,

$$\nabla_\phi D_{KL}[q_\phi(z|x) || p_\theta(z)]$$

$$\nabla_\phi \left[ -\frac{1}{2} \log |\epsilon_\phi(x)| - d + \text{tr}[\epsilon_\phi(x)] + \|\mu_\phi(x)\|_2^2 \right] = \nabla_\phi(KL)$$

## 2. The Decoder



one training step for the decoder:

$$\theta^{t+1} \leftarrow \theta^t + \alpha \nabla_\theta J_\phi(q_\phi)$$