

Variational Autoencoders

Reference

- Auto-encoding variational Bayes. Diederik P Kingma and Max Welling, ICLR, 2014.
- Tutorial on Variational Autoencoders, Carl Doersch, arxiv:1606.05908v2.
- Variational Inference: A Review for Statisticians, David M Blei, Alp Kucukelbir, Jon D McAuliffe, arXiv:1601.00670v9.
- <https://www.youtube.com/watch?v=uaaqyVS9-rM&feature=youtu.be&t=19m42s> (mathematical derivation)
- <https://www.jeremyjordan.me/variational-autoencoders/> (good read)

Motivation

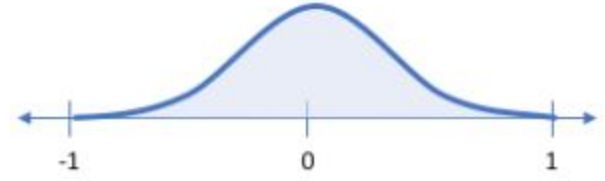
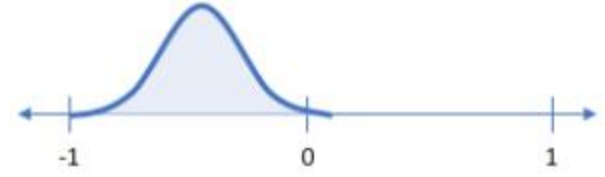
- Autoencoder has “encoding” part and “decoding” part.
 - Encoding part can be “learning” latent space
 - From latent space, observation is decoding the latent space
 - Latent space can only be observed “probabilistically”
 - Assign probability distribution to the interface of encoder and decoder
 - Simple autoencoder is learning definite value of encoder and decoder
 - Variational autoencoder is learning probability distribution of latent space (interface between encoder and decoder)



Smile (discrete value)



Smile (probability distribution)



vs.

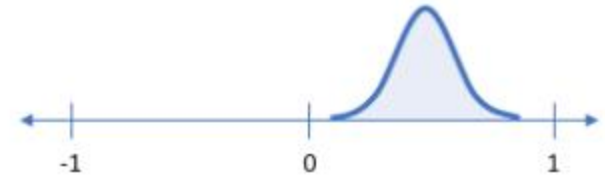
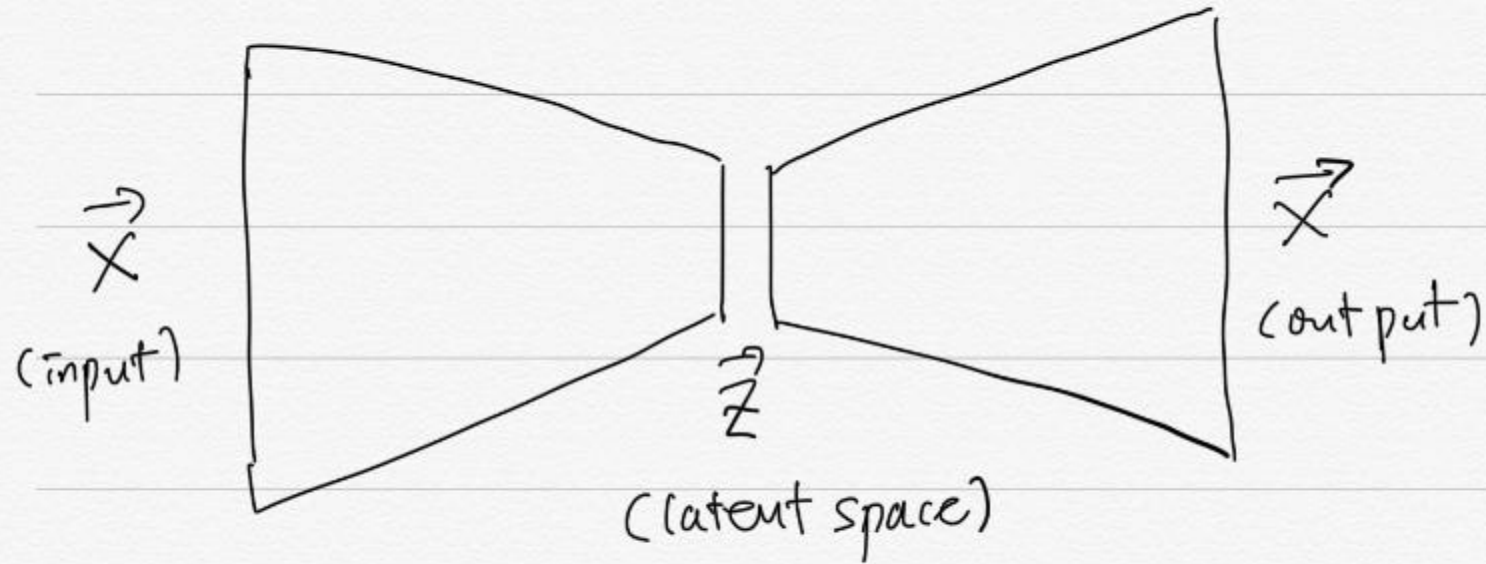


Image from <https://www.jeremyjordan.me/variational-autoencoders/>

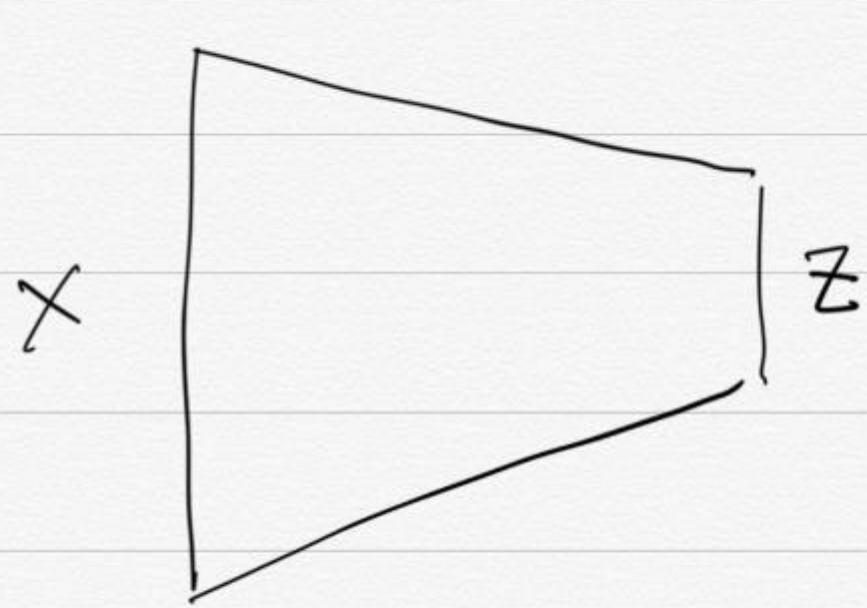
Autoencoder



encoding →

→ decoding

"Encoding"



x z :

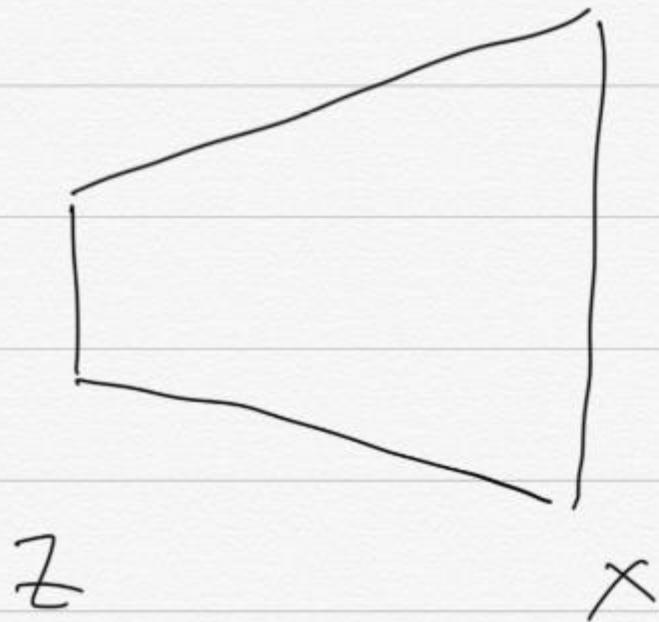
$$P(z|x) = \frac{P(x|z) P(z)}{P(x)}$$

↓
posterior

Prior
↙

Goal: we want to learn $P(z|x)$
posterior distribution

"Decoding" : reconstruct data



$$p(x|z)$$

Ⓢ sample \vec{z} from
distribution

$p(z|x)$ and then

insert to "deterministic"

decoder $p(x|z)$

Conventional autoencoder:

doesn't learn the distribution of

$P(Z|X)$ but simply train $P(Z|X) \rightarrow P(X|Z)$

$\Rightarrow P(X|X)$

deterministically.

Challenge:

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

← prior: given

learnable

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

intractable integration

Solution:

(1) Monte Carlo (evaluate $\int p(x|z)p(z)dz$
directly \rightarrow computationally
expensive)

(2) Variational Inference:

approximate $p(z|x)$ with $q(z|x)$.

choose $q(z|x)$ in tractable form

Measure of "closeness"

KL divergence

$$KL(q(x) \parallel p(x)) = - \sum_x q(x) \log \left[\frac{p(x)}{q(x)} \right]$$

$$\geq 0$$

$$\neq KL(p(x) \parallel q(x))$$

Goal: minimize KL divergence

$$KL(q(z|x) || p(z|x)) = - \sum_z q(z|x) \log \left[\frac{p(z|x)}{q(z|x)} \right]$$

↑ ↑
x is given

$$= - \sum_z q(z|x) \log \left[\frac{p(x, z)}{q(z|x)} \frac{1}{p(x)} \right]$$

$$= - \sum_z q(z|x) \log \left[\frac{p(x, z)}{q(z|x)} \right] + \underbrace{\sum_z q(z|x) \log p(x)}_{= 1}$$

$$= - \sum_z q(z|x) \log \left[\frac{p(x, z)}{q(z|x)} \right] + \log p(x)$$

$$= KL(q(z|x) \parallel p(z|x))$$

$$\Rightarrow KL(q(z|x) || p(z|x)) + \sum_z q(z|x) \log \left[\frac{p(x, z)}{q(z|x)} \right]$$

$= \log p(x)$: constant. (independent of
choice of $q(z|x)$)

\Rightarrow minimize KL divergence is equivalent to

$$\text{maximize } \mathcal{L} = \sum_z q(z|x) \log \left[\frac{p(x, z)}{q(z|x)} \right]$$

Since KL divergence ≥ 0 ,

$L \leq \log p(x)$: L is a lower bound
of $\log p(x)$

$\Rightarrow L$ is a "Variational" lower bound of
 $\log p(x)$.

Now, maximize L :

$$\sum_z q(z|x) \log \frac{p(x, z)}{q(z|x)}$$

$$p(x, z) = p(x|z)p(z)$$

$$= \sum_z q(z|x) \log \frac{p(x|z)p(z)}{q(z|x)} \quad \leftarrow \begin{array}{l} \text{prior} \\ \text{of latent space } z \end{array}$$

$$= \sum_z q(z|x) \log p(x|z) - \underbrace{\left(- \sum_z q(z|x) \log \frac{q(z|x)}{p(z)} \right)}_{= KL(q(z|x) || p(z))}$$

$$= \mathbb{E}_{z \sim q} [\log p(x|z)] - KL(q(z|x) || p(z))$$

Maximize :

$$\underbrace{\mathbb{E}_{z \sim q} [\log p(x|z)]}_{\Downarrow} - \underbrace{KL(q(z|x) \parallel p(z))}_{\Downarrow}$$

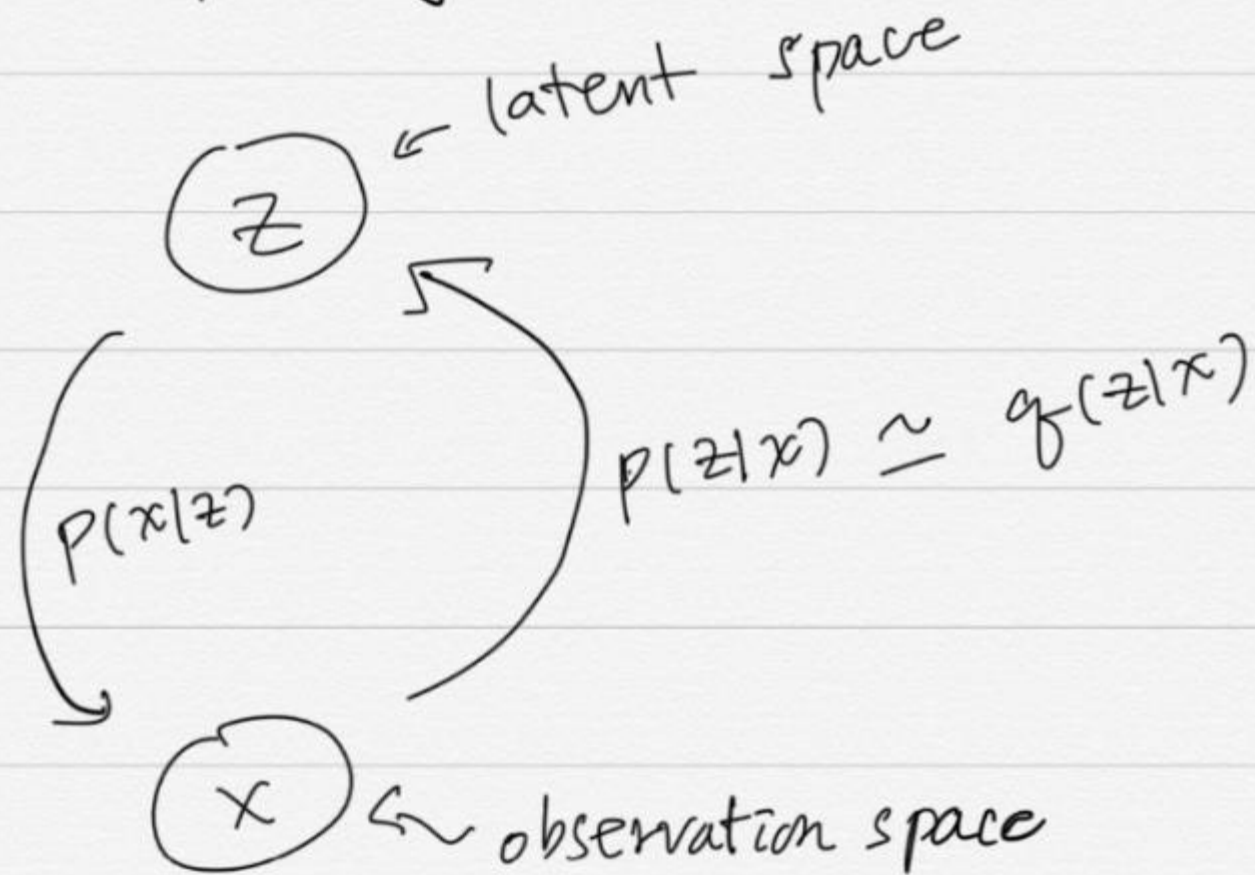
maximize log likelihood
of reconstruction

minimize KL

→ make $q(z|x)$

as similar as to prior $p(z)$

Graphically :



$p(x|z)$

→ learning is doable

by maximizing
log expectation

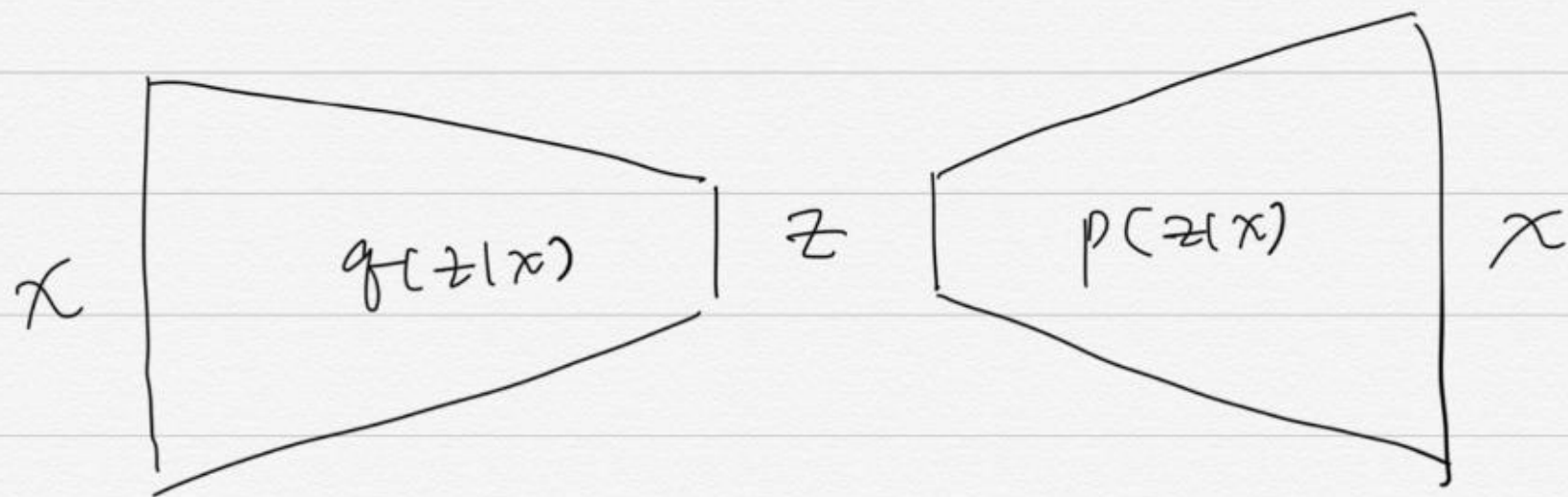
$p(z|x)$

→ learning is hard

so approximate by

$q(z|x)$

Variational Autoencoder



encoder: $-KL(q(z|x) || p(z))$

make $q(z|x)$ similar
to "my choice" $p(z)$

decoder: $\mathbb{E}_z [\log p(x|z)]$

minimize reconstruction error
deterministic algorithm.

$$(1) \text{ maximize } \mathbb{E}_z [\log p(x|z)]$$

$$\text{if } p(x|z) \sim \mathcal{N}(f(z), I) : f(z) = \text{decoder learning} \\ - |x - f(z)|^2 \\ \sim e$$

$$\Rightarrow \min \left[\sum_z |x - f(z)|^2 \right] \rightarrow \text{Standard square loss minimization}$$

$$\text{if } p(x|z) \sim \text{Bernoulli} (\text{learning } f(z) \in [0, 1])$$

$$\Rightarrow \min \left[\sum \left[x \log f(z) + (1-x) \log (1-f(z)) \right] \right]$$

\rightarrow standard cross entropy minimization

$$(2) \quad \text{minimize} \quad KL(q(z|x) || p(z))$$

: Choose $p(z) \sim \mathcal{N}(0, I) \rightarrow$ prior choice

$$q(z|x) \sim \mathcal{N}(\mu, \Sigma)$$

\hookrightarrow choose $D \times D$ diagonal

\Rightarrow Tractable choice (Gaussian)

$$KL(q(z|x) \parallel p(z))$$

$$= \frac{1}{2} \left(\text{tr}(\Sigma) + \mu^T \mu - D - \log \det(\Sigma) \right)$$

(general form of $q(z|x) \sim \mathcal{N}(\vec{\mu}, \Sigma)$
and $p(z) \sim \mathcal{N}(0, I)$)

learn μ & Σ from gradient descent algorithm

and generate $\vec{z} \sim \mathcal{N}(\mu, \Sigma)$

feed generated \vec{z} in to decoder $p(x|z)$

Or, generate z from $N(0, I)$ and feed into $p(x|z)$ to generate artificial output.

"Generative Model"

challenge: how to train the network?

Each time, input $X \rightarrow$ compute μ & Σ

\rightarrow "sample" z from $N(\mu, \Sigma)$

\rightarrow feed z to decoder $p(x|z)$

\rightarrow compute log likelihood

\rightarrow Backpropagate: impossible for
"sampling" stage

Solution = Reparametrization Trick:

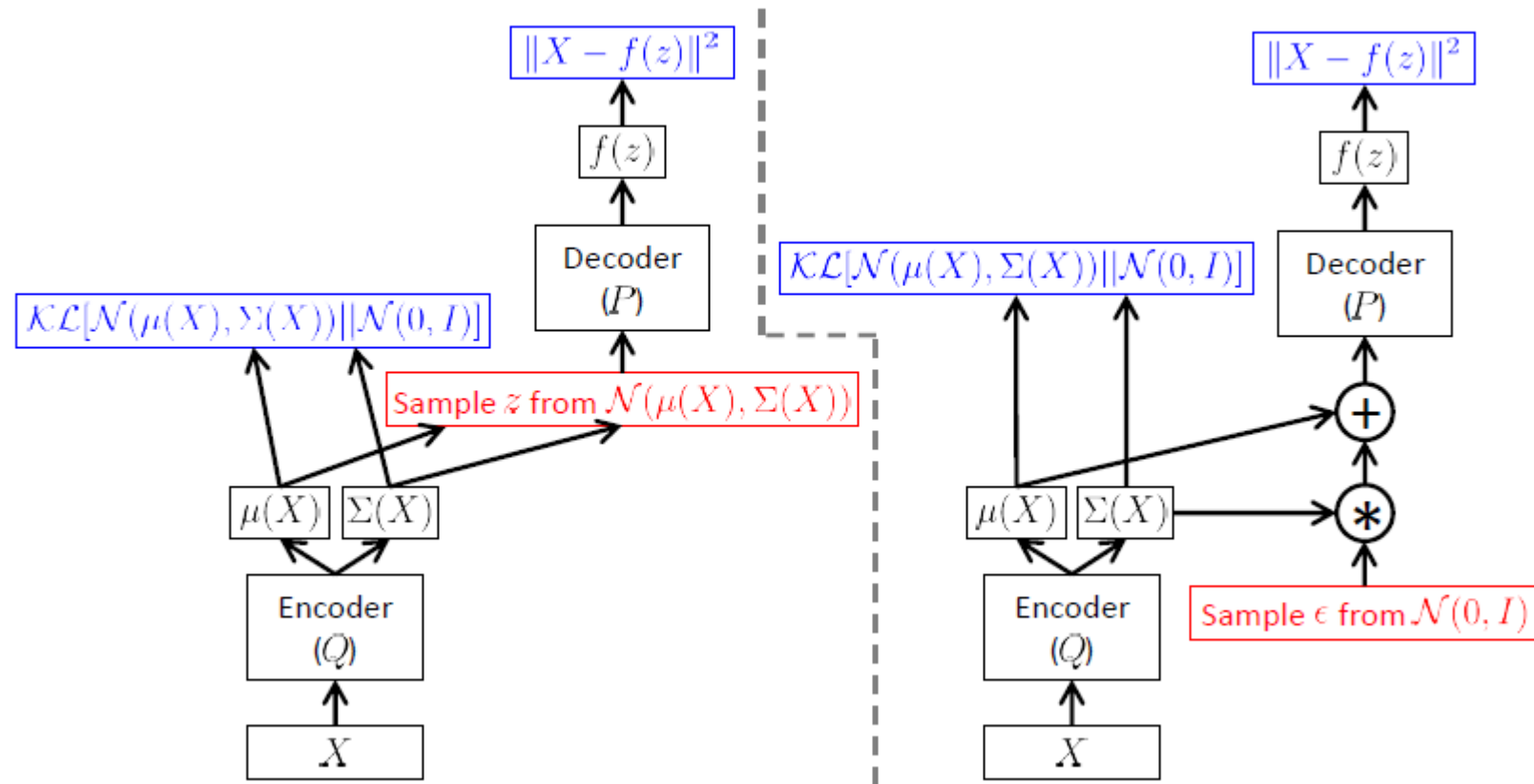


Image: Figure 4 from arXiv:1606.05908v2

$$\bar{z} = \mu + \Sigma^{1/2} \cdot \epsilon \rightarrow z \sim N(\mu, \Sigma)$$

Since Σ is positive semi-definite,
train $\log(\Sigma)$ instead and exponentiate.

Learn μ and σ from input X and treat ϵ as another input variables to generate random z

Q: Why $P(z) \sim N(0, I)$ prior would be enough to generate complex output?

A: With enough non-linearity in $P(x|z)$ (decoder network), we can learn complex structure from simple $N(0, I)$

Example :

$$z \sim N(0, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix})$$

$$u = \frac{z}{10} + \frac{z}{\|z\|} \rightarrow \text{distributed on a ring.}$$

Meaning of $\underbrace{\mathbb{E}_z [\log p(x|z)]}_{\textcircled{1}} - \underbrace{KL(q(z|x) || p(z))}_{\textcircled{2}}$

① : Try to reconstruct input as much as possible (usual learning algorithm)

② : Enforce $q(z|x)$ to have simple distribution
→ Regularization.