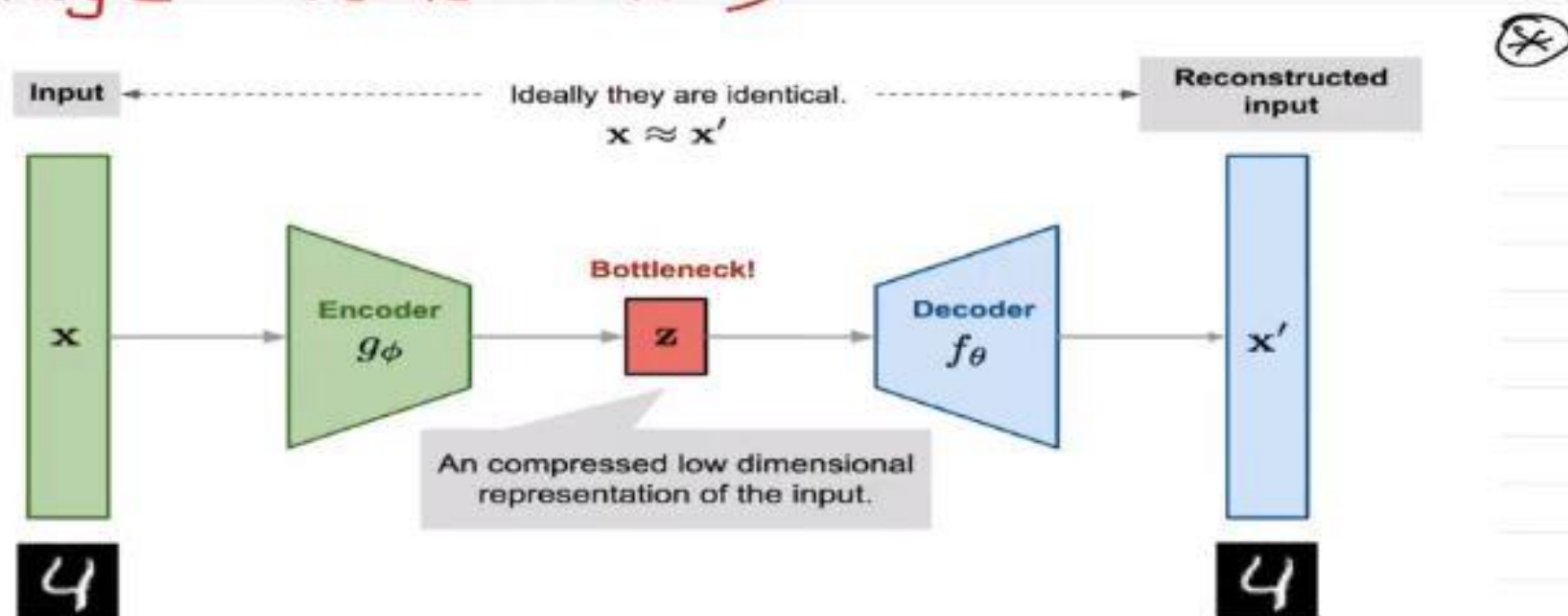


Variational Autoencoders

- ① Review of Stacked Autoencoders
- ② Basics of Probability
- ③ K L Divergence & its significance
- ④ Derivation of Loss function for Variational Autoencoders

Stacked Autoencoders

(Image Reconstruction)



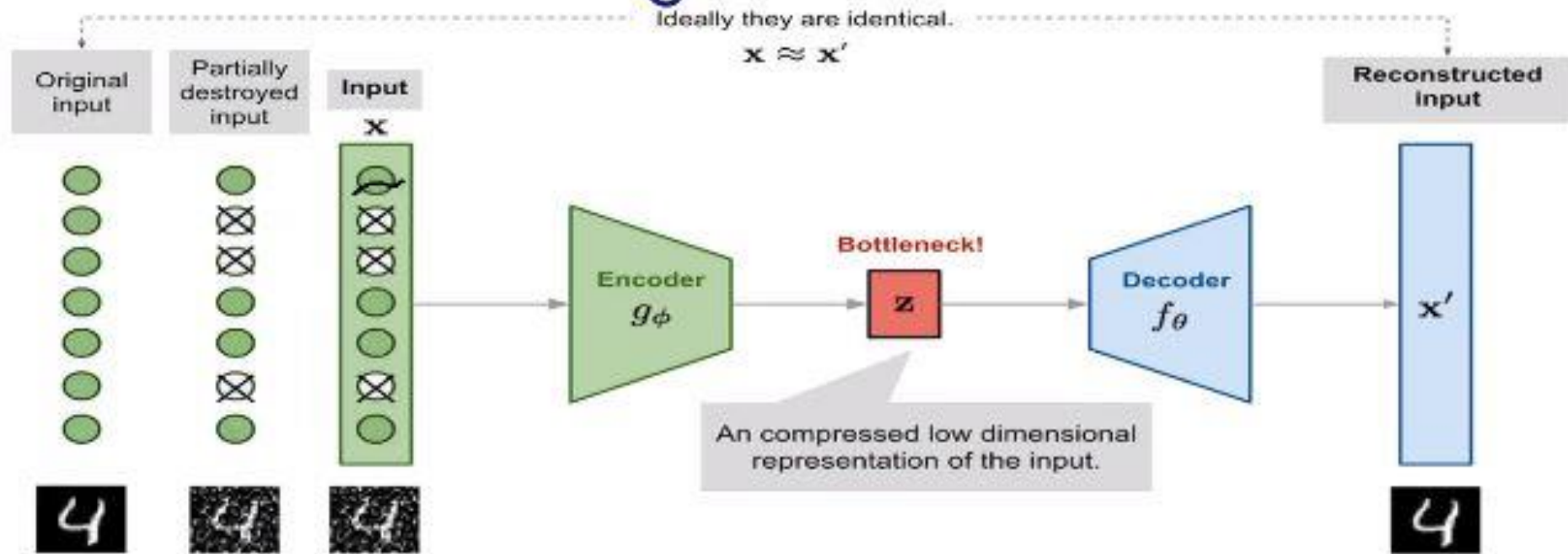
Cost function:
$$L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^{(i)} - f_\theta(g_\phi(x^{(i)}))]^2$$

⊗ from lilianweng github account

11

DL-16

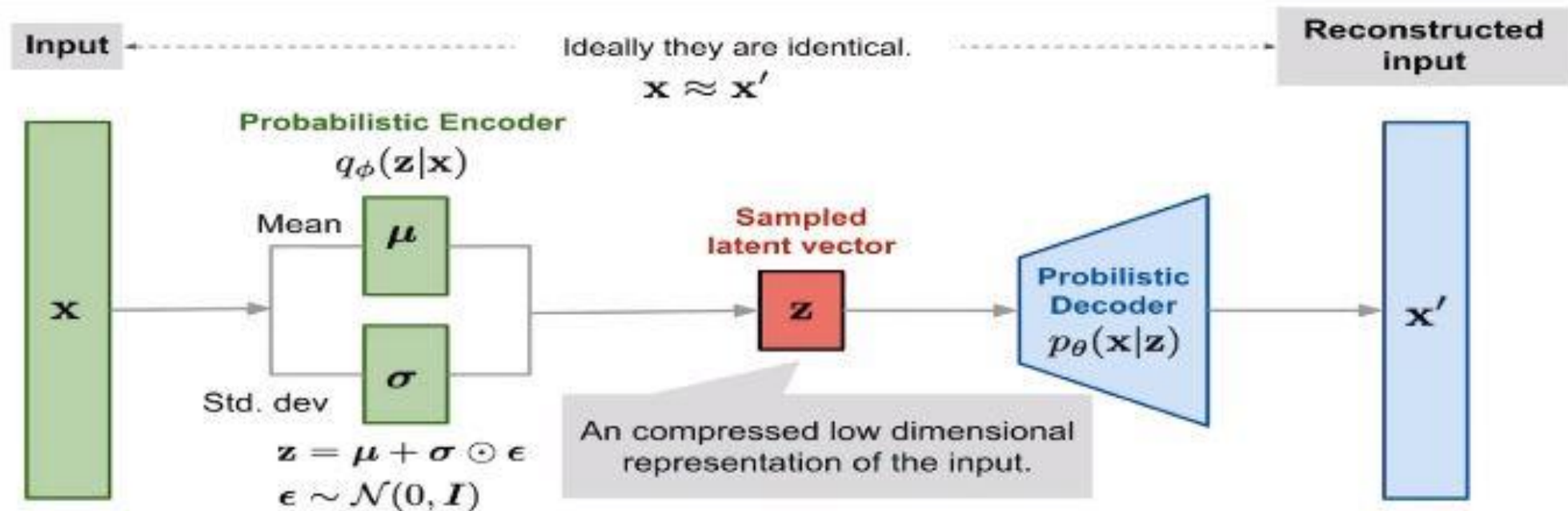
Denoising Autoencoder



$$\tilde{x}^{(i)} \sim \mathcal{X}(\tilde{x}^{(i)} | x^{(i)})$$

$$Loss: L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^{(i)} - f_\theta(g_\phi(\tilde{x}^{(i)}))]^2$$

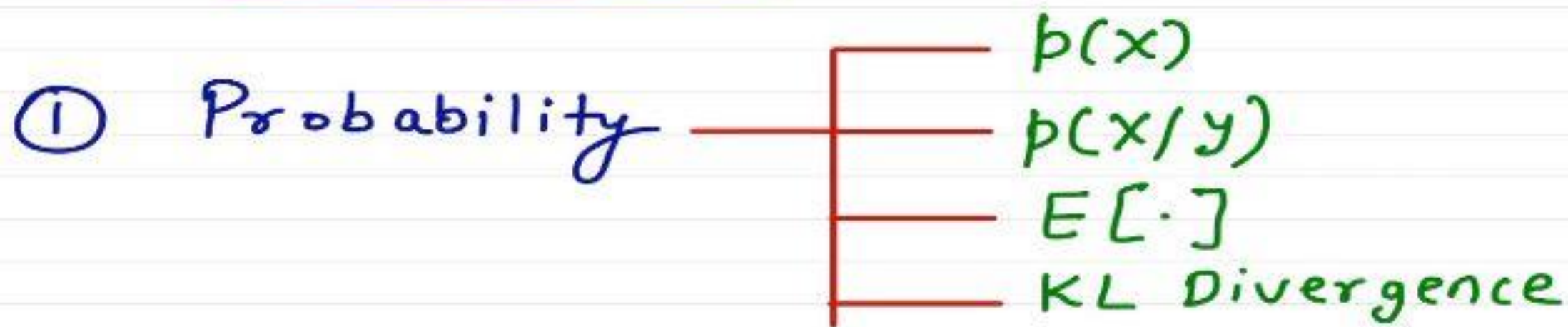
Variational Autoencoders



$$\text{Loss} = \mathcal{L}(\theta, \phi) = -\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log p_\theta(x|z) \right] + D_{KL}(q_\phi(z|x) \| p_\theta(z))$$

Pre-requisite

VA



$p(x)$: defines the probability of random variable x

$p(x/y)$: defines as the probability of random variable x provided y has happened
Also called as conditional probability

Kun

$$p(y/x) = \frac{p(x/y) p(y)}{p(x)} \rightarrow \text{Baye's Theorem}$$

Likelihood ratio

prior probability

posterior probability

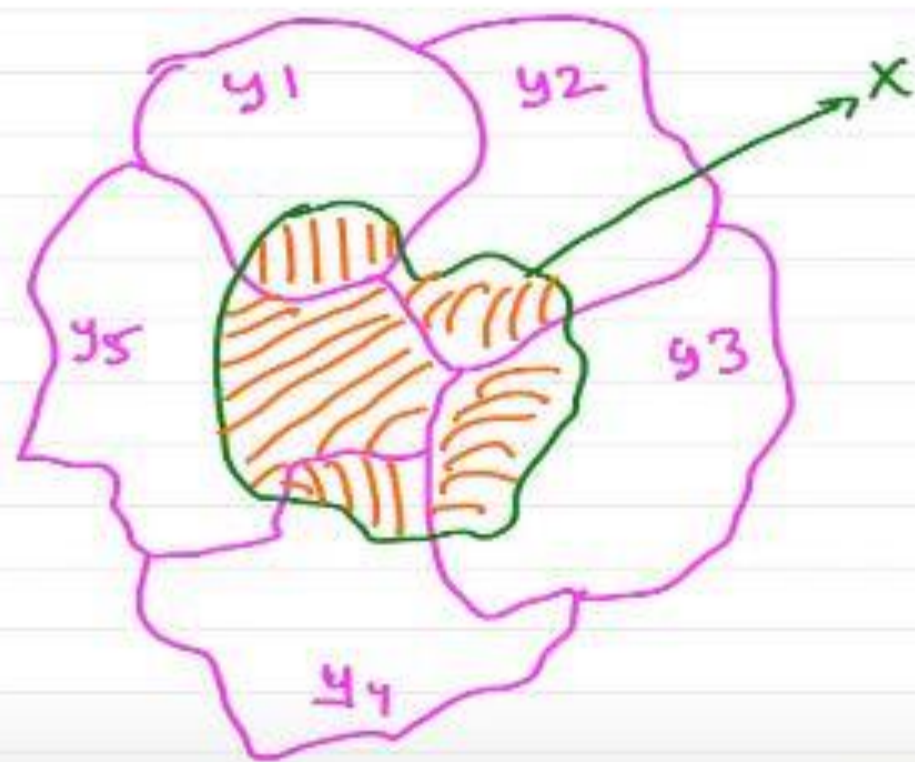
$$= \frac{p(x, y)}{p(x)} \rightarrow \text{joint distribution}$$

①

→ Theorem of Total probability.

Let y_1, y_2, \dots, y_N be a set of mutually exclusive events (i.e. $y_i \cap y_j = \emptyset$) & event X is the union of N mutually exclusive events, then

$$P(X) = \sum_{i=1}^N P(X/y_i) P(y_i) \quad \text{--- ②}$$



$$P(x) = \sum_{i=1}^4 P(x, y_i)$$

$$= \sum_{i=1}^4 P(x/y_i) P(y_i)$$

y_1, y_2, \dots, y_4



So substituting ② in ① results in

$$\left\{ p(y/x) = \frac{p(x/y) p(y)}{\sum_{i=1}^n p(x/y_i) p(y_i)} \right\}$$

Expectation of random variable X i.e. $E(X)$

Expected value of random variable is a weighted average of the possible values of X can take, each value being weighted according to the probability of that

event defined as

$$E(x) = \sum_{i=1}^k x_i P(x=x_i)$$

Q When a die is tossed once. What is the probability of getting 3.

Ans Sample space = $\{1, 2, 3, 4, 5, 6\}$, $P(3) = \frac{1}{6}$

Q2 In tossing a fair die, what is the probability the 3 has occurred conditioned on the toss being odd.

A Since, we are given that odd number has occurred the sample space reduces from $\{1, 2, 3, 4, 5, 6\}$ to $\{1, 3, 5\}$. Hence the probability of 3 in this reduced sample space is $\frac{1}{3}$. It can be observed there is increase in the probability compared to the earlier case. Why?

Q3 Let X represent the outcome of a fair six sided die. What is the $E(X)$?

Ans $X = \{1, 2, 3, 4, 5, 6\}$ $P(X) = 1/6$

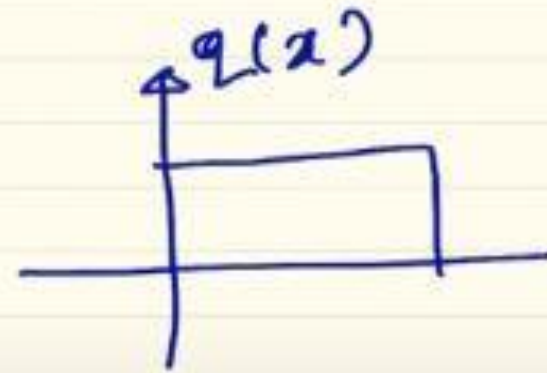
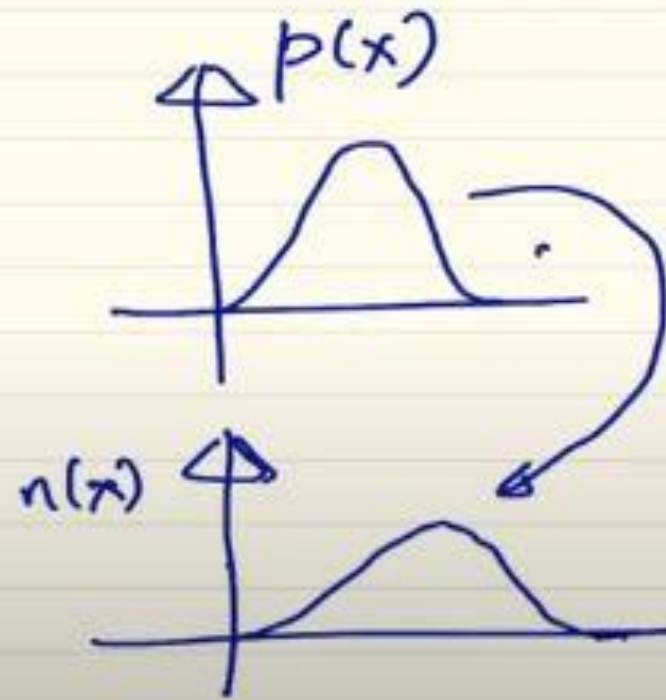
$$E(X) = \sum_{i=1}^6 x P(X) = 1 \times \frac{1}{6} + 2 \times \frac{1}{6} + 3 \times \frac{1}{6} + 4 \times \frac{1}{6} + 5 \times \frac{1}{6} + 6 \times \frac{1}{6} = 3.5$$

K-L Divergence



Kullback-Leibler divergence (K-L) is a measure of how one probability distribution is different from the second. For the discrete probability distribution P & Q , the K-L divergence between P & Q is defined as

K-L Divergence

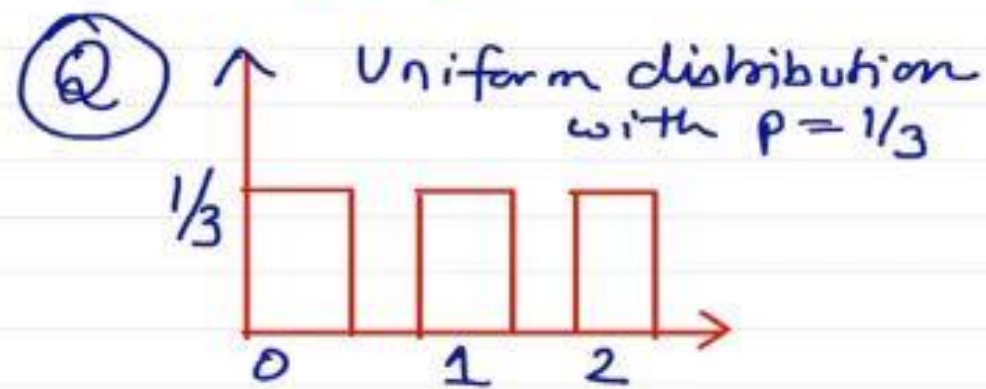
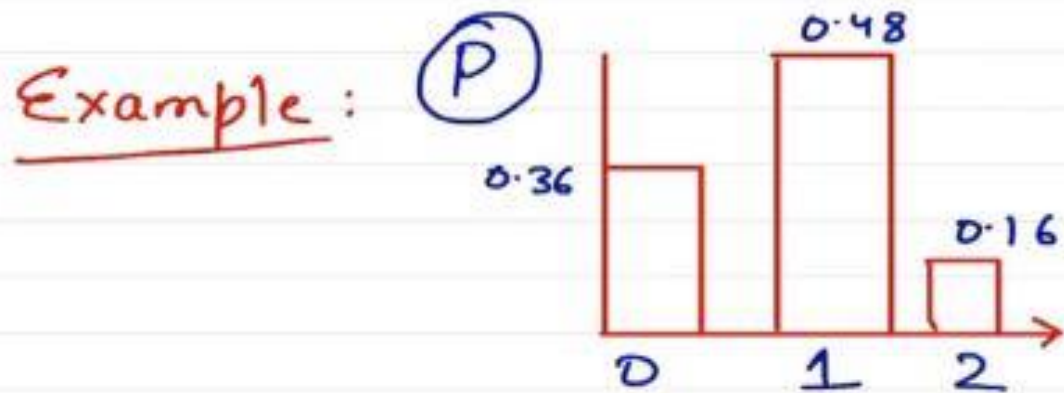


high value

Zero

$$\rightarrow D_{KL}(P \parallel Q) = \sum_x P(X=x) \log \left(\frac{P(X=x)}{Q(X=x)} \right)$$

$$\equiv \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$



$$D_{KL}(Q \parallel P) = \frac{1}{3} \ln \left(\frac{0.333}{0.36} \right) + \frac{1}{3} \ln \left(\frac{0.333}{0.48} \right) + \frac{1}{3} \ln \left(\frac{0.333}{0.16} \right)$$

$$= 0.09637 \text{ nats}$$

Properties:

- ① $KL(P||Q) \text{ or } KL(Q||P) \geq 0$
- ② $KL(P||Q) \neq KL(Q||P)$ (Not symmetric)

Suppose we have two multivariate normal distributions defined as

$$p(x) = N(x; \mu_1, \Sigma_1)$$

$$q(x) = N(x; \mu_2, \Sigma_2)$$

where μ_1 & μ_2 are the means & Σ_1, Σ_2 are the covariance matrix

And the multivariate normal density is defined as

$$N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

if the two distributions have the same dimension k .

$$D_{KL}(p(x) || q(x)) = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right]$$

Prove?

Proof: We know

$$KL(P(x) || Q(x)) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \quad \text{--- (1)}$$

We know

$$P(x) = \frac{1}{\sqrt{(2\pi)^K |\Sigma_1|}} \exp \left(-\frac{(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}{2} \right)$$

$$\Rightarrow \log P(x) = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \quad \text{--- (2)}$$

Similarly

$$\log Q(x) = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \quad \text{--- (3)}$$

Eq ① can be rewritten as

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) \log P(x) - \log(Q(x))$$

Substituting ② & ③ in ① results in

$$KL(P(x) \parallel Q(x)) = \sum_x p(x) \left\{ -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| \right.$$

$$- \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{k}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2|$$

$$\left. + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\}$$

which on simplification results in

$$KL(P(x) \parallel Q(x)) =$$

$$\sum_x p(x) \left\{ \frac{1}{2} \log \left[\frac{\Sigma_2}{\Sigma_1} \right] + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right\} \quad \text{--- (7)}$$

Now, let consider part by part

$$\sum_x p(x) \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \equiv E_p \left[\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right]$$

$$\begin{aligned}
 \rightarrow E(x^T A x) &= E(\text{tr}(x^T A x)) \text{ --- (a)} \\
 &= E(\text{tr}(A x x^T)) \text{ --- (c)} \\
 &= \text{tr}(E(A x x^T)) \text{ --- (e)}
 \end{aligned}$$

Let's rewrite again

$$\frac{1}{2} E_p \left[\overbrace{(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}^{\text{scalar}} \right]$$

$$E_p \left[\text{tr} \left(\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right) \right]$$

$$E_p \left[\text{tr} \left(\frac{1}{2} (x - \mu_1) (x - \mu_1)^T \Sigma_1^{-1} \right) \right] \rightarrow \text{(d)}$$

Trace & Expectation trick.

\rightarrow if x is scalar then $E(x) = E(\text{tr}(x))$
Since trace of x is scalar

$$\rightarrow \text{tr}(AB) = \text{tr}(BA) \text{ --- (b)}$$

$$\rightarrow \text{tr}(ABC) = \text{tr}(BCA) \text{ --- (c)}$$

$$= \text{tr}(CAB) \text{ --- (d)}$$

$$\rightarrow \text{tr}(ABC) \neq \text{tr}(ACB)$$

$$E(\text{tr}(x)) = \text{tr}(E(x)) \text{ --- (e)}$$

→ $\rightarrow \text{tr} \left[E_p \left(\frac{1}{2} (x - \mu_1) (x - \mu_1)^T \Sigma_1^{-1} \right) \right] \text{ --- } \textcircled{a}$

→ $\text{tr} \left\{ E_p \left[(x - \mu_1) (x - \mu_1)^T \right] \frac{1}{2} \Sigma_1^{-1} \right\}$

→ Covariance matrix

$\Rightarrow \text{tr} \left[\Sigma_1 \frac{1}{2} \Sigma_1^{-1} \right]$

$\text{tr} \left[I_K \right] \equiv K \text{ --- } \textcircled{2}$

Now Consider the second part

$$\sum_x p(x) \left[\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] \quad \checkmark$$

$$\sum_x p(x) \left\{ \frac{1}{2} [(x - \mu_1) + (\mu_1 - \mu_2)]^T \Sigma_2^{-1} [(x - \mu_1) + (\mu_1 - \mu_2)] \right\}$$

$$\sum_x p(x) \left[\frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + \frac{2}{2} (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$\Rightarrow E_p \left[\frac{1}{2} (x - \mu)^T \Sigma_2^{-1} (x - \mu) + (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

Now Consider the second part

$$\begin{aligned}
 & \sum_x p(x) \left[\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] \quad \checkmark \frac{(A+B)^T \Sigma_2^{-1} (A+B)}{(A^T + B^T) \Sigma_2^{-1} (A+B)} \\
 & \sum_x p(x) \left\{ \frac{1}{2} [(x - \mu_1) + (\mu_1 - \mu_2)]^T \Sigma_2^{-1} [(x - \mu_1) + (\mu_1 - \mu_2)] \right\} \quad \begin{matrix} A^T \Sigma_2^{-1} A \\ B^T \Sigma_2^{-1} B \end{matrix} \\
 & \sum_x p(x) \left[\frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + \frac{2}{2} (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right. \\
 & \quad \left. + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \quad \begin{cases} B^T \Sigma_2^{-1} A \\ A^T \Sigma_2^{-1} B \end{cases} \\
 & \Rightarrow E_p \left[\frac{1}{2} (x - \mu)^T \Sigma_2^{-1} (x - \mu) + (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right. \\
 & \quad \left. + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]
 \end{aligned}$$

Expanding we get

$$E_P \left\{ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) \right\} + E_P \left[\underbrace{(x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}_{\text{Constant}} \right] + E_P \left[\underbrace{(\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}_{\text{Constant}} \right]$$

$$= \text{tr} \left\{ \frac{\Sigma_2^{-1} \Sigma_1}{2} \right\} + \underbrace{(\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}_{E[\text{Constant}] = \text{Constant}} + 0$$

similar to earlier derivation

$\rightarrow \beta$

0 proved on next slide

$$\rightarrow E_p \left[(x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$\left[(E_p(x) - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$\equiv (\mu_1 - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) = 0 \text{ (Proved)}$$

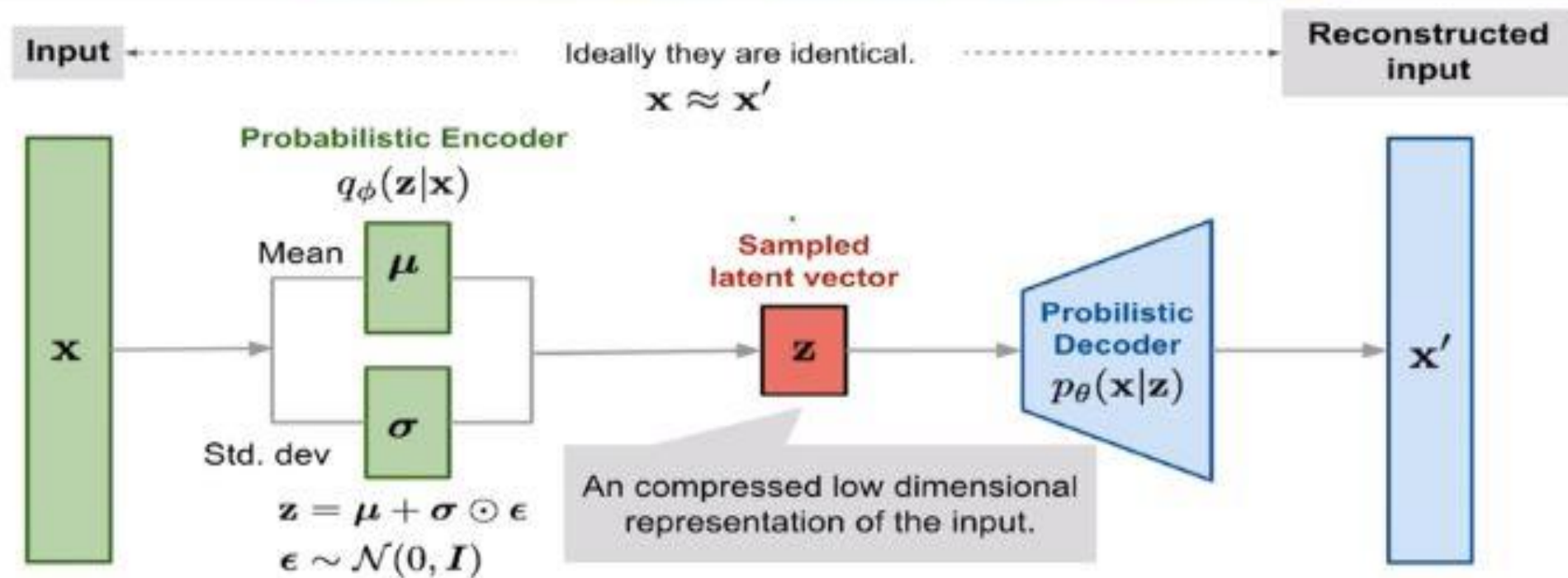
So substituting (2), (4) in (1) we obtain

$$KL(p(x) \parallel q(x)) = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

proved

Variational Autoencoder

#

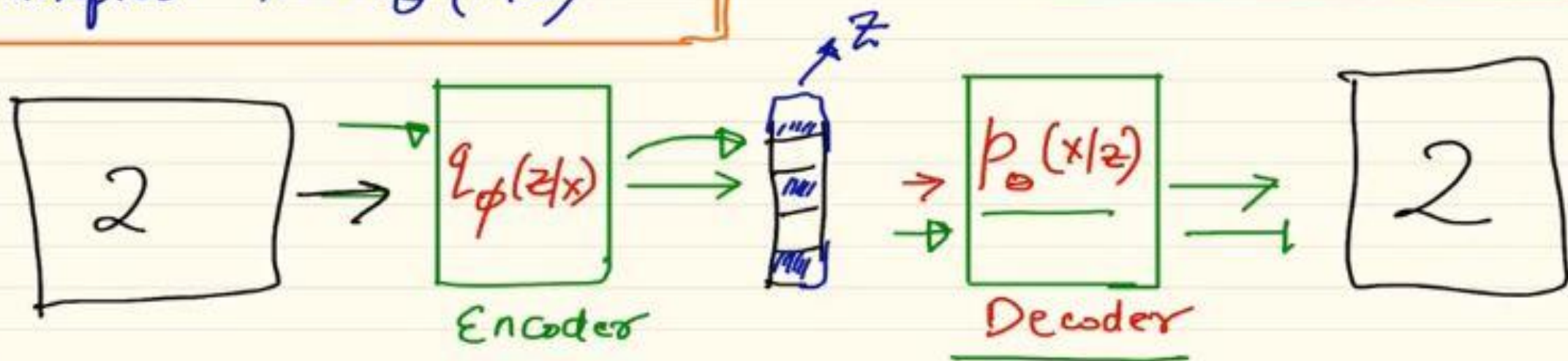


$$\text{Loss} = \mathcal{L}(\theta, \phi) = -\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log p_\theta(x|z) \right] + D_{KL}(q_\phi(z|x) \| p_\theta(z))$$

The goal of VAE

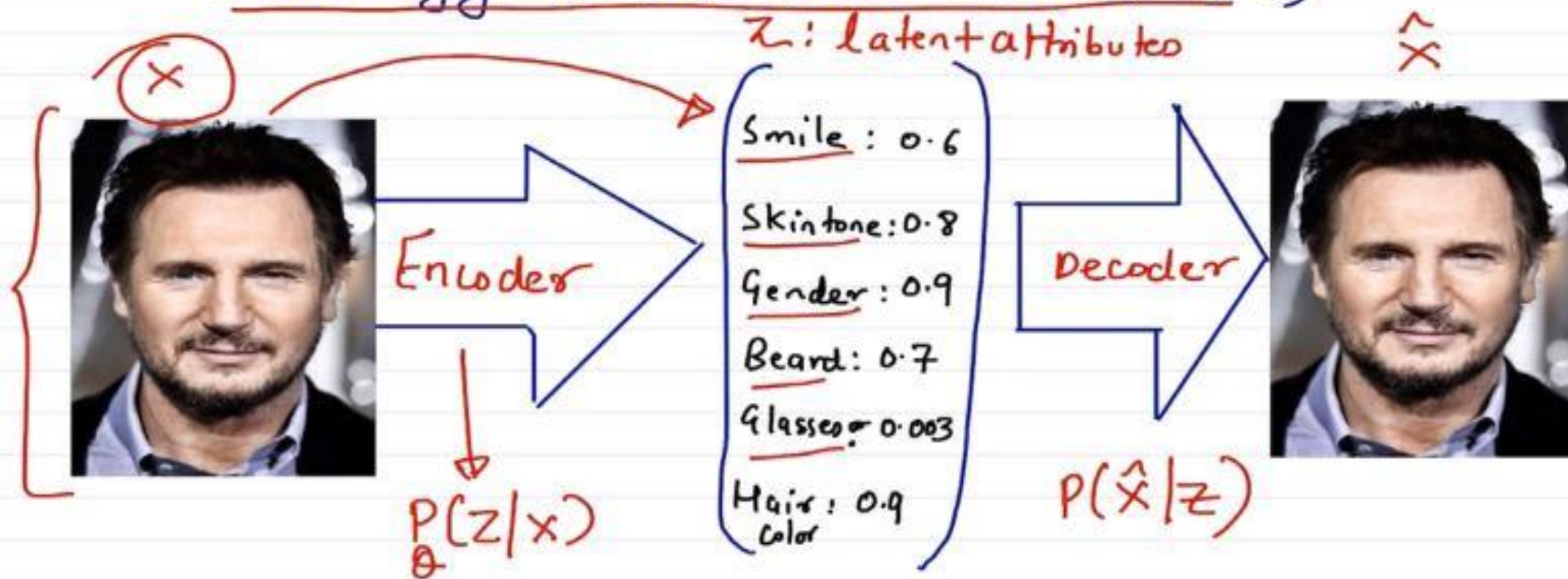
The goal of VAE is to find a distribution $q_{\phi}(z/x)$ of some latent variables which we can sample from $z \sim q_{\phi}(z/x)$ to generate new samples $x' \sim p_{\theta}(x/z)$

Typical Autoencoder



Latent variables can be placed in 2 categories

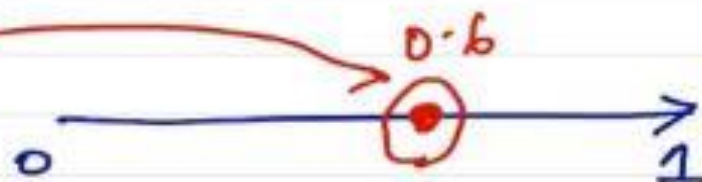
- ① ^{Latent (2)} Variables corresponding to a real feature of the object that have not been measured (may be technology is not available to do that)



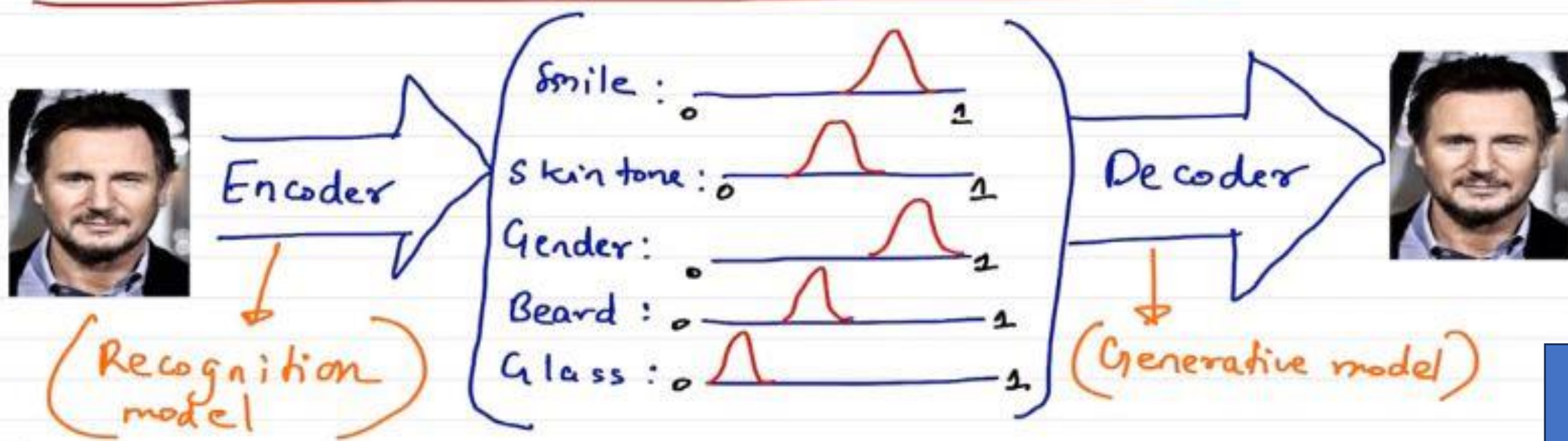
Using VE, we define latent attributes in probabilistic terms.

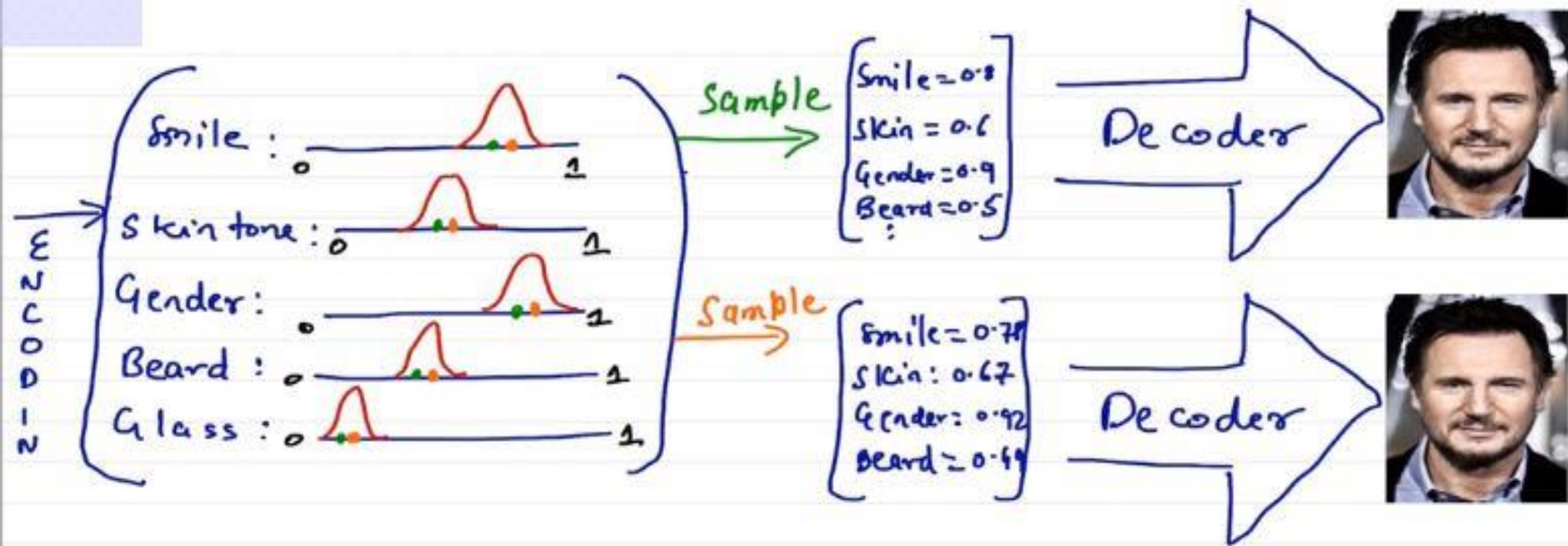
AE
Smile (discrete)

VE
Smile (Probabilistic)



With this approach, we now represent each latent attribute for a given input as a probability distribution. When decoding, we will randomly sample from each latent state distribution to generate a vector as i/p for our decoder model.

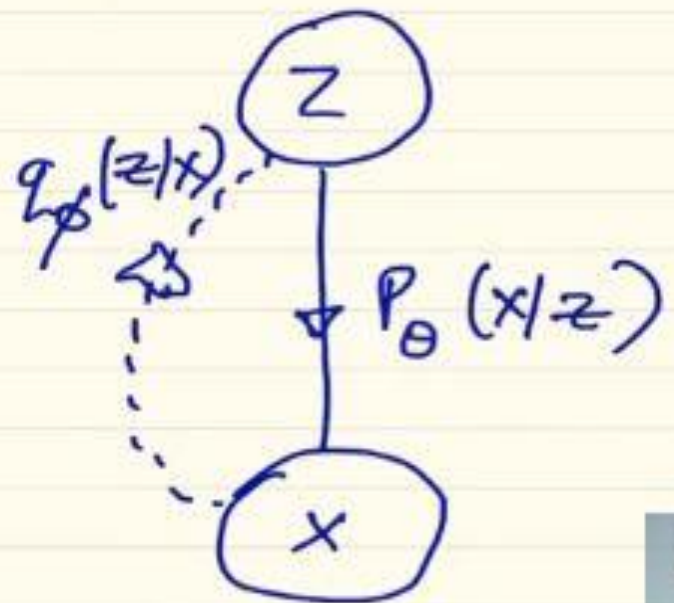
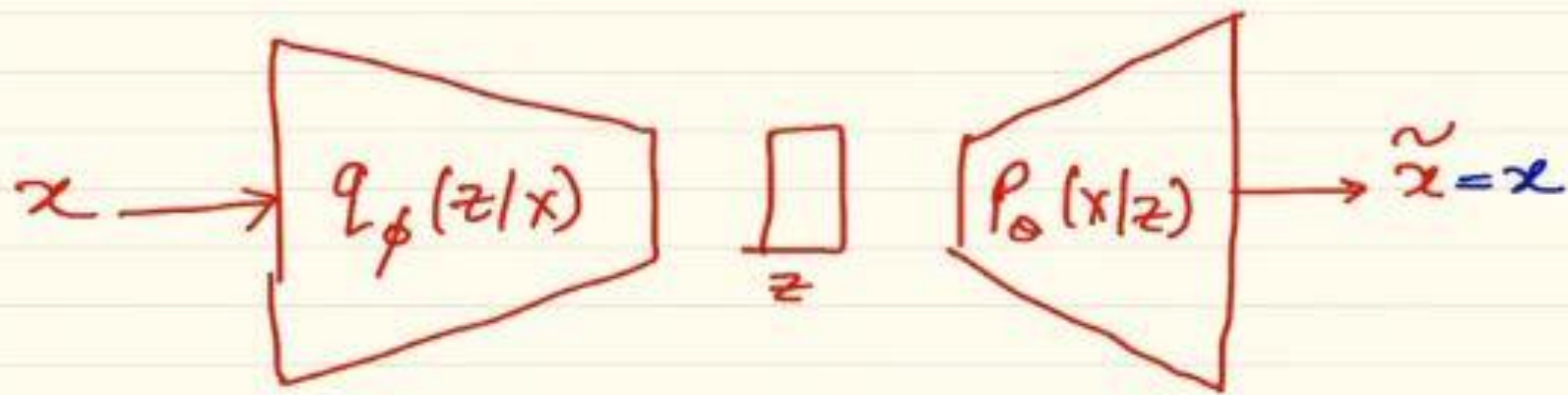




By constructing our encoder model to output range of possible values (a statistically distribution) from which we will randomly sample to feed into our decoder model, The the values which are nearby each other in latent space must correspond ^{to} similar reconstruct

Let's recall the Goal of VAE

The goal of VAE is to find a distribution $q_{\phi}(z/x)$ of some latent variables, which we can sample from $z \sim q_{\phi}(z/x)$, to generate new samples x' from $p_{\theta}(x/z)$

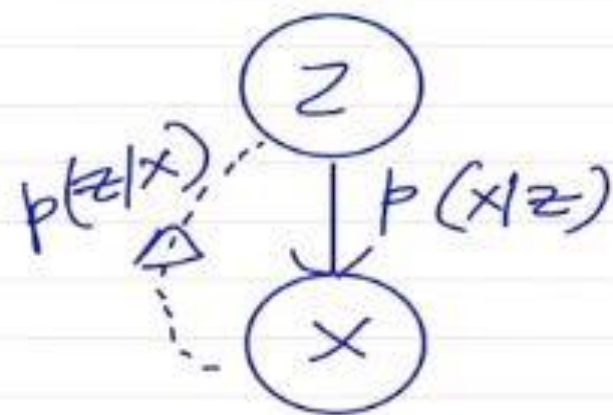


The problem of Approximate Inference

Let x be a set of observed variables & let z be the set of latent variables with joint distribution $p(z, x)$. Then the inference problem is to compute the conditional distribution of latent variables given the observations i.e $p(z/x)$

→ We can write it as

$$p(z/x) = \frac{p(x/z) p(z)}{p(x)} \quad \text{--- (A)}$$



Evaluating (A) is difficult because $p(x)$ cannot be solved

Reason :

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

this integral is not available in closed form or is intractable (i.e. requires exponential time to compute) due to multiple integrals involved for latent variable vector z .

Alternative? The alternative is to approximate $p(z|x)$ by another distribution $q(z|x)$ which is defined in such a way that it has tractable solution. This is done using variational inference (VI). The main idea of

VII is to pose the inference problem as an optimization problem. How? By modelling $P(z|x)$ using $Q(z|x)$ where $Q(z|x)$ has a simple distribution such as Gaussian.

As discussed let us calculate KL b/w $P(z|x)$ & $Q(z|x)$

$$\begin{aligned} D_{KL}(Q(z|x) \parallel P(z|x)) &= \sum_z Q(z|x) \log \left(\frac{Q(z|x)}{P(z|x)} \right) \\ &= E_{z \sim Q(z|x)} \left[\log \left(\frac{Q(z|x)}{P(z|x)} \right) \right] \\ &= E_{z \sim Q(z|x)} \left[\log(Q(z|x)) - \log(P(z|x)) \right] - \textcircled{B} \end{aligned}$$

Substituting (A) in (B) results in

Here $\boxed{\mathbb{Z} = z \sim Q(z/x)}$

$$D_{KL}[Q(z/x) \parallel P(z/x)] = \mathbb{E}_{\mathbb{Z}} \left[\log Q(z/x) - \log \frac{P(x/z) P(z)}{P(x)} \right]$$

$$= \mathbb{E}_{\mathbb{Z}} \left[\log(Q(z/x)) - \log P(x/z) - \log P(z) + \log P(x) \right]$$

Since the expectation is over z & $P(x)$ does not involve z , it can be moved out

$$D_{KL}[Q(z/x) \parallel P(z/x)] - \log P(x) = \mathbb{E}_{\mathbb{Z}} \left[\log Q(z/x) - \log P(x/z) - \log P(z) \right]$$

Rearranging the equations we obtain

$$\textcircled{2}: z \sim Q(z|x)$$

$$\log p(x) - D_{KL}[Q(z|x) \parallel p(z|x)]$$

$$= E_{\textcircled{2}}[\log(p(x|z))] - E_z[\log Q(z|x) - \log p(z)]$$

$$= \boxed{E_{\textcircled{2}}[\log(p(x|z))] - D_{KL}[Q(z|x) \parallel p(z)]}$$

→ This is VAE objective function, where the first term represents the reconstruction likelihood & the second term ensures that our learned distributions Q is similar to the prior distribution P



& loss = - Objective function

$$L(\theta, \phi) = -E_{z \sim Q_{\phi}(z/x)} \left[\log \left(P_{\theta}(x/z) \right) \right] + D_{KL} \left[Q_{\phi}(z/x) \parallel P_{\theta}(z) \right]$$

Proved !!

Also $\log P_{\theta}(x) - D_{KL} \left[Q_{\phi}(z/x) \parallel P_{\theta}(z/x) \right] = -L(\theta, \phi)$

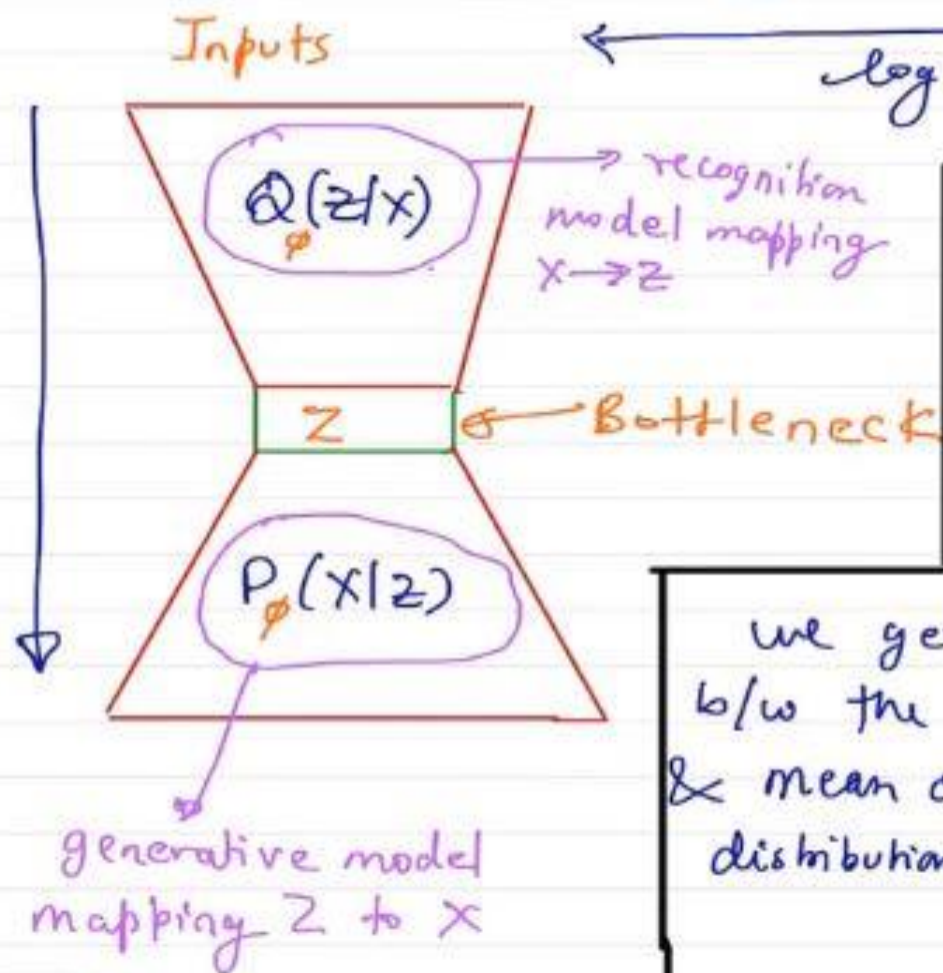
Reconstruction Regularizer

So our target is to find optimal θ , ϕ such that

$$\underbrace{\theta^*}_{\text{Recognition}}, \underbrace{\phi^*}_{\text{generation}} = \left\{ \arg \min_{\theta, \phi} L(\theta, \phi) \right\}$$

MORE INTUITION ABOUT LOSS FUNCTION Kumar

$$L(\theta, \phi) = -E_{z \sim Q_{\phi}(z/x)} \left[\log \left(P_{\theta}(x/z) \right) \right] + D_{KL} \left[Q_{\phi}(z/x) \parallel P_{\theta}(z) \right]$$



log-likelihood

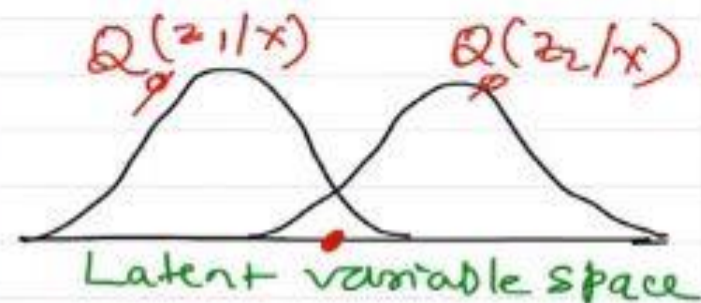
$$P_{\theta}(x/z) = N(\psi_{\theta}(z), \Sigma_{\theta}(z))$$

So when we take log of Gaussian

we get a square error b/w the data sample x & mean of the Gaussian distribution.

Regularizer

Role



KL divergence

not allows pdf of latent variables not collapse with zero variance but penalize if deviates from $N(0,1) = P_{\theta}(z)$

$$\underbrace{p_{\theta}(x|z)} = \frac{1}{\sqrt{(2\pi)^k |\Sigma_{\theta}(z)|}} \exp \left[\frac{(x - \mu_{\theta}(z))^T \Sigma_{\theta}^{-1}(z) (x - \mu_{\theta}(z))}{2} \right]$$

$$\log p_{\theta}(x|z) \propto \underbrace{(x - \mu_{\theta}(z))^T \Sigma_{\theta}^{-1}(z) (x - \mu_{\theta}(z))}$$

Squared Reconstruction error

Let's discuss more about $\{D_{KL}(Q(z|x) || P(z))\}$

Here, $P(z)$ is the latent variable distribution. $P(z)$

→ The easiest choice is $N(0, 1)$. We want $Q(z|x)$ to be as close to $Q(z|x)$ so that we could sample it easily.

→ Having $P(z) = N(0, 1)$ adds another benefit. Let's say if we want $Q(z|x)$ to be Gaussian with parameter $\mu(x)$ & $\Sigma(x)$, the KL divergence has the closed form as derived earlier. ~~In that derivation,~~

$$D_{KL}[N(\mu(x), \Sigma(x)) \parallel N(0, I)] =$$

$$\rightarrow \left\{ \frac{1}{2} \left[\text{tr}(\Sigma(x)) + \mu(x)^T \mu(x) - K - \log |\Sigma(x)| \right] \right\}$$

→ Here K , is the dimension of the Gaussian.

→ $\text{tr}(\Sigma(x))$ is the trace function ; which is sum of the diagonal matrix of $\Sigma(x)$

→ Moreover, determinant $|\Sigma(x)|$ of a diagonal matrix is product of its diagonal so above equation can be simplified as

$$\rightarrow E_p[(x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2)]$$

$$\left[(E_p(x) - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$$\equiv (\mu_1 - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) = 0 \text{ (Proved)}$$

$$\Sigma_1 = \Sigma_{\phi}(x)$$

$$\mu_1 = \mu_{\phi}(x)$$

$$\Sigma_2 = 1$$

$$\mu_2 = 0$$

So substituting (2), (3) in (1) we obtain VAE

$$\left\{ \underbrace{KL(p(x) || q(x))}_{\substack{N \\ \mu_1, \Sigma_1} \quad \substack{\phi \\ \mu_2, \Sigma_2}} = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr}(\Sigma_2^{-1} \Sigma_1) \right] + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right\}$$

proved

$$D_{KL}[N(\mu_\phi(x), \Sigma_\phi(x)) || N(0, I)] =$$

$$\frac{1}{2} \left(\sum_K \Sigma_\phi(x) + \sum_K \mu_\phi^2(x) - \sum_K 1 - \log \prod_K \Sigma_\phi(x) \right)$$

$$= \frac{1}{2} \left(\sum_K \Sigma_\phi(x) + \sum_K \mu_\phi^2(x) - \sum_K 1 - \sum_K \log(\Sigma_\phi(x)) \right)$$

$$= \frac{1}{2} \sum_K \left(\Sigma_\phi(x) + \mu_\phi^2(x) - 1 - \log(\Sigma_\phi(x)) \right)$$

In practice, it is better to model $\Sigma_\phi(x)$ as $\log(\Sigma_\phi(x))$

as it is more numerically stable to take exponent compared to computing log. Hence substituting $\exp(\Sigma(x))$ instead of $\Sigma(x)$ we obtain,

$$D_{KL}[N(\mu(x), \Sigma(x)) || N(0, 1)] = \frac{1}{2} \sum_K [\exp(\Sigma(x)) + \mu^2(x) - 1 - \Sigma(x)]$$

So the final loss function of VAE becomes

$$L(\theta, \phi) = -E_z [\log(p(x|z))] + \frac{1}{2} \sum_K [\exp(\Sigma(x)) + \mu^2(x) - 1 - \Sigma(x)]$$

$\downarrow z \sim Q(z|x)$