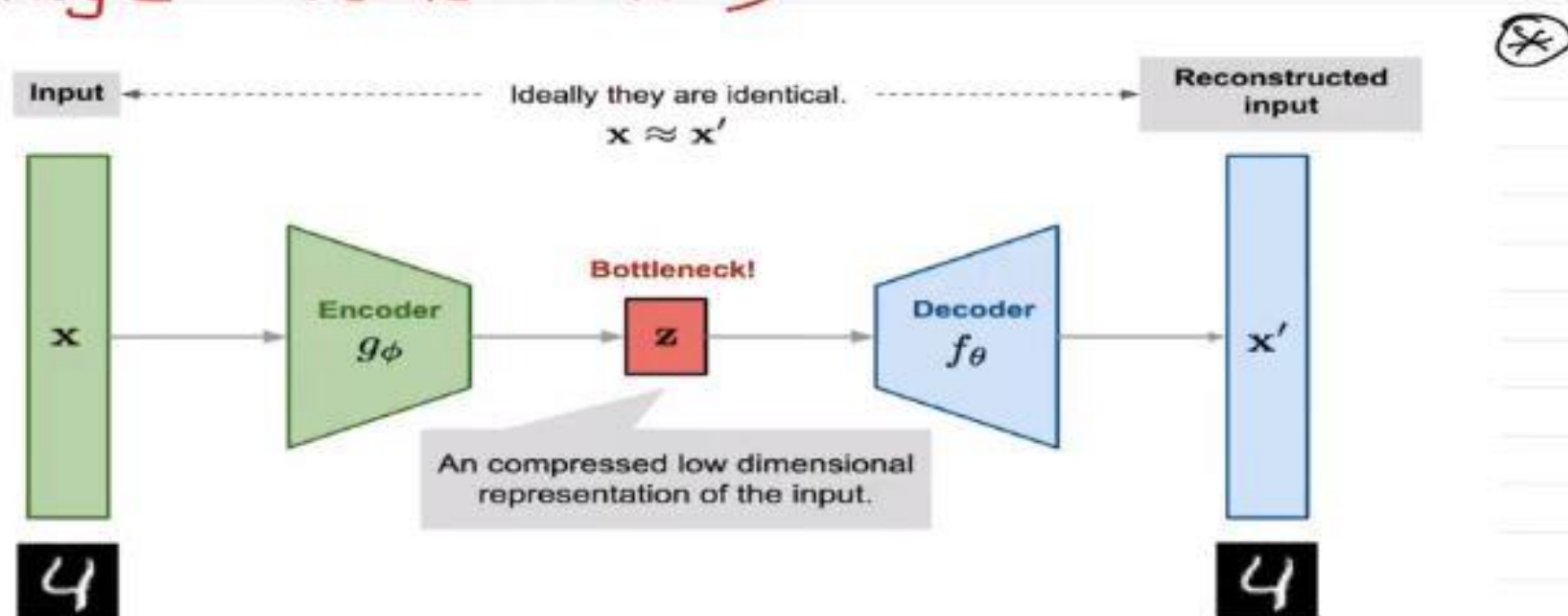


## Variational Autoencoders

- 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)



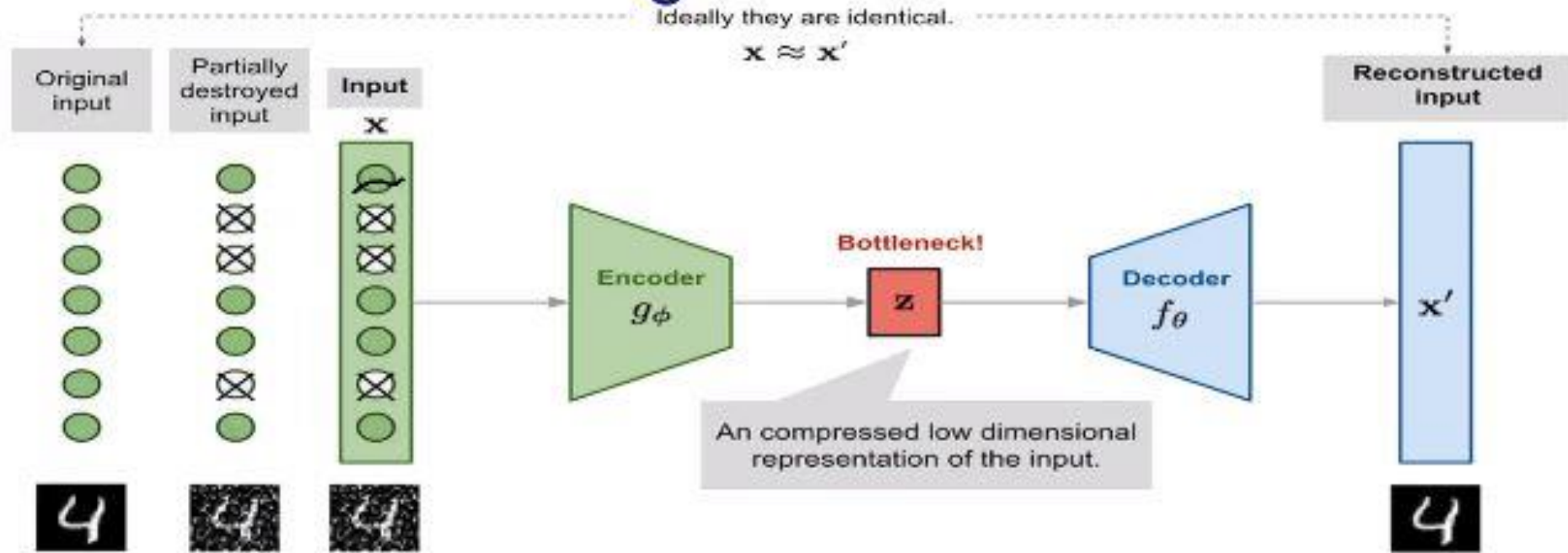
Cost function: 
$$L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^{(i)} - f_\theta(g_\phi(x^{(i)}))]^2$$

⊗ from lilianweng github account

11

DL-16

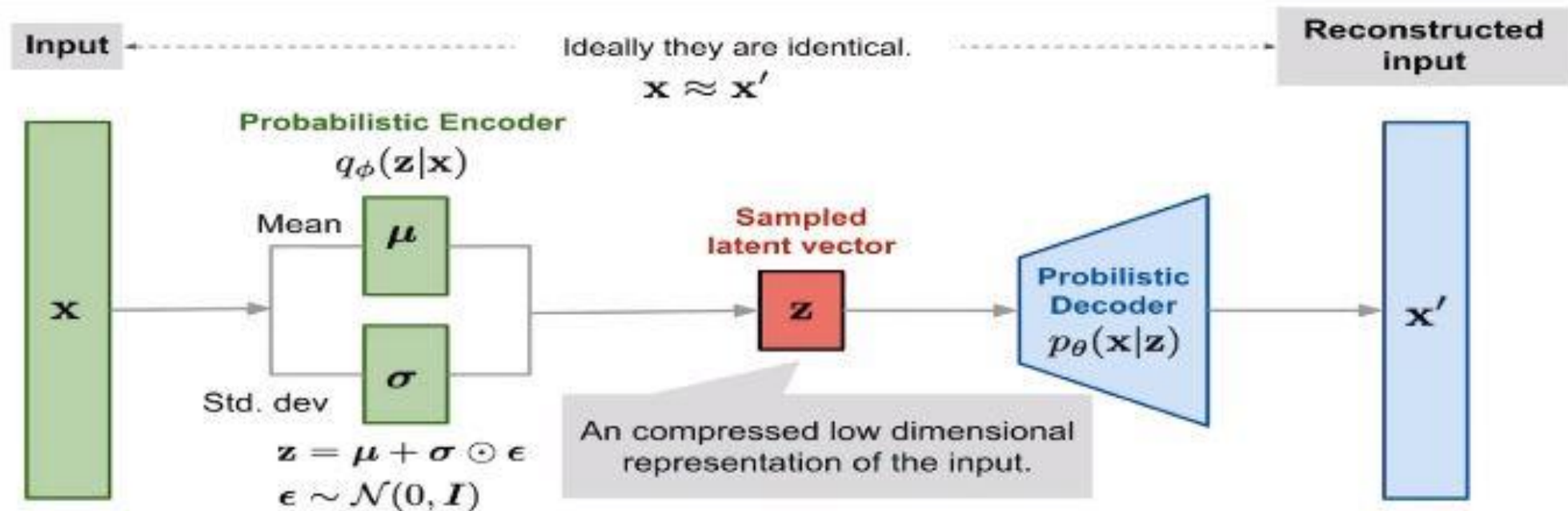
# Denoising Autoencoder



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

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

# Variational Autoencoders

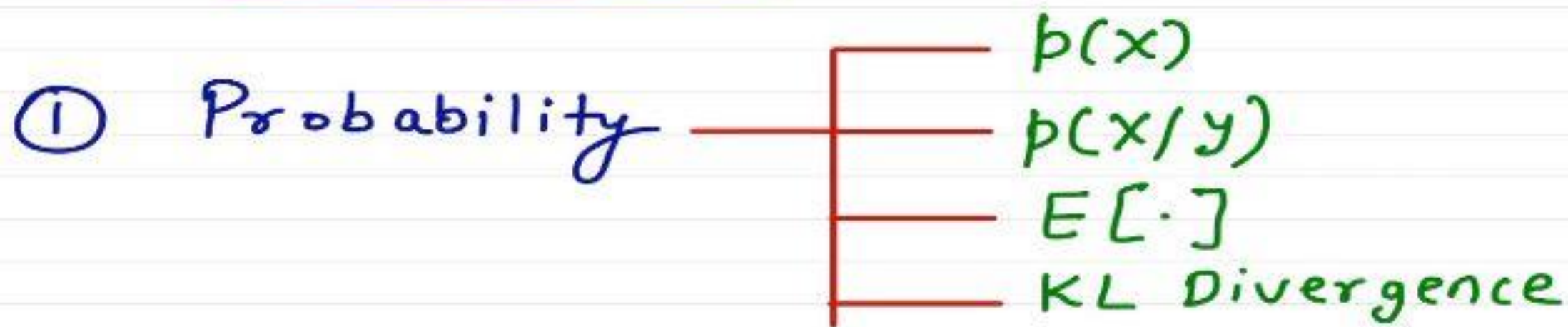


$$\text{Loss} = \mathcal{L}(\theta, \phi) = -\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] + D_{KL}(q_\phi(z|x) \| p_\theta(z))$$



# 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

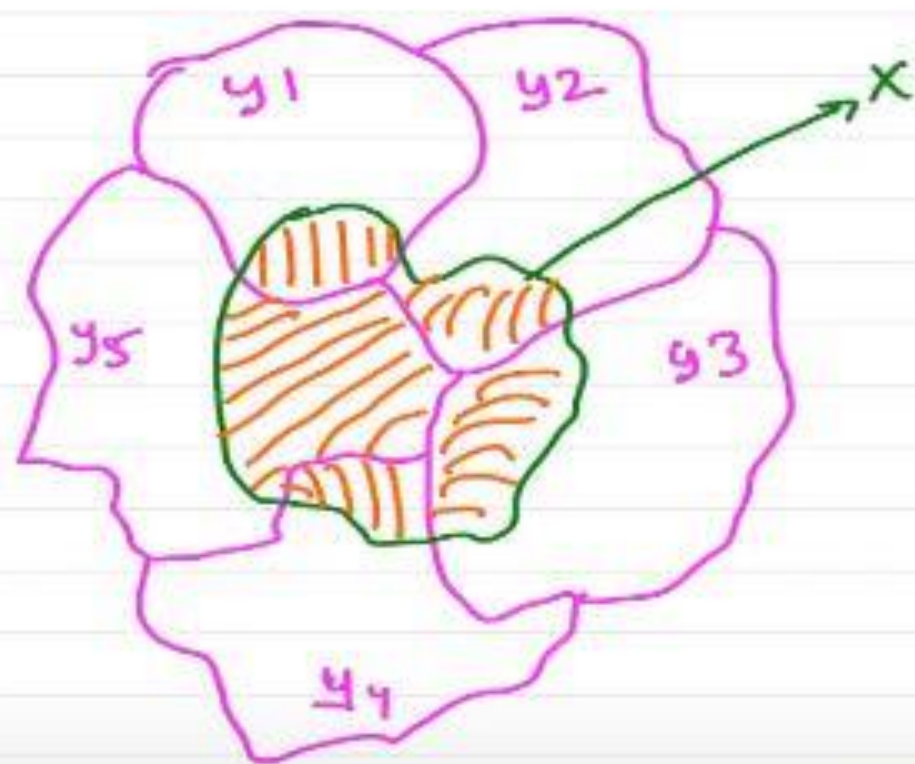
$$= \frac{P(x, y)}{P(x)} \rightarrow \text{joint distribution}$$

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

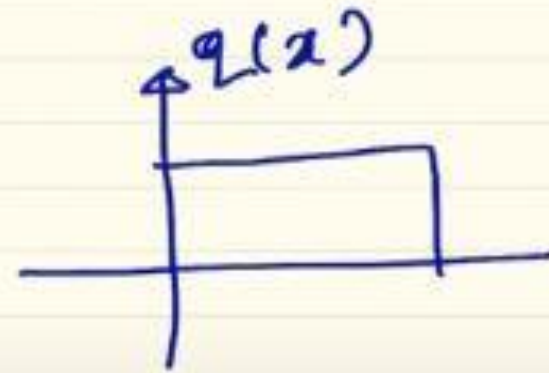
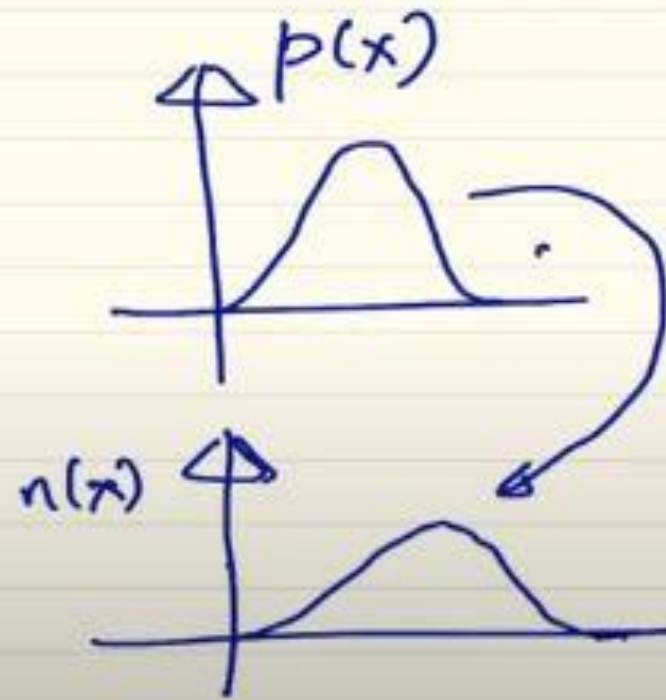


Kullback-Leibler divergence (K-L) is a measure of how one probability distribution is different from the second. For the discrete probability distribution  $P$  &  $Q$ , the K-L divergence between  $P$  &  $Q$  is defined as



# K-L Divergence

---

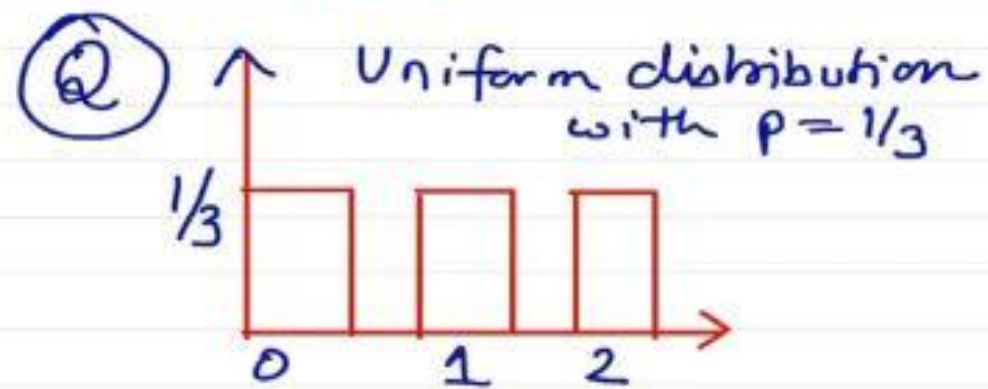
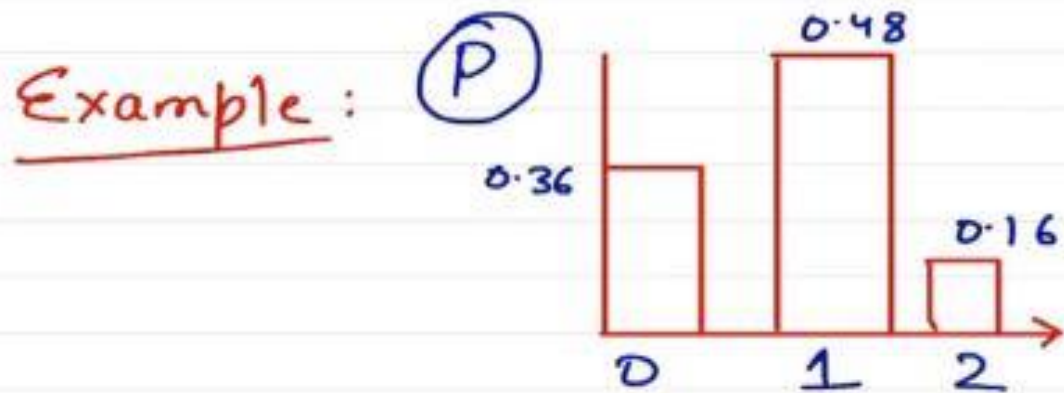


high value

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  $\mu_1$  &  $\mu_2$  are the means &  $\Sigma_1, \Sigma_2$  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

$$KL(P(x) \parallel Q(x)) = \sum_x p(x) \left\{ -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| \right. \\ \left. - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{k}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| \right. \\ \left. + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\}$$



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 → scalar

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

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

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

Trace & Expectation trick.

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

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

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

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

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

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

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

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

→ 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

Kumar

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

$A^T \Sigma_2^{-1} A$   
 $B^T \Sigma_2^{-1} B$   
 $\begin{cases} B^T \Sigma_2^{-1} A \\ A^T \Sigma_2^{-1} B \end{cases}$



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



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

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

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

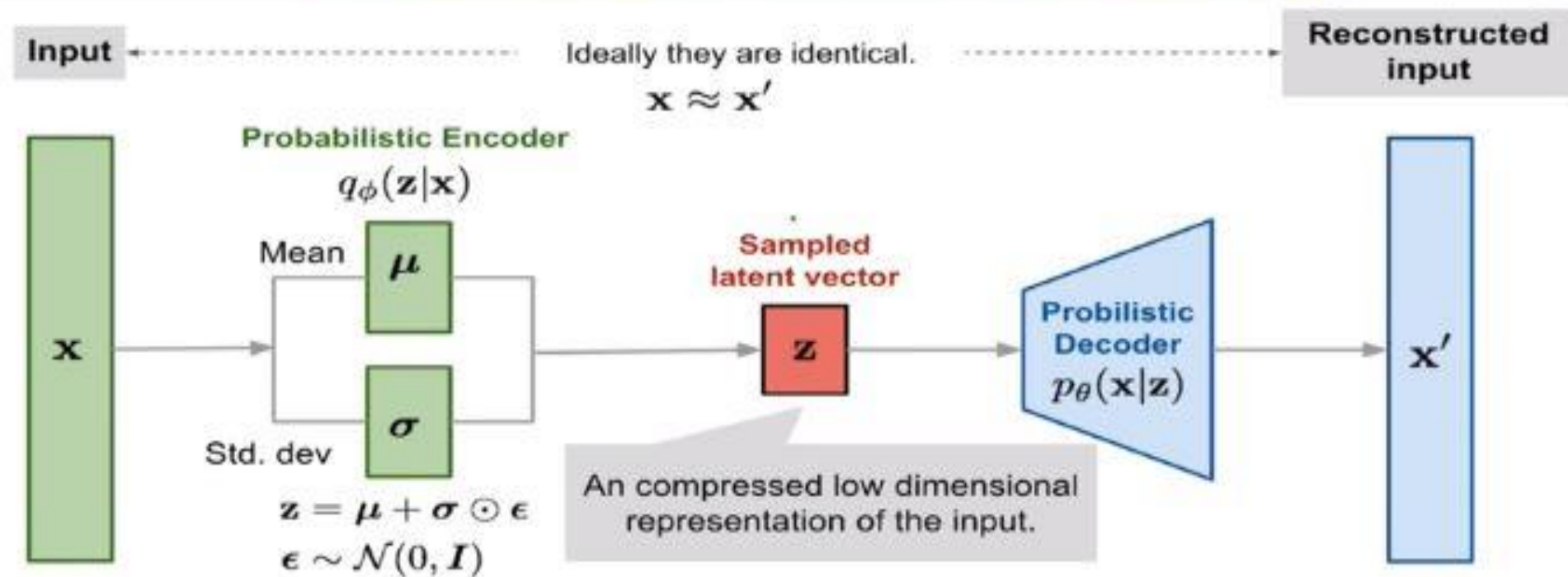
So substituting (2), (4) in (1) we obtain

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

proved

# Variational Autoencoder

#

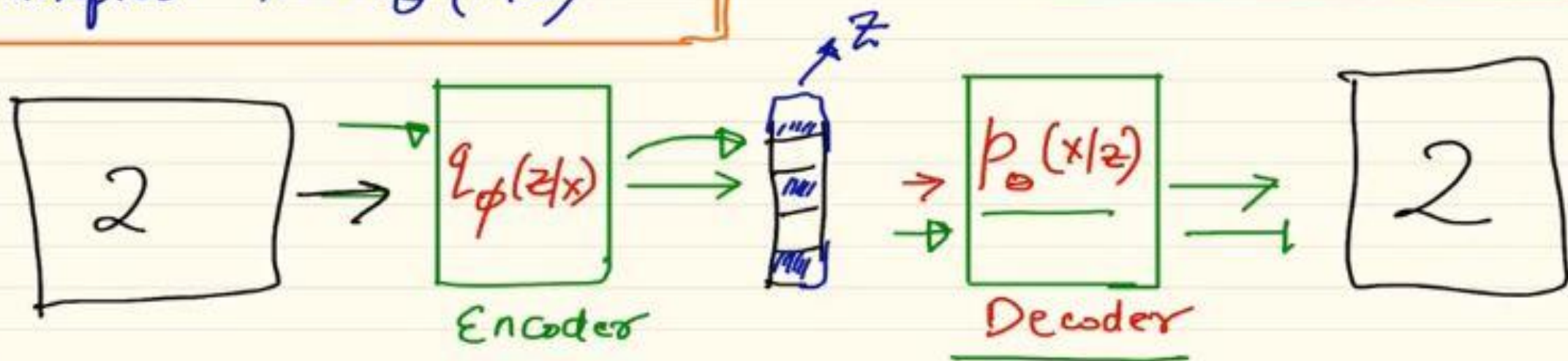


$$\text{Loss} = \mathcal{L}(\theta, \phi) = -\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] + D_{KL}(q_\phi(z|x) \| p_\theta(z))$$

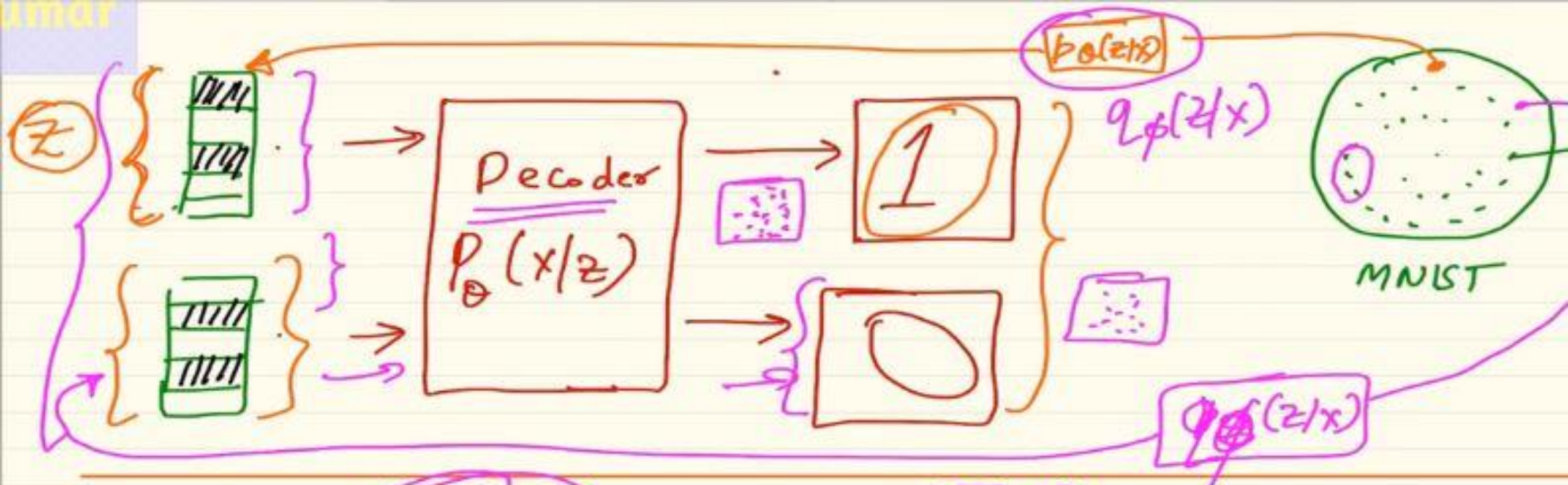
## The goal of VAE

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  $x' \sim p_{\theta}(x/z)$

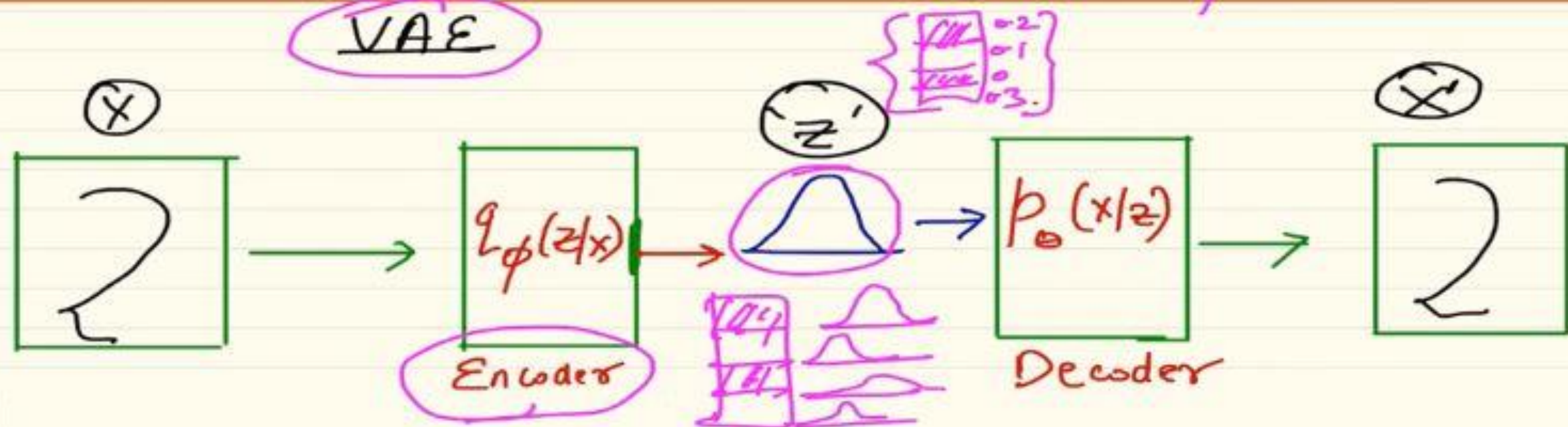
Typical Autoencoder





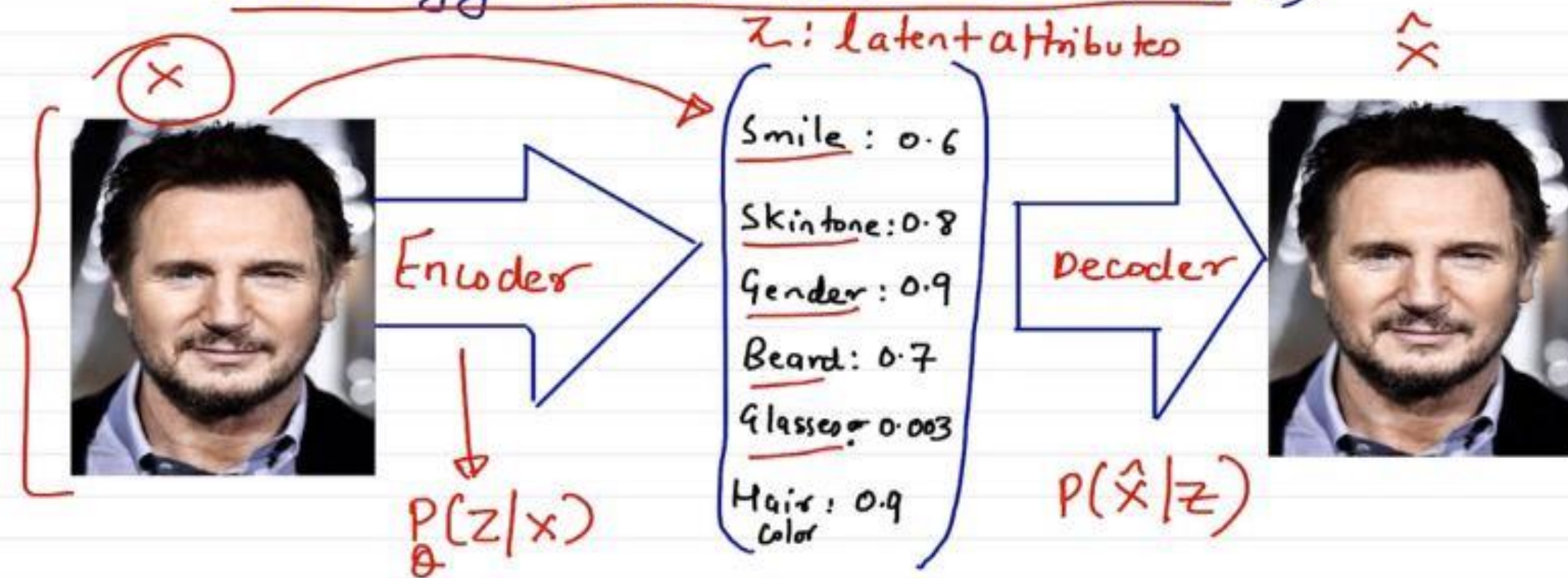


VAE



Latent variables can be placed in 2 categories

- ① <sup>Latent (2)</sup> Variables corresponding to a real feature of the object that have not been measured (may be technology is not available to do that)

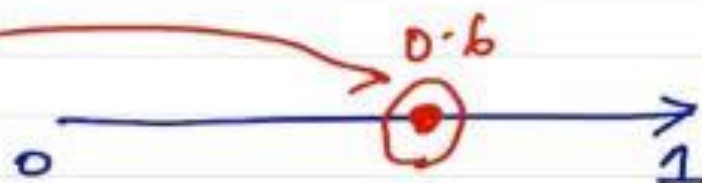




Using VE, we define latent attributes in probabilistic terms.

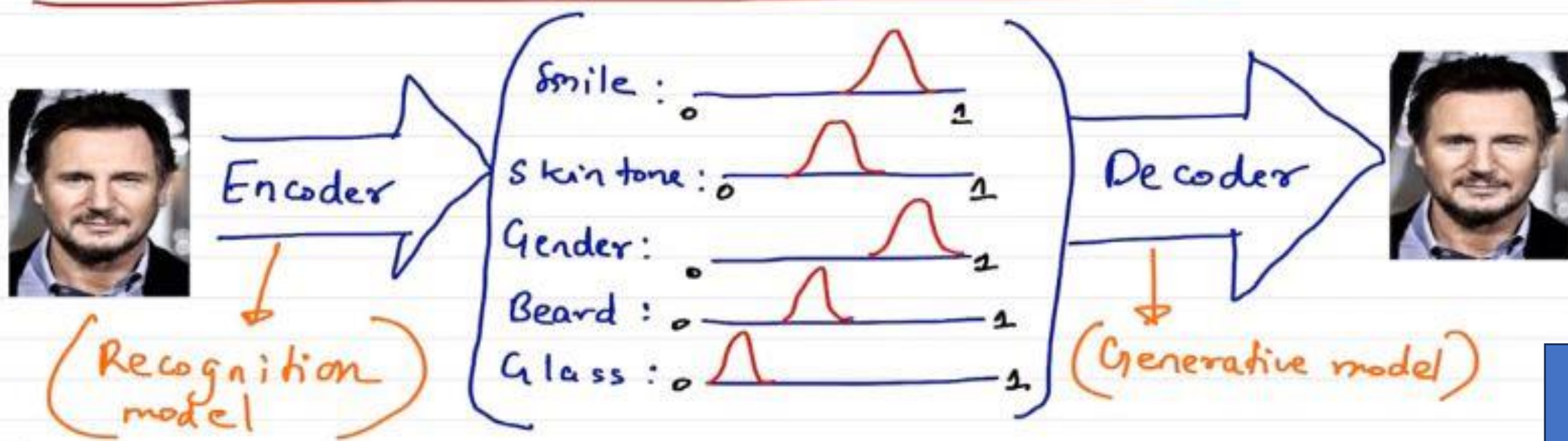
AE  
Smile (discrete)

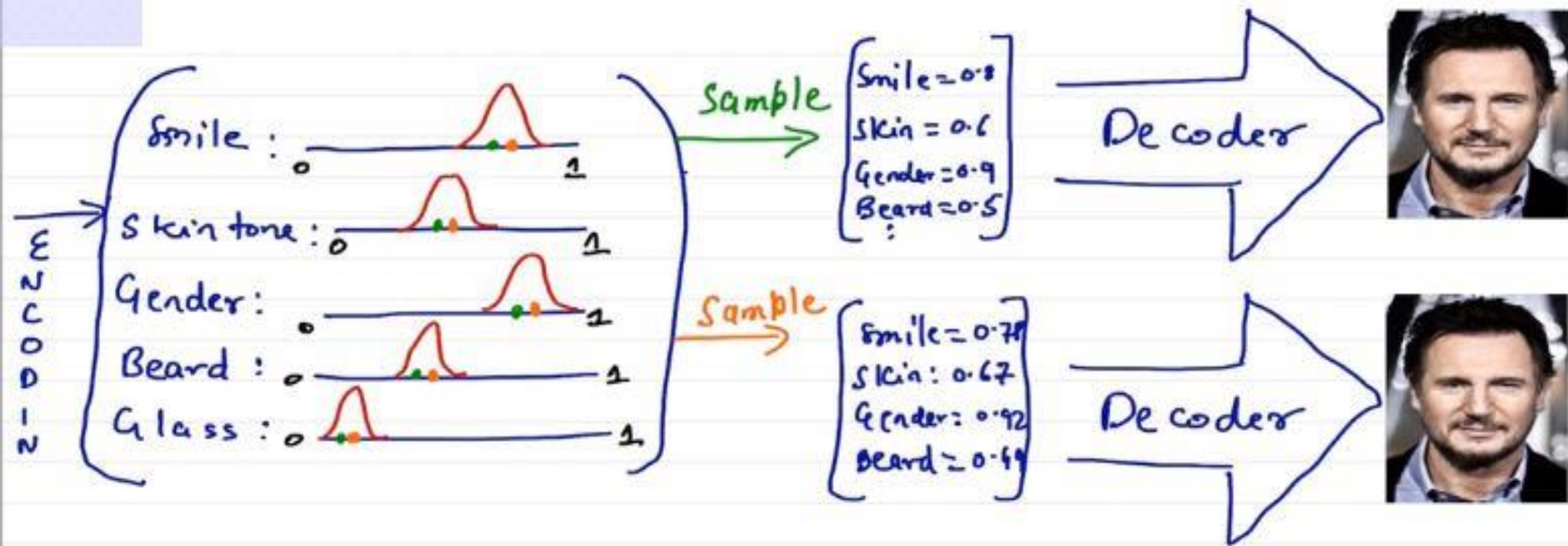
VE  
Smile (Probabilistic)





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 input 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 <sup>to</sup> similar reconstruct