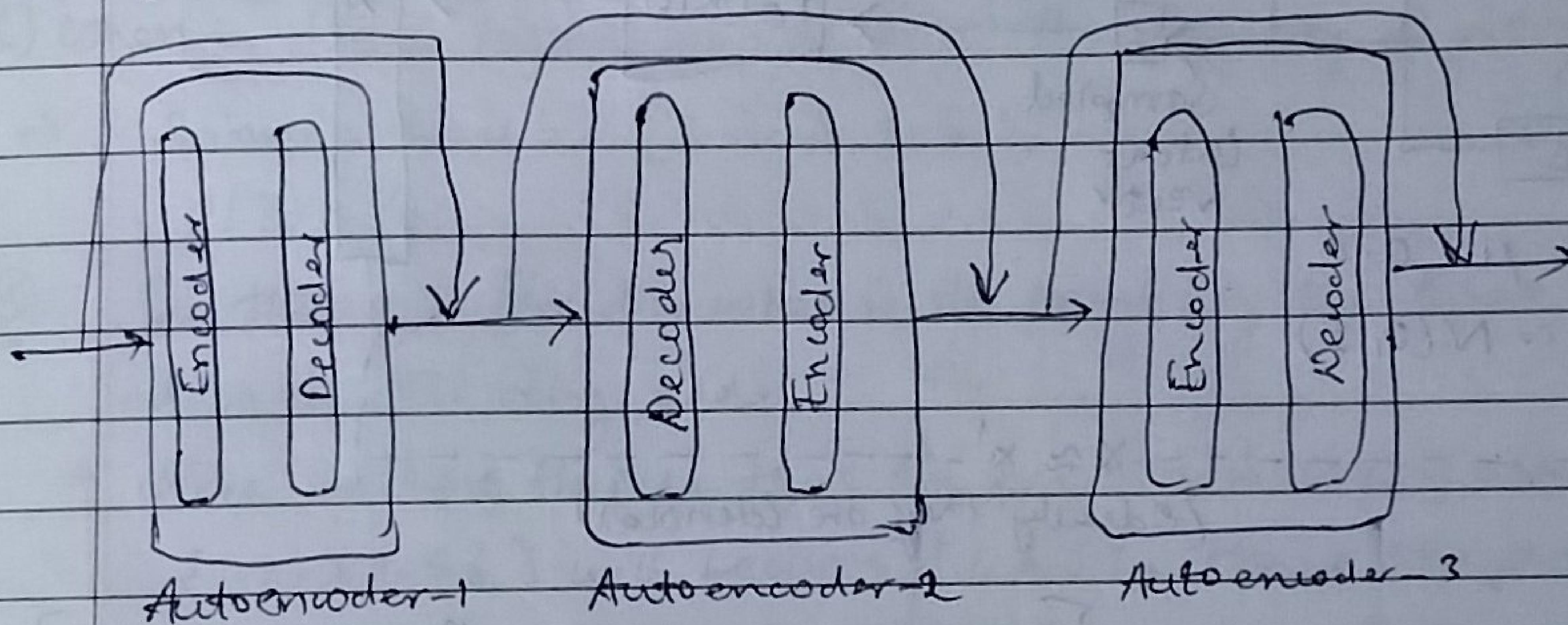


Variational Auto Encoders :-

We will first learn about the stacked autoencoder. As we know, that some datasets have a complex relationship within the features. Thus, ^{Page No.} during only one Autoencoder is not sufficient. A single autoencoder might be unable to reduce the dimensionality of the input features. Therefore, for such use cases, we use stacked autoencoders.

The stacked autoencoders are, as the name suggests, multiple encoders stacked on top of one another. A stacked ~~auto~~ autoencoder with three encoder stacked on top of each other is shown in the following figure



According to the architecture shown in the figure above, the input data is given to autoencoder 1. The output of the autoencoder 2 and the input of the autoencoder 1 is then given as an input to autoencoder 2. Similarly, the output of autoencoder 2 and the input of autoencoder 3 are given as input to autoencoder 3. Thus, the length of the input vector for autoencoder 2. This technique also helps to solve the problem of insufficient data to some extent.

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

Let's learn about Denoising Autoencoders. \rightarrow The purpose of DAE is to remove the noise.

DAE is primarily used for unsupervised learning. It is a variant of the standard autoencoder, but here is a critical twist : \rightarrow It is trained to reconstruct clean input data from corrupted (noisy) versions.

Simply, we deliberately corrupt the input data x to get a noisy version \tilde{x} , then train the network to reconstruct the original, clean data x from \tilde{x} .

Why Noise \rightarrow Preventing Overfitting, Encourage network to learn the more general features & improves the performance

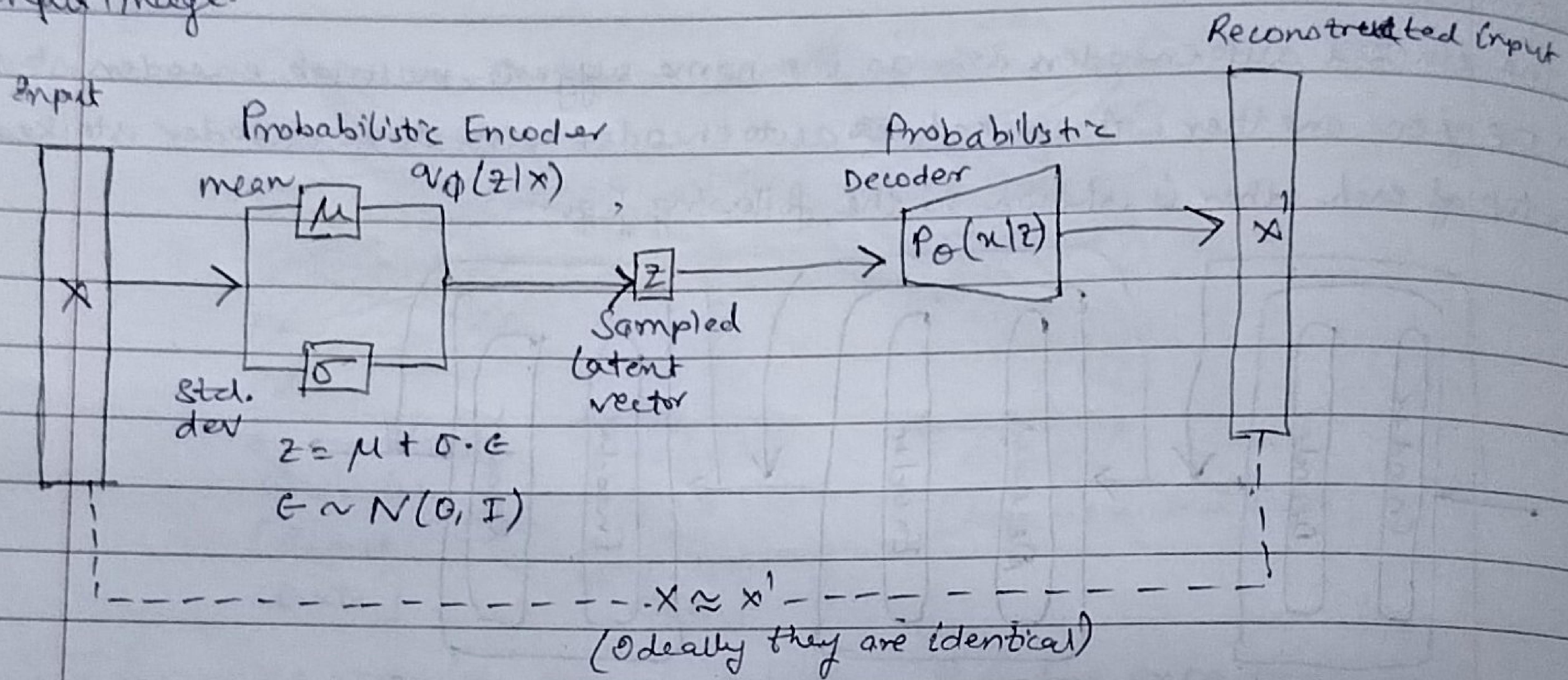
$$x^{(i)} \sim g_{\phi}(z^{(i)} | \theta^{(i)})$$

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

Page No.

Date

The bottleneck in the autoencoder contains the compressed version of the image that is input in a very optimised form that is the reduced dimension of the input image.



The only thing that is added on are \rightarrow [Probabilistic Encoder & Probabilistic Decoder]

Prerequisite:

- Probability $\rightarrow p(x) \rightsquigarrow$ probability
- $\rightarrow p(x|y) \rightsquigarrow$ conditional probability
- $\rightarrow E[\cdot] \rightsquigarrow$ expectation
- \rightarrow KL Divergence

- $p(x)$: defines the probability of random variable x $[p(x) \leftarrow x]$
- $p(x|y)$: defines as the probability of a random variable x provided y has happened, it also called conditional probability.

$$p(y|x) = \frac{p(x|y) p(y)}{p(x)} \rightsquigarrow \begin{array}{l} \text{Likelihood} \\ \text{Posterior Probability} \end{array} \quad \begin{array}{l} \text{prior probability} \\ \text{Bayes theorem} \end{array}$$

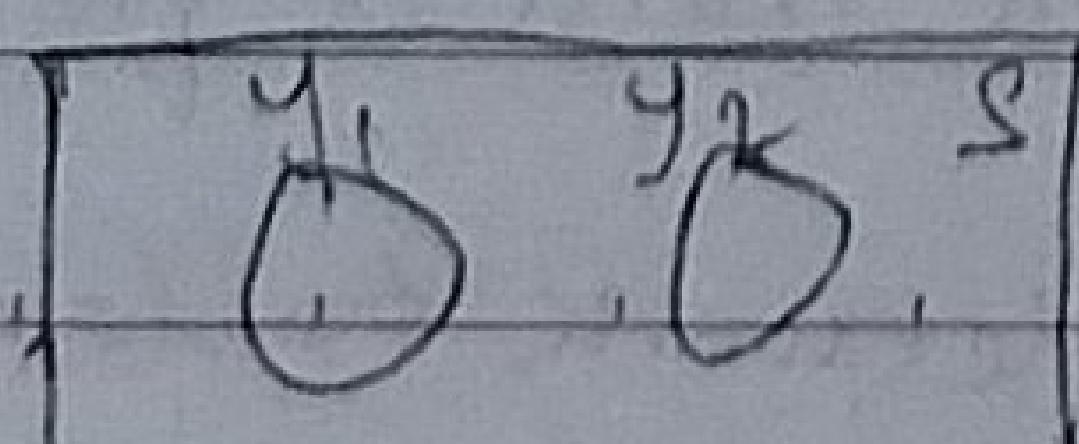
\downarrow

$\frac{p(x,y)}{p(y)} \rightsquigarrow \text{Joint distribution} \quad \text{--- (1)}$

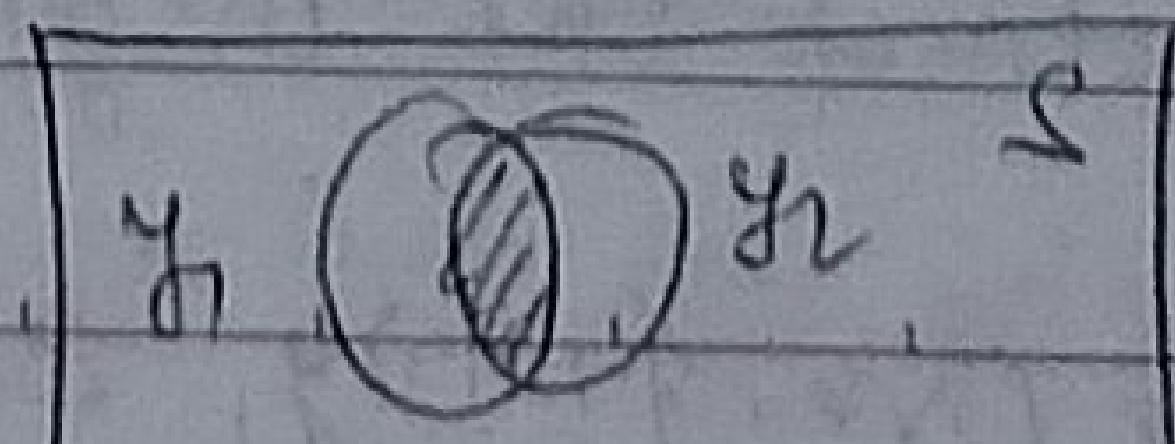
(ii) Theorem of total probability :-

Let $y_1, y_2, y_3, \dots, y_N$ 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

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



Mutual Exclusive



No mutual exclusive

from (2) in (1) $\left\{ P(Y|X) = \frac{P(X|Y) P(Y)}{\sum_{i=1}^n P(X|Y_i) P(Y_i)} \right\}$

Page No.

Date

(4) Expectation of a 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 :-

$$E(X) = \sum_{i=1}^K x_i P(x=x_i) \approx E_p(x)$$

(5) When a die is tossed once. what is the Probability of getting 3

→ Sample Space = {1, 2, 3, 4, 5, 6}
 $P(3) = 1/6$

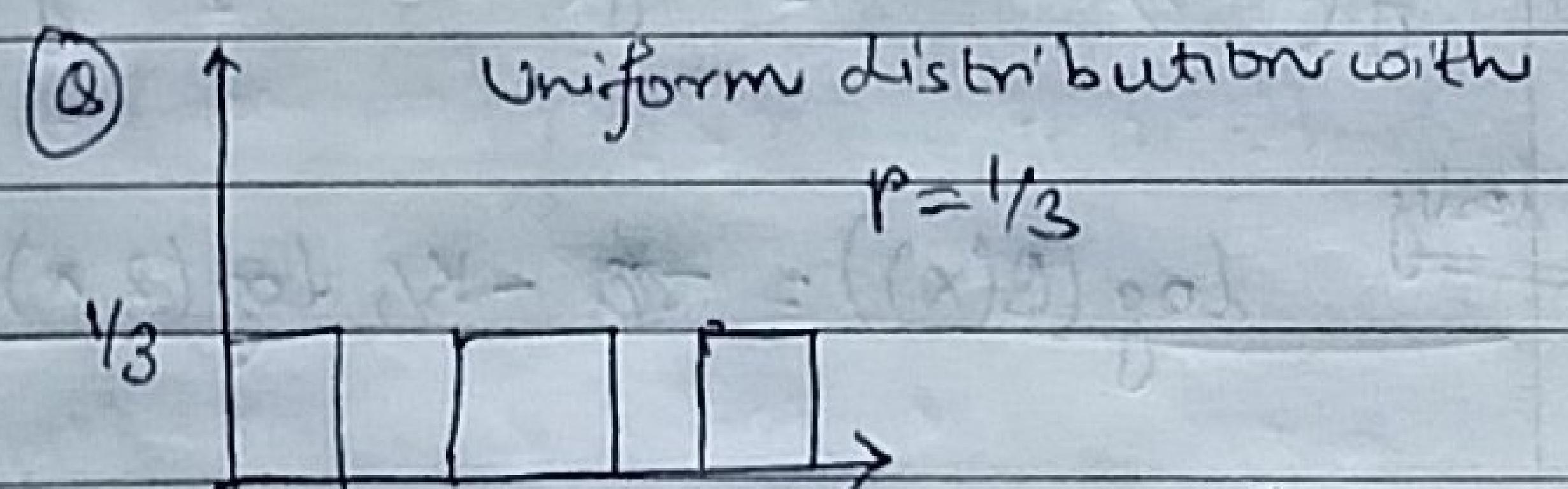
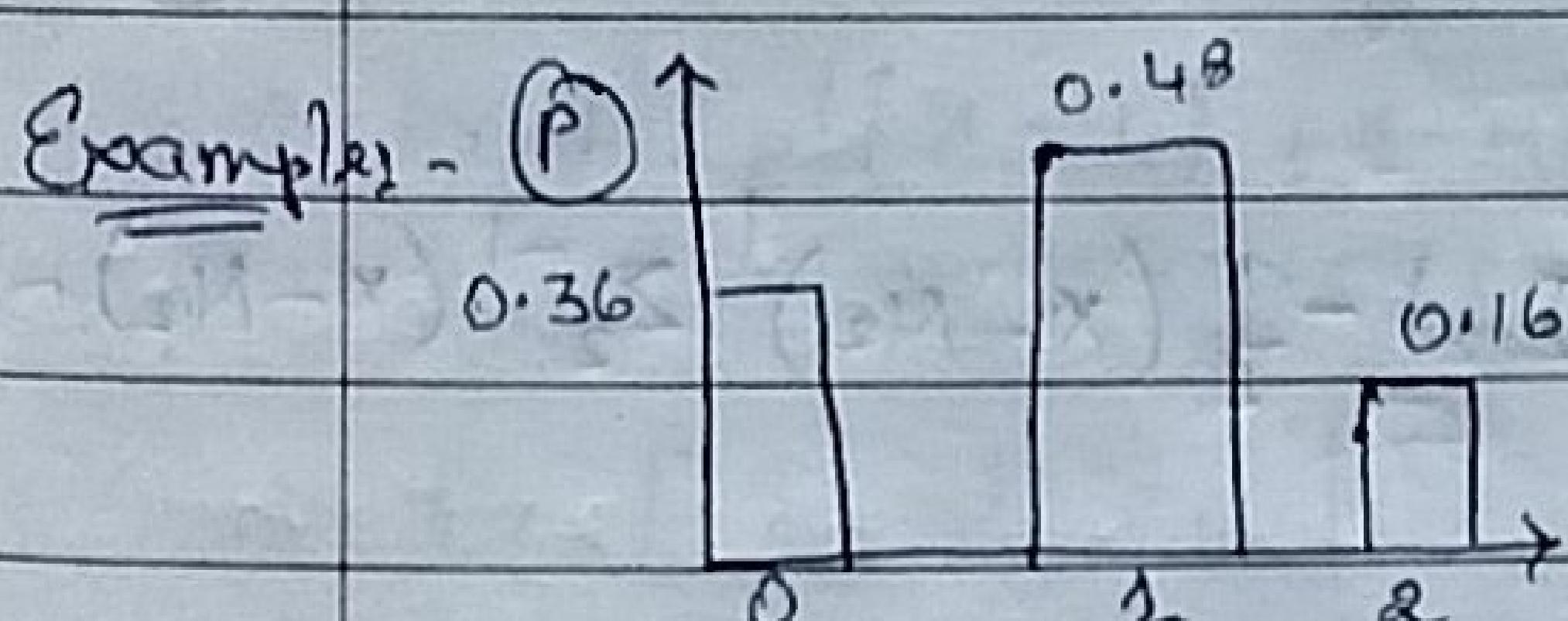
(6) In tossing a fair die, what is the probability the 3 has occurred conditionally on the toss being odd.

→ Since we have the idea that the question is saying about the odd no. {1, 3, 5} will become {1, 3, 5}. Hence the probability will be reduced to $1/3$.

(7) KL Divergence:- (Kullback - Leibler divergence) → It is measure of how one probability distribution is different from the second. For the discrete probability distribution P & Q, the KL divergence b/w P & Q is defined as

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

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



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

$$= 0.09637 \text{ nats} \rightarrow \text{unit}$$

Properties

- ① $KL(P||Q)$ 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 \rightarrow

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

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

$N \Rightarrow$ Normal distribution

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

Recall, the multivariate normal density is defined as \rightarrow

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

if the two distributions have the same dimension 'K'.

$$KL(P(x)||Q(x)) = \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1)$$

Derivation : \rightarrow We know, $KL(P(x)||Q(x)) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \quad \textcircled{1}$

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

\rightarrow Taking log both sides \rightarrow $\textcircled{2}$

$$\log(P(x)) = -K/2 \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \quad \textcircled{3}$$

Similarly

$$\log(Q(x)) = -K/2 \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \quad \textcircled{4}$$

We can write eqn $\textcircled{1}$ as \rightarrow

$$KL(P(x)||Q(x)) = \sum_x P(x) (\log(P(x)) - \log(Q(x))) \quad \textcircled{1}$$

Substituting $\textcircled{3}$ & $\textcircled{4}$ in $\textcircled{1}$ results in: \rightarrow

$$KL(P(x) || Q(x)) = \sum_n P(n) \left\{ -\frac{\kappa}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{\kappa}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\}$$

$$\text{or } KL(P(x) || Q(x)) = \sum_n P(n) \left\{ + \frac{1}{2} \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} \log (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right\} - \textcircled{5}$$

Now, let's consider part by part.

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

\nearrow we can write this by the Rule of the Expectation

FROM LINEAR ALGEBRA

$$\rightarrow E(x^T A x) = E(\text{tr}(x^T A x)) - \textcircled{6}$$

$$\sum E(\text{tr}(A x x^T)) - \textcircled{7}$$

$$\rightarrow \text{tr}(E(A x x^T)) - \textcircled{8}$$

Trace & Expectation trick

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

• trace of x is scalar

$$\rightarrow \text{tr}(AB) = \text{tr}(BA) - \textcircled{9}$$

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

$$= \text{tr}(CAB) - \textcircled{10}$$

$$\rightarrow \text{tr}(ABC) \neq \text{tr}(ACB) - \textcircled{11}$$

$$\rightarrow E(\text{tr}(x)) = \text{tr}(E(x)) - \textcircled{12}$$

\hookrightarrow trace

trace: \rightarrow sum of all diagonal elements in a matrix

Let's rewrite again,

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

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

$$E_p \left[\text{tr} \left(\frac{1}{2} (x - \mu_1)(x - \mu_1)^T \Sigma_1^{-1} \right) \right] - \textcircled{13}$$

$$\rightarrow \text{tr} \left[E_p \left(\frac{1}{2} (x - \mu_1)(x - \mu_1)^T \Sigma_1^{-1} \right) \right] - \textcircled{14}$$

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

\hookrightarrow covariance matrix

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

$$\Rightarrow \text{tr} [I_K] = K - \textcircled{15}$$

$$\therefore E_p \left[\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right] = K \rightarrow \textcircled{16}$$

Nao, consider the second part

$$\rightarrow \sum_x P(x) \left[\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] \xrightarrow{\substack{(A+B)^T \Sigma_2^{-1} (A+B) \\ (A^T + B^T) \Sigma_2^{-1} (A+B)}} =$$

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

(Adding μ_1 &
subtracting μ_1)

$$\rightarrow \sum_x P(x) \left[\frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \xrightarrow{\substack{B^T \Sigma_2^{-1} A \\ A^T \Sigma_2^{-1} B}}$$

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

Expanding we get

$$\rightarrow E_P \left\{ \frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) \right\} + E_P \left\{ (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right\} +$$

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

$$\rightarrow \underbrace{\text{tr} \left\{ \frac{\Sigma_2^{-1} \Sigma_1}{2} \right\}}_{\substack{\text{similar to} \\ \text{earlier derivation}}} + \underbrace{(\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}_{E[\text{constant}] = \text{constant}} + \underbrace{0}_{\substack{\text{prove} \\ \text{(next)}}}$$

(8)

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

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

$$\rightarrow (\mu_1 - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) = 0 \quad \underline{\text{proved}}$$

(*) $\boxed{E_P(x) \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = k} - (7)$

$$\boxed{E_P(x) \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) = \text{tr} \{ \Sigma_2^{-1} \Sigma_1 \} + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}$$

L (8)

from (7) & (8) in (5)

$$\boxed{KL(P(x) || B(x)) = \frac{1}{2} \log \left| \frac{\Sigma_1}{\Sigma_2} \right| + k + \text{tr} \{ \Sigma_2^{-1} \Sigma_1 \} + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)}$$

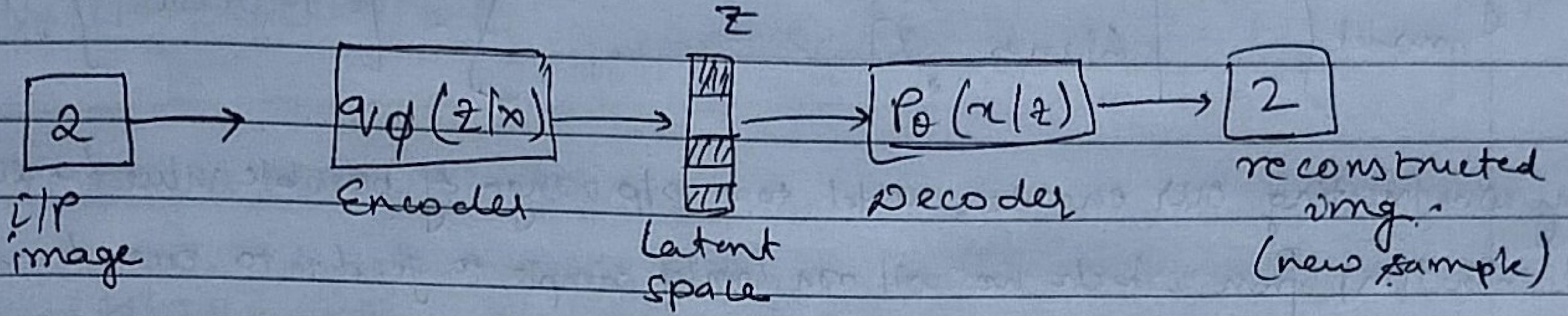
(proved)

The goal of VAE

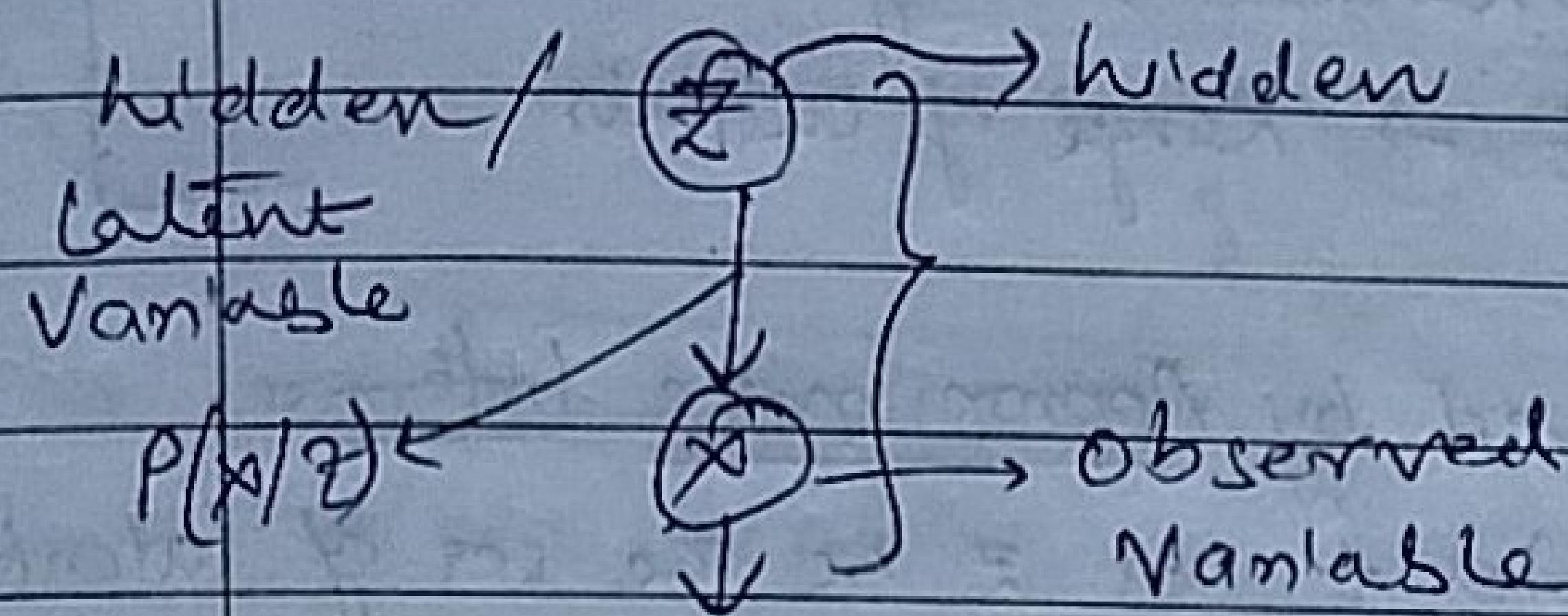
The goal of VAE is to find a distribution $q(z|x)$ of some latent variables which we can sample from $z \sim q(z|x)$ to generate new samples $x' \sim p_\theta(x|z)$

Page No.

Date

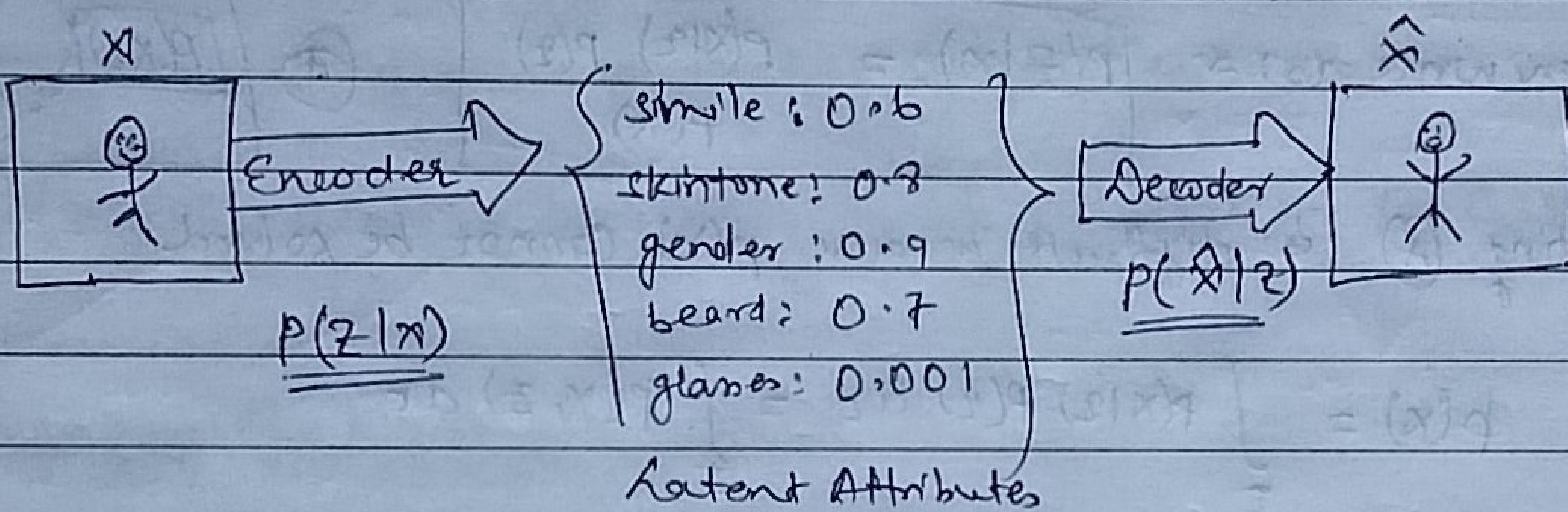


Latent Variable (z) :- We model ~~system~~ as a collection of random variables.



The edge drawn from $z \rightarrow x$ is the conditional distribution $p(x|z)$

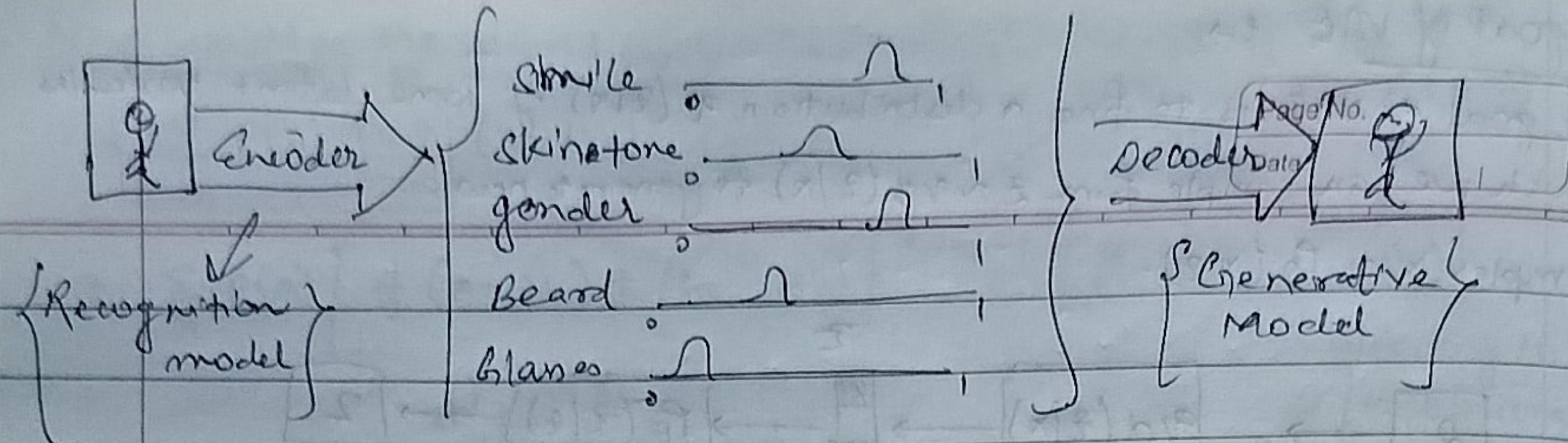
Latent Variables corresponding to a real feature of the object that have not been measured (may be technology is not available to do that)



→ In the above examples, we trained autoencoder on a large dataset of faces with encoding dimension of 6. An ideal autoencoder will learn the descriptive attributes of faces such as skin color, smile etc. In order to describe observation in some compressed form.

→ In above example, we have described the i/p image in terms of latent variables using single value to describe each attribute. For instance, what single value you will assign for photo of Mona Lisa?

→ With the approach of PD, we can represent each latent attribute for a given i/p as a probability distribution. When decoding, we will randomly sample from each latent state distribution to generate a vector as i/p for an decoder model.



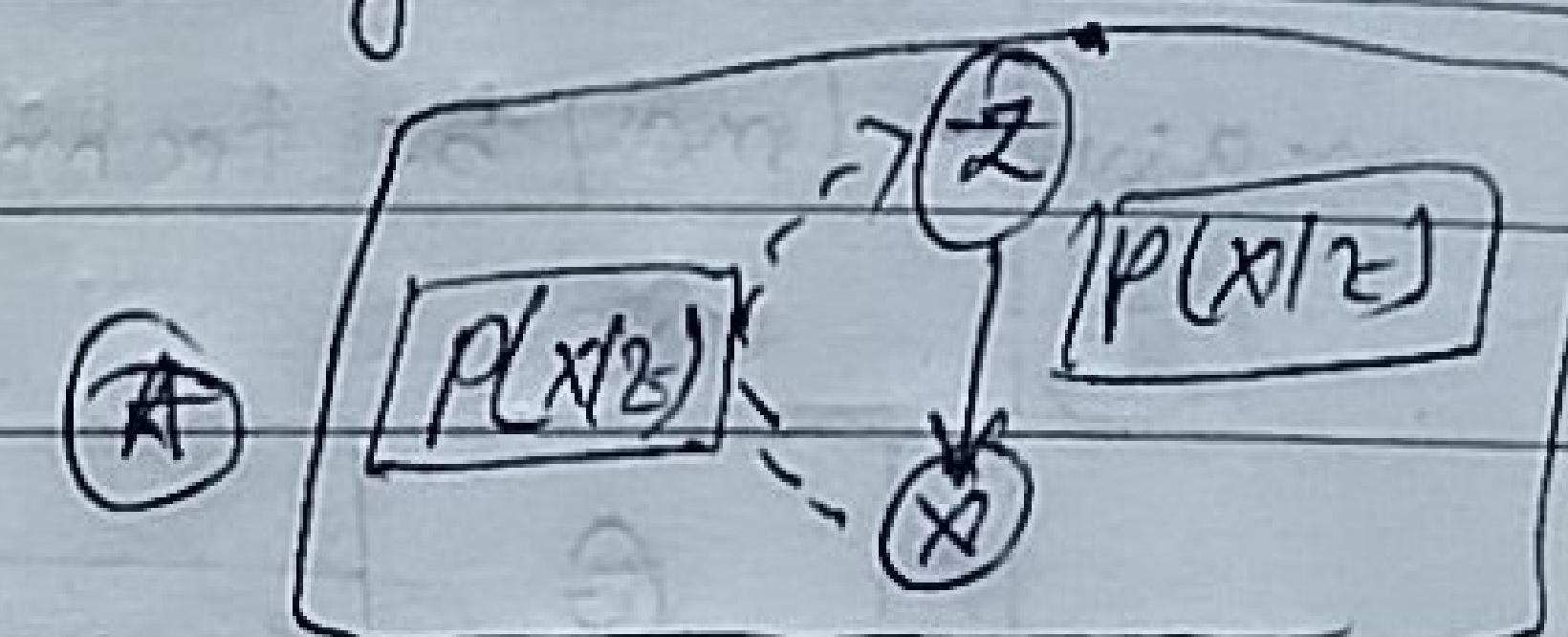
By constructing our encoder model to a o/p range of possible value (a statistically distribution) from which we will randomly sample to feed into our decoder model. The values which are nearby each other in latent space must corresponds to similar reconstruction.

Auto Encoder → Only one output
VAE → Range of outputs

- The derivation of loss function can be solved by 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 as:
$$p(z|x) = \frac{p(x|z) p(z)}{p(x)}$$



Evaluating \textcircled{A} is difficult because $p(x)$ cannot be solved

Reason:-

$$p(x) = \int_{\mathbb{Z}} p(x|z) p(z) dz = \int_{\mathbb{Z}} p(x, z) dz$$

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 . → we have 6 variable so we need 6 integrals

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 via variational inference. The main idea of VI is to pose the inference problem as an optimisation 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 $p(z|x)$ & $Q(z|x)$

$$D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) = \sum_z Q_\phi(z|x) \log \left(\frac{Q_\phi(z|x)}{P_\theta(z|x)} \right)$$

$$= E_{z \sim Q_\phi(z|x)} \left[\log \frac{Q_\phi(z|x)}{P_\theta(z|x)} \right]$$

Page No.

Date

$$\therefore D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) = E_{z \sim Q_\phi(z|x)} \left[\log(Q_\phi(z|x)) - \log(P_\theta(z|x)) \right] - B$$

Putting A to B (Substituting)

$$D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) = \sum_z Q_\phi(z|x) \log \left(\frac{Q_\phi(z|x)}{P_\theta(z|x)} \right)$$

$$= E_{z \sim Q_\phi(z|x)} \log \left(\frac{Q_\phi(z|x)}{P_\theta(z|x)} \right)$$

$$D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) = E_z \left[\log Q_\phi(z|x) - \log \frac{P_\theta(x|z) P_\theta(z)}{P_\theta(x)} \right]$$

$$= E_z \left[\log(Q_\phi(z|x)) - \log P_\theta(x|z) - \log P_\theta(z) + \log P_\theta(x) \right]$$

Since the expectation is over z & $P(x)$ doesn't involve z , it can be moved out

$$D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) - \log P_\theta(x) = E_z \left[\log(Q_\phi(z|x)) - \log P_\theta(x|z) - \log P_\theta(z) \right]$$

Rearranging the Equations:-

$$\log P_\theta(x) - D_{KL}(Q_\phi(z|x) \parallel P_\theta(z|x)) =$$

$$E_z [\log(P_\theta(x|z))] - E_z [\log(Q_\phi(z|x)) - \log P_\theta(z)]$$

$$= E_z [\log(P_\theta(x|z))] - D_{KL}[Q_\phi(z|x) \parallel P_\theta(z|x)]$$

This is the VAE objective function, where the first term represents the reconstruction likelihood &

the second term ensures that our learned distribution

Q is similar to the prior distributions P .

loss function = - Objective function

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

Page No.

Proved!!!

$$\text{on } \log P_\theta(x) - [D_{KL} Q_\phi(z|x) || P_\theta(z)] = -L(\theta, \phi)$$

Now, our target is to find optimal θ, ϕ such that $\theta^*, \phi^* = \underset{\theta, \phi}{\operatorname{arg\min}} (L(\theta, \phi))$

Understanding each Variable

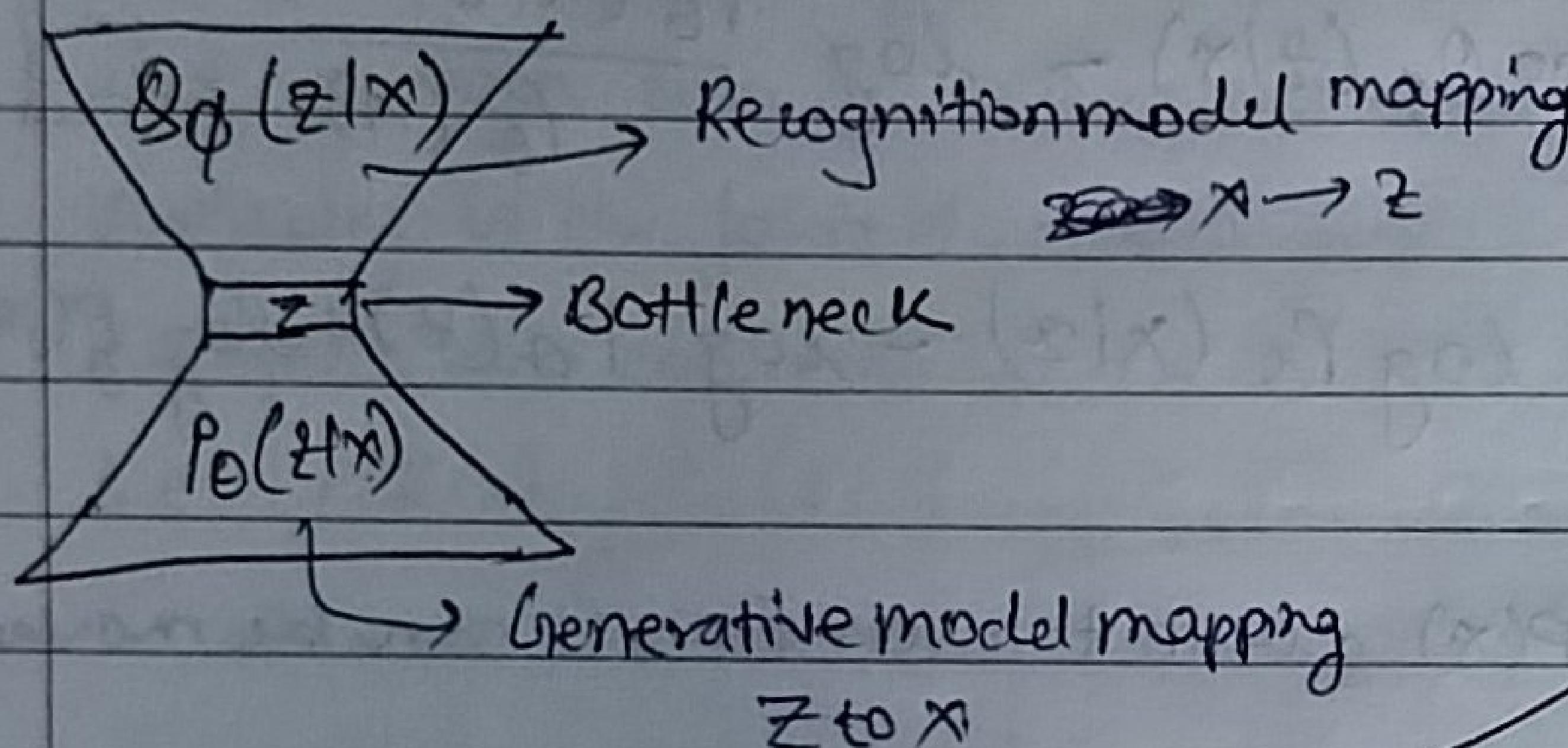
$$L(\theta, \phi) = E_{z} [\log P_\theta(x)] - [D_{KL} Q_\phi(z|x) || P_\theta(z)]$$

Inputs

\uparrow

\log
likelihood

Regulator



$$P_\theta(x|z) = N(\mu_\theta(z), \Sigma_\theta(z))$$

So, when we take log of Gaussian we get a square error b/w the data sample x & mean of the Gaussian distribution

Role

$$Q_\phi(z_1|x) \quad Q_\phi(z_2|x)$$

Latent Variable Space

KL divergence not allows pdf of latent variables not collapse w/ 0 variance but penalize if deviates from $N(0, 1) = P_\theta(z)$

$$P_\theta(x|z) = \frac{1}{\sqrt{(2\pi)^K |\Sigma_\theta(z)|}}$$

$$\exp \left[\frac{(x - \mu_\theta(z))^T \Sigma_\theta^{-1}(z) (x - \mu_\theta(z))}{2} \right]$$

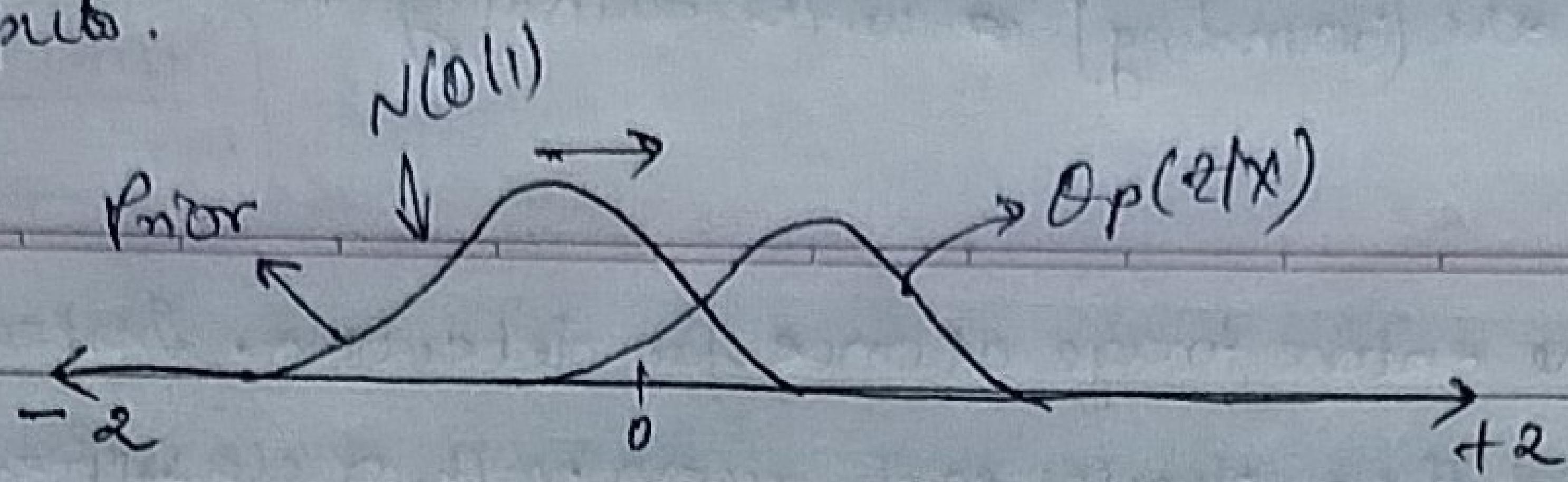
$$\text{on } \log P_\theta(x|z) \propto \left[(\hat{x} - \mu_\theta(z))^T \Sigma_\theta^{-1}(z) (x - \mu_\theta(z)) \right]$$

squared Reconstruction error

→ data fidelity in inverse problem

$$L = \text{data fidelity term} + \text{KL div term}$$

① Penalizing reconstruction loss (without data fidelity) :- encourages distribution to describe input.

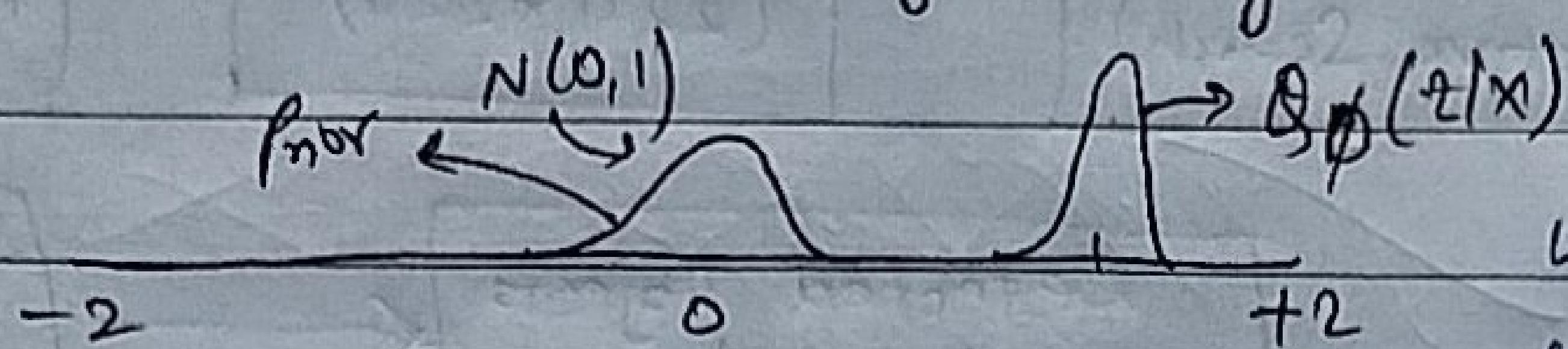


Page No. _____

Date _____

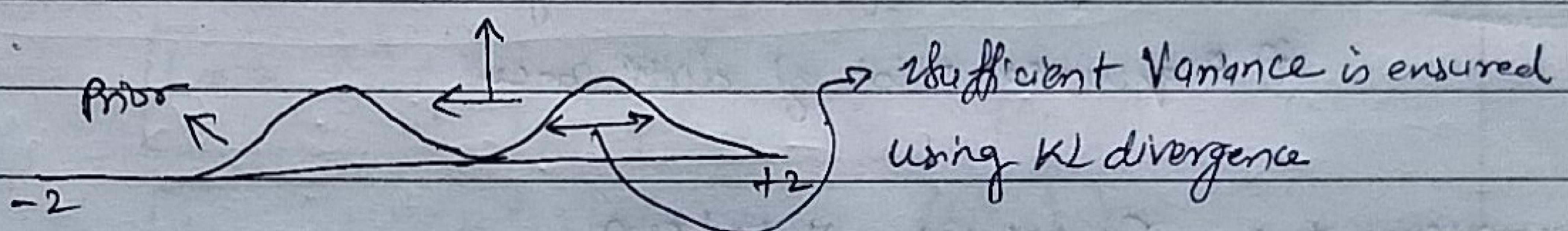
Our distribution deviates from the prior to describe some characteristics of data

② Without regularization networks cheat by learning narrow distributions,



with the small variance, this distribution is eff rep single value.

③ Attraction b/w the two distribution is due to 1st term



- ORIGINAL $\rightarrow N(0,1)$
- OUR DISTRIBUTION $\rightarrow N(\mu_\phi(x), \Sigma_\phi(x))$

Final loss (No derivation) :-

$$L(\theta, \phi) = -E_{\mathbb{Z}} \left[\log P_\theta(x|z) \right] + \frac{1}{2} E_K \left[\exp \left(\sum_p (\cdot) \right) + \mu_\phi^2(x) - 1 - \Sigma_\phi(x) \right]$$

$\downarrow z \sim P_\phi(z|x)$

Optimised loss :-

$$\hat{\phi}_i = -\frac{1}{S} \sum_{s=1}^S V_\phi \left[\log P_\theta(x|z^{(s)}) \right]$$

$$\text{where, } z^{(s)} = \mu_\phi(x) + E \odot \sigma_\phi(x)$$

$$E \odot \approx N(0,1)$$