

probability of z given x

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

Our YouTube Channel :-

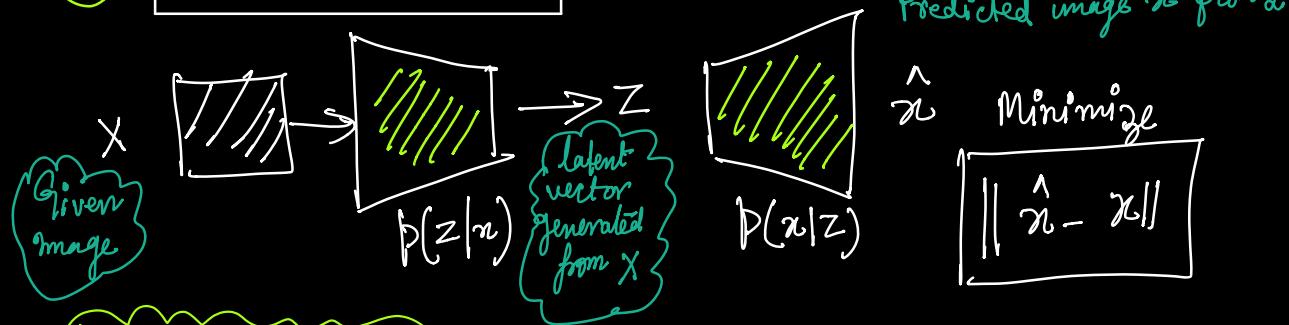
①

Deep learning for all at Manas
lab at IIT Mandi.

Aditya Nigam - I.I.T Mandi.

→ Can be converted into ② [Bottleneck]

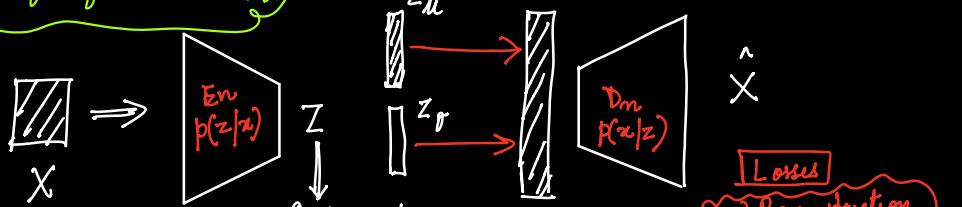
(A) Basic Autoencoder



Here our main question is ① How to make sampling backpropagable using reparameterization technique

⑥ Why we want to train encoder so as to produce \hat{x} 's that follow a specific fixed distribution

(B) Variational Autoencoders



Instead of the ② we want our encoder to get μ, σ

z_u, z_σ

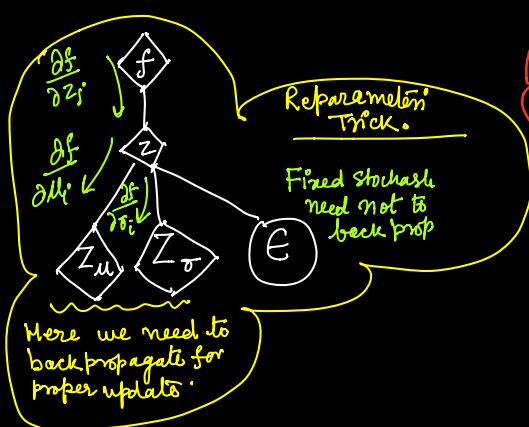
Sampling layer

sample a point from $G(z_u, z_\sigma)$

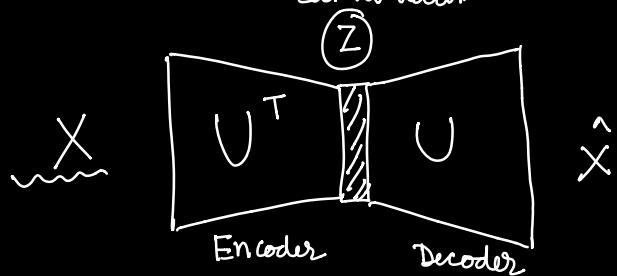
$(n \times 1)$

$(n \times n) \Rightarrow$ take only diagonal (variance)

$$\begin{aligned} & \text{KL-Div}(G(z_u, z_\sigma), \mathcal{N}(0, 1)) \\ &= \frac{1}{2} \sum_{i=1}^n (\mu_i^2 + \sigma_i^2 - \log(1e-8 + \sigma_i^2)) - 1 \end{aligned}$$



In basic autoencoder data goes into the bottleneck and reconstructed (2)



$$\text{Loss} = \min \|x - \hat{x}\|$$

(reconstruction error.)

If there is no non-linearity (i.e. w/o any activation fn) and there is only one hidden layer then this is very similar to PCA analysis.

This does not ensure that such A.E and PCA both learn the identical basis but map the similar space

Encoder (E)

$$\textcircled{1} \quad \underbrace{Z}_{p \times 1} = \underbrace{U^T}_{p \times d} \underbrace{X}_{d \times 1}$$

$$X \in \mathbb{R}^d \quad (\text{d-dim vector})$$

$$Z \in \mathbb{R}^p \quad (p\text{-dim vector})$$

Encoder is learning some transformation that can convert $\underbrace{X}_{\text{i/p}}$ to $\underbrace{Z}_{\text{latent vector}}$

Decoder (D)

$$\textcircled{2} \quad \underbrace{\hat{X}}_{d \times 1} = \underbrace{U}_{d \times p} \underbrace{Z}_{p \times 1}$$

$$\text{applying } \textcircled{1} \quad \hat{X} = U U^T X$$

Hence our loss fn has to be $\min \|X - \hat{X}\|$

$$\therefore \min \|X - UU^T X\|$$

main difference b/w PCA and A.E can be that in PCA

$$UU^T = I \quad [U \text{ is orthonormal by construction}]$$

In A.E $[U]$ may not be learned as orthonormal.

One can train deep autoencoders with non-linearity in order to learn better representation.

Important Concepts

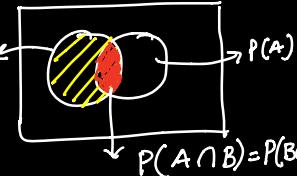
(3)

(A) Bayes Theorem:

① Conditional probability $\hat{=}$ events are $A \cap B$

$$\therefore P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$



$$\Rightarrow P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$\Rightarrow P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Assuming 1 student in the class of 20 has flu.
Event A: Student is A B C $\Rightarrow [P(A) = \frac{1}{20}]$

Evidence B: 5 Girls & 15 Boys

Now given this new evidence what is the probability that A B C has flu.

$$P(A|B) = \frac{1}{5} \text{ (goes up) if girl } \Rightarrow \textcircled{O} \text{ (if student is a Boy)}$$

New evidences are going to influence the Hypothesis

(B) Information (I): How one can estimate the amount of information in a sentence/expression
 \Rightarrow (An event) (X)

Here we can have 3 things :

<u>X Event</u>	<u>P(X) Probability</u>	<u>-log(P(X)) Information</u>
① Virat scored a century.	\uparrow (highly probable event)	\downarrow (less information)
② Kenya wins Cricket World Cup	\downarrow (rare event)	\uparrow (high information)
③ Tomorrow it rain or don't	1 (Certain event)	① [no information]

{So basically rare events carry more information}

(C) Average of Information is Entropy (H): (4)

The expected value of information w.r.t any event \hat{x}
averaged over all values \hat{x} can attain is Entropy H .

$$H = - \sum p(x) \log p(x)$$

This is the expected value of $\log p(x)$ w.r.t $p(x)$.
 Summation over all x 's
 Probability of that \hat{x} to happen
 Information Content in any \hat{x}

(D) KL-Divergence (KL-Div): In order to compute the similarity between two distributions say P and Q
KL-Div $(P \parallel Q)$ can be used defined as the KL-Div of Q distribution w.r.t P .

(i) Entropy of Q - Entropy of P

Amount of information in Q distribution

Amount of information in P distribution

$$-\sum q(x) \log q(x) + \sum p(x) \log p(x)$$

This expectation wrt $q(x)$ This expectation is wrt $p(x)$

KL-Div is almost this except that the expectation is always computed w.r.t $p(x)$ as KL-Div is w.r.t $p(x)$

$$(ii) \quad -\sum p(x) \log q(x) + \sum p(x) \log p(x) \quad (5)$$



This is the cross entropy between (p) and (q_r) distributions.



This is (ave) entropy of (p) distribution

Now both expectations are wrt $p(x)$

Hence,

$KL\text{-Div}(p(x) \mid q_r(x))$ can be formally defined as the difference between average information of $q_r(x)$ wrt $p(x)$ and that of $p(x)$ wrt $p(x)$.

$$\begin{aligned} KL\text{-Div}(p(x) \mid q_r(x)) &= -\sum p(x) \log q_r(x) + \sum p(x) \log p(x) \\ &= \sum p(x) \log \frac{p(x)}{q_r(x)} \\ &= -\sum p(x) \log \frac{q_r(x)}{p(x)} \end{aligned}$$

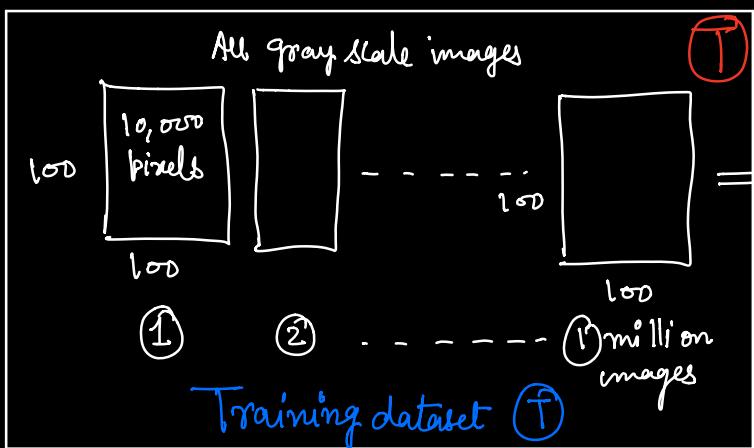
⊗ $KL\text{-Div}$ is not Symmetric as $KL\text{-Div}(p \mid q) \neq KL\text{-Div}(q \mid p)$

⊗ $KL\text{-Div} \geq 0$ it is always +ve

↳ Hence it is a distance measure b/w a divergence.

$$\therefore KL\text{-Div}(q_r(z) \mid\mid p(z|x)) = -\sum q_r(z) \log \frac{p(z|x)}{q_r(z)}$$

(we will come back to this.)



Let us assume that we have a very complex and huge training dataset T of 1 million images of size $100 \times 100 = 10,000$ pixels per image. (6)

But any 100×100 image is not just any random 2D matrix and all gray values are not equiprobable at each location.

In an image pixel's have dependence, specific gray level co-occurrence patterns.

Each image can be seen as 10K pixels each being sampled independently from $\{0-255\}$. Hence there can be $(256)^{10,000}$ possible such images. Very big.

We are assuming this pixel by pixel image sampling (i.e. gray value selection) experiment as a stochastic process and can be modeled using random variables.

Collection of random variables where each of them uniquely associated with an element in the set

$$\therefore P(X) = P(X_1, X_2, X_3, \dots, X_{10,000})$$

Joint probability distribution

Depending upon our training dataset T , we wanted to estimate $P_\theta(X)$ where $P_\theta(X) = \text{Probability of } (X \in T)$, with the distribution parameterized over $[\theta]$.

$$\theta^* = \underset{\theta}{\operatorname{arg\,max}} [P(X \in T)]$$

minimizing the error.

But why we are interested to compute $P_{\theta}(x)$

(7)

→ (a) Classification: Helps us to discriminate b/w images that are coming from \mathcal{T} or rest.

→ (b) Generative modeling: It can help us to sample new/unseen x_i^s from $P_{\theta}(x)$ distribution that are not even present. Such as non-trivial views, poses, interpolations b/w 2 views/poses.

For this image sampling experiment $p(x) = p(x_1, x_2, \dots, x_{10,000})$ is multivariate probability distribution. If we estimate it we know how to sample a new image (basically 10,000 values) from this joint distribution.

But such probability distribution estimation is intractable and very complex.

$$P(x) = P(x_1, x_2, \dots, x_{10,000}) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots \cdot p(x_{10,000}|x_1, x_2, \dots, x_{9,999})$$

M1

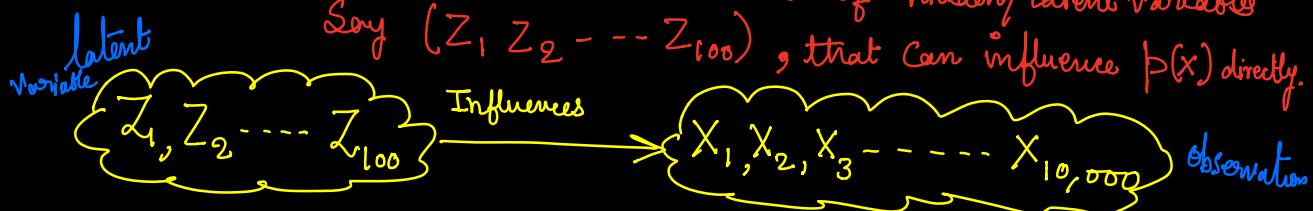
→ Computation is infeasible.

→ Since (x) is an image these x_i^s are not independent

→ There is huge amount of dependency between random variables, (x_i^s)

We can assume another set of hidden/latent variables

Say $(z_1, z_2, \dots, z_{100})$, that can influence $p(x)$ directly.



Now our observation (x) got dependent upon latent variables (z)

Basically our image (X), got influenced by few factors (Z), such as pose, illumination, noise...⁽⁸⁾

Since dimensions of (Z) is far lesser than (X) it is easier to get hold of $P_{\theta}(x)$ via (Z)

for example our line dataset, length, color, angle, width...

external parameters

\circledast We have observed some (X) that is influenced by (Z)

M_2 ∴ Distribution of observed data (X) is significantly influenced by (Z)

$$p(x) = \int_Z p(x, z) dz = \int_Z p(x|z) \cdot p(z) dz$$

Hence this marginalized estimation is also not feasible

Since this marginalization has to happen over all Z 's (It's can also be huge) may be dimensionally lesser than X [This is also intractable]

Now both M_1 & M_2 are not feasible, let us try something else that is convenient

Bayes Theorem:

$$P(H|E) = \frac{P(E|H) \times P(H)}{P(E)}$$

Annotations for Bayes Theorem:

- Hypothesis whose probability may be affected by evidence (E) i.e. data
- Probability of hypothesis given the observed evidence Posterior
- Probability of observing (E) given (H) Likelihood
- Prior probability distribution w/o evidence / data
- Marginal likelihood (Remains to be the same for all possible Hypothesis, if data is fixed)

9

$$x = x_1 - x_N \quad \text{Observed variable}$$

$$z = z_1 - z_M \quad \text{Set of latent variables}$$

with $N \gg M$

and

$p(z, x)$ is joint distribution

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

Given any x , we need to estimate a suitable \hat{q}_z that can be used as a bottleneck of the observed.

This makes the real estimation again intractable (exponential in time) w/o any close form.

This is again marginalized over our latent variable

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

Instead of statistically estimating $p(z|x)$ we wanted to approximate $p(z|x)$ via. $q_z(z) \in \mathcal{Q}$. \Rightarrow family of some known tractable distribution (Say unit normal, Gaussians...)

$$\therefore q_z^*(z) = \arg \min_{q_z(z) \in \mathcal{Q}} KL\text{-Div}(q_z(z) \mid p(z|x))$$

$p(z|x)$

We wanted to approximate this mapping function

Via.

$q_z(z)$

This is a distribution chosen from the family of tractable distributions (easy to estimate)

(*) It is chosen to be tractable as well as close to $p(z|x)$

(*) we need to estimate the best set of parameters of (*) that make it as close as above condition.

$$\therefore KL\text{-Div}(q_z(z) \mid p(z|x)) = - \sum_z q_z(z) \log \frac{p(z|x)}{q_z(z)}$$

(*) We wanted to maximize $P(X)$ that seems intractable. (15)
 Instead we started to approximate $P(z|x)$ via $q(z) \in \mathcal{L}$

$$\begin{aligned}
 \text{KL-Div}(q(z) \mid\mid p(z|x)) &= - \sum_z q(z) \log \frac{p(z|x)}{q(z)} \quad [\text{WHY}] \text{ of tractable} \\
 &= - \sum_z q(z) \log \frac{p(x,z)}{p(x) \cdot q(z)} \quad \text{Remember estimation has been carried out for a given } x \text{ here!} \\
 &= - \sum_z q(z) \left[\log \frac{p(x,z)}{q(z)} - \log p(x) \right] \quad x \text{ is fixed.} \\
 &= - \sum_z q(z) \log \frac{p(x,z)}{q(z)} + \log p(x) \quad \sum_z q(z) \rightarrow 1 \\
 &= - \sum_z q(z) \log \frac{p(x,z)}{q(z)} + \log p(x)
 \end{aligned}$$

$$\log p(x) = \text{KL-Div}(q(z) \mid\mid p(z|x)) + \sum_z q(z) \log \frac{p(x,z)}{q(z)}$$

for any given x
 "i.e. data", no matter
 what is $\{z\}$ {that we are trying several hypothesis}
 $\log p(x) = \text{Constant}$

KL-Div is always ≥ 0 (l): Variational lower bound

If $p(x)$ is the evidence
 This will be called as
 ELBO: Evidence lower Bound

Minimize (B)

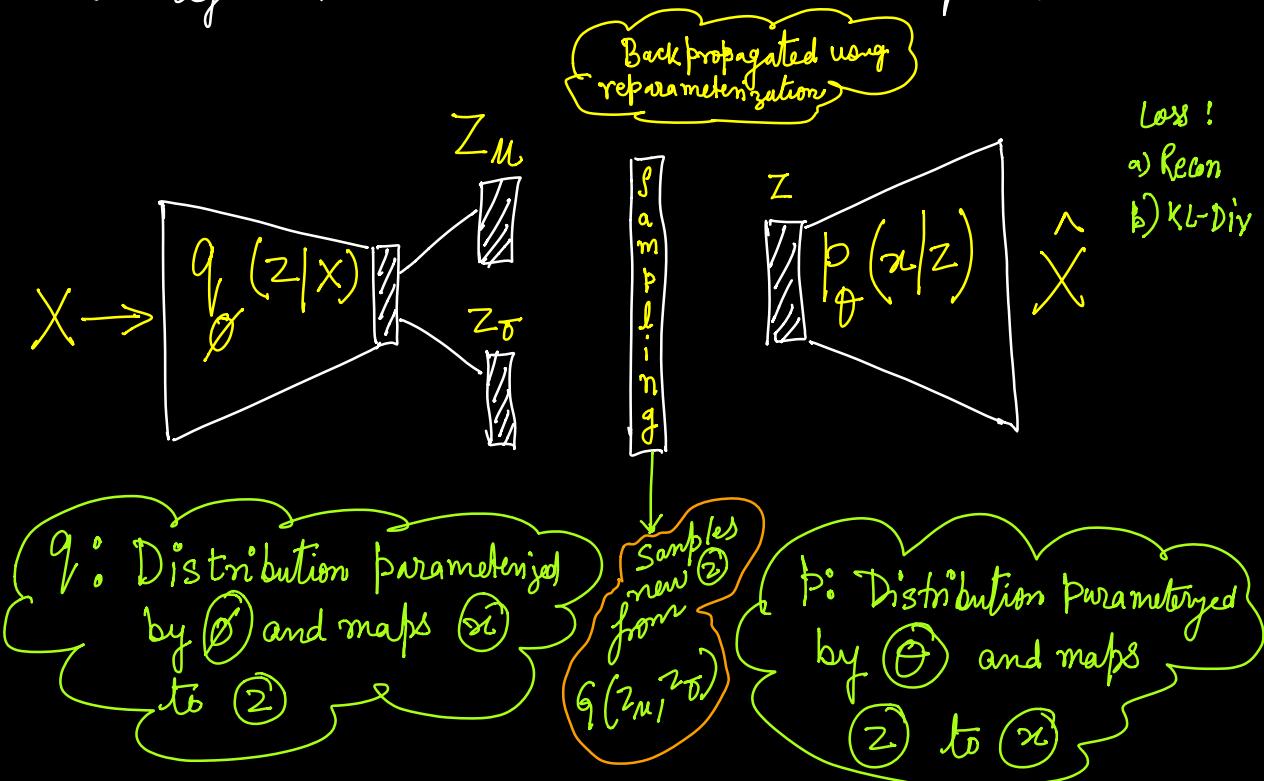
as a result this
 Maximized (C)

Constant (A)

This also means that minimizing the lower bound of a fn implies we are minimizing the fn $[\log p(x)]$. Minimizing (C) is equivalent to minimizing (B)

Let us again redraw variational autoencoder diagram:

(10)



ELBO

$$L(\theta, \phi; x, z) = \mathbb{E}_{\theta} \log p_{\theta}(x|z) - D_{KL}(q_{\theta}(z|x) \parallel p(z))$$

Reconstruction

KL-Div

This need to be minimized

maximized

minimized

Known tractable dist'n

$N(0, 1)$

Minimizing this KL-Div makes $q_{\theta}(z)$ as close to $p(z) = N(0, 1)$

But these z 's are coming from Encoder hence averaged over $q_{\theta}(z|x)$

This is the log likelihood of x given z . Maximizing it ensures z should have to generate x .

Just copy $\mathbb{E} \text{LBO}(\mathcal{L})$ from above:

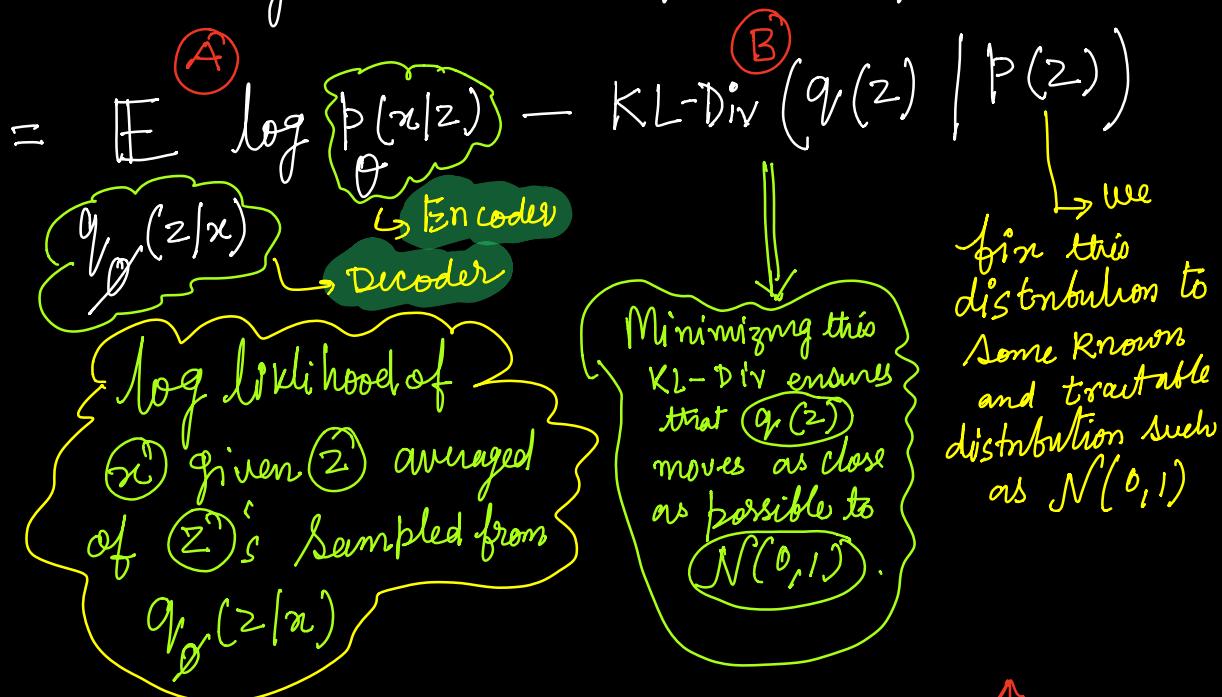
(ii)

$$\mathcal{L} = \sum_{\mathcal{Z}} q_{\theta}(z) \log \frac{p(x, z)}{q_{\theta}(z)}$$

Evidence lower bound
 $(z \in q_{\theta}(z|x))$

$$= \sum q_{\theta}(z) \log \frac{p(x|z) * p(z)}{q_{\theta}(z)}$$

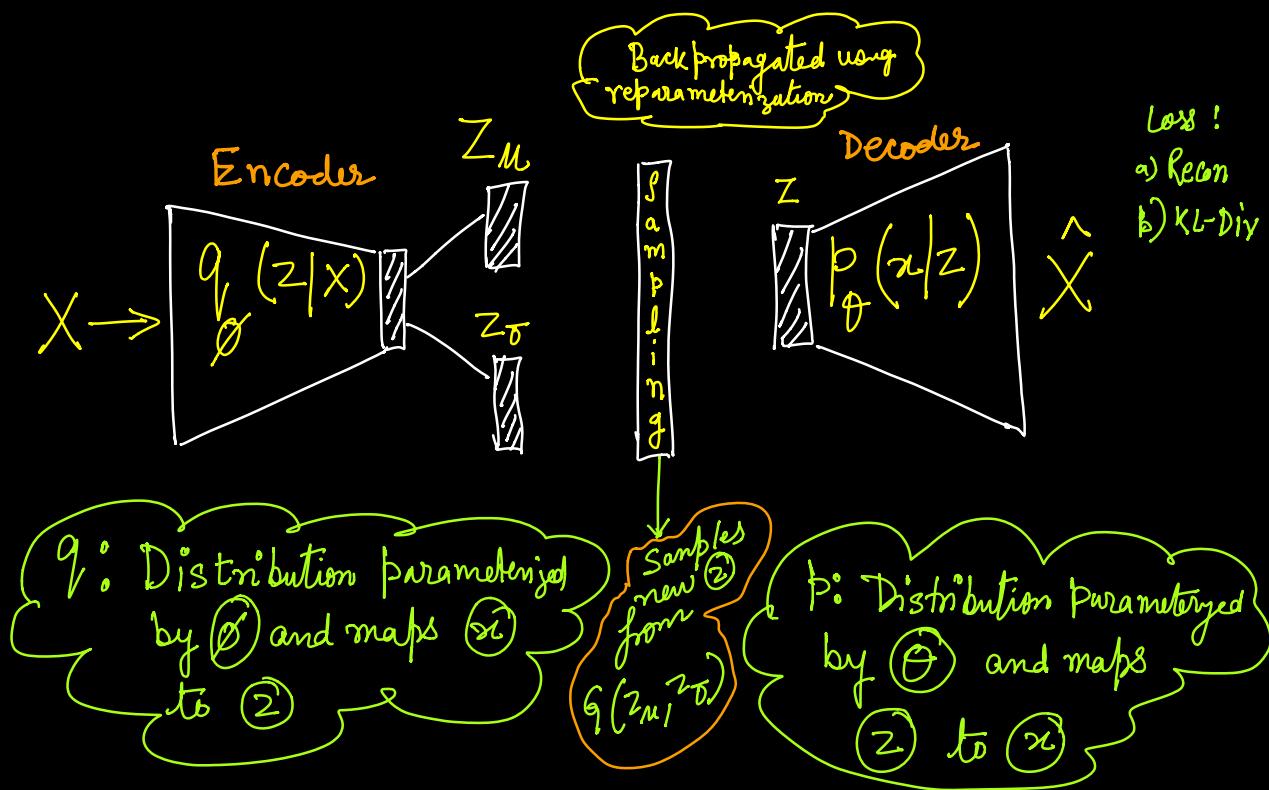
$$= \sum q_{\theta}(z) \log p(x|z) + \sum q_{\theta}(z) \log \frac{p(z)}{q_{\theta}(z)}$$



Still ELBO has one term to be maximized and other to be minimized (ELBO itself have to be Maximized). Hence converting it to a loss fn that can be minimized optimally.

$$\mathbb{E}_{q_{\theta}(z|x)} \left[\log p_{\theta}(x|z) \right]$$

- ① It is log likelihood of x
given z ,
② Averaged over z 's,
③ Sample from $q_{\theta}(z|x)$



∴ This likelihood minimization can be suitably realized as reconstruction loss minimization

- ① Get a sampled z from given x (via. Encoder $q_{\theta}(z|x)$)
- ② Feed z to Decoder (i.e. $p_{\theta}(x|z)$) and get \hat{x} .
- ③ In order to maximize the above averaged log likelihood, we expect that \hat{x} must be as close to x .

(13)

Observing the transitivity.

$$\begin{array}{ccc} \textcircled{X} & \xrightarrow{\quad} & \textcircled{Z} \\ \xrightarrow{\quad} p(z|x) & & \xrightarrow{\quad} p(\hat{x}|z) \end{array}$$

This ultimately becomes $p(\hat{x}|x)$

assuming this to be
a gaussian distribution

$$\Rightarrow p(\hat{x}|x) = e^{-\|x - \hat{x}\|^2}$$

taking log both sides

$$\textcircled{A} \quad \underbrace{\log p(x|z)}_{\text{Maximizing}} \rightarrow \log p(\hat{x}|z) \rightarrow \log (e^{-\|x - \hat{x}\|^2}) \downarrow -\|x - \hat{x}\|^2 \quad \textcircled{B}$$

 \textcircled{A} is equivalent tominimizing $\textcircled{-B}$ Hence final VAE loss : $\|x - \hat{x}\|_2$

$$L(\emptyset, \theta; x, z) = \min \quad \|x - \hat{x}\|_2 + \text{KL} \left(q_\emptyset(z) \middle\| p(z) \right)$$

Averaged over z 's
Sampled from the distribution

Predicted by $q_\emptyset(z|x)$ [ENCODER] $\mathcal{N}(0, I)$

Fed back to $p_\theta(x|z)$ [DECODER]
to get an $\textcircled{\hat{x}}$. Distribution of z 's sampled from $q_\emptyset(z_m, z_o)$

RECAP :

(14)

- ① First we observed that estimating $p(x)$ directly is Intractable.

$$p(x) = p(x_1, x_2, \dots, x_{10,000}) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots \cdot p(x_{10,000}|x_1, x_2, \dots, x_{9,999})$$

- ② Then we try to estimate $p(x)$ indirectly via latent variables (z). But that too observed to be Intractable.

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

(marginalized over all z 's)

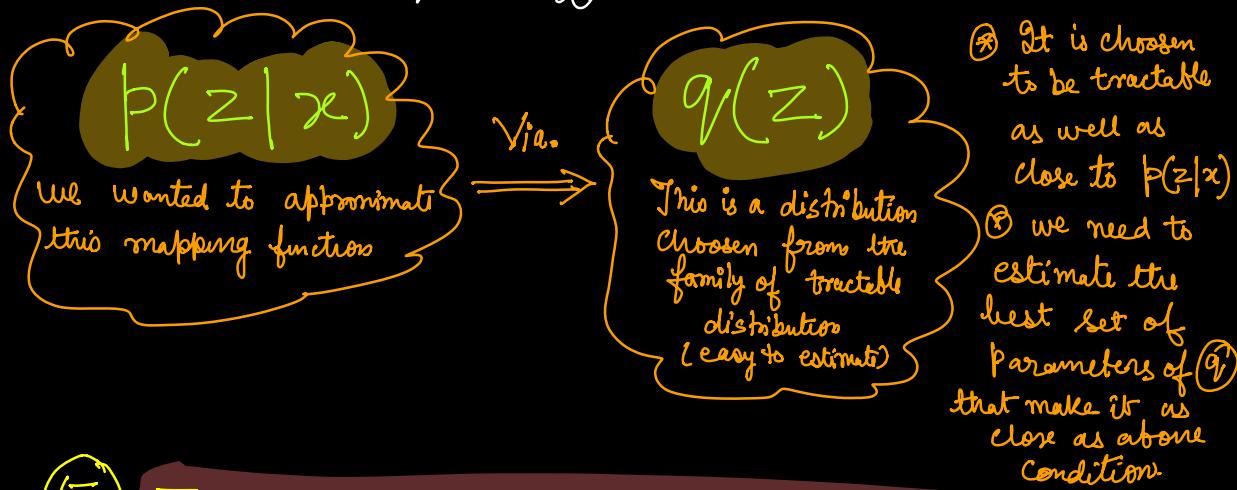
Hence this marginalized estimation is also not feasible.

Since this marginalization has to happen over all z 's (z 's can also be huge) [This is also intractable]
 may be dimensionally lesser than x

- ③ Later we tried to estimate $p(x)$ via.
- $$p(z|x) = \frac{p(z, x)}{p(x)}$$
- which was also not tractable due to $p(x)$ marginalization.

(4) Then we decided to approximate $p(z|x)$
 via. $q_v(z) \in \{ \text{family of tractable fn} \}$ (15)

$$q_v^*(z) = \arg \min_{q_v(z) \in \mathcal{L}} \text{KL-Div}(q_v(z) \mid p(z|x))$$



(5) The minimization of above KL-Div has been observed to be equivalent to maximizing ELBO(L)

$$\log p(x) = \text{KL-Div}(q_v(z) \mid p(z|x_i)) + \sum_z q_v(z) \log \frac{p(x_i, z)}{q_v(z)}$$



This also means that minimizing the lower bound of a fn implies we are minimizing the fn $[\log p(x)]$.

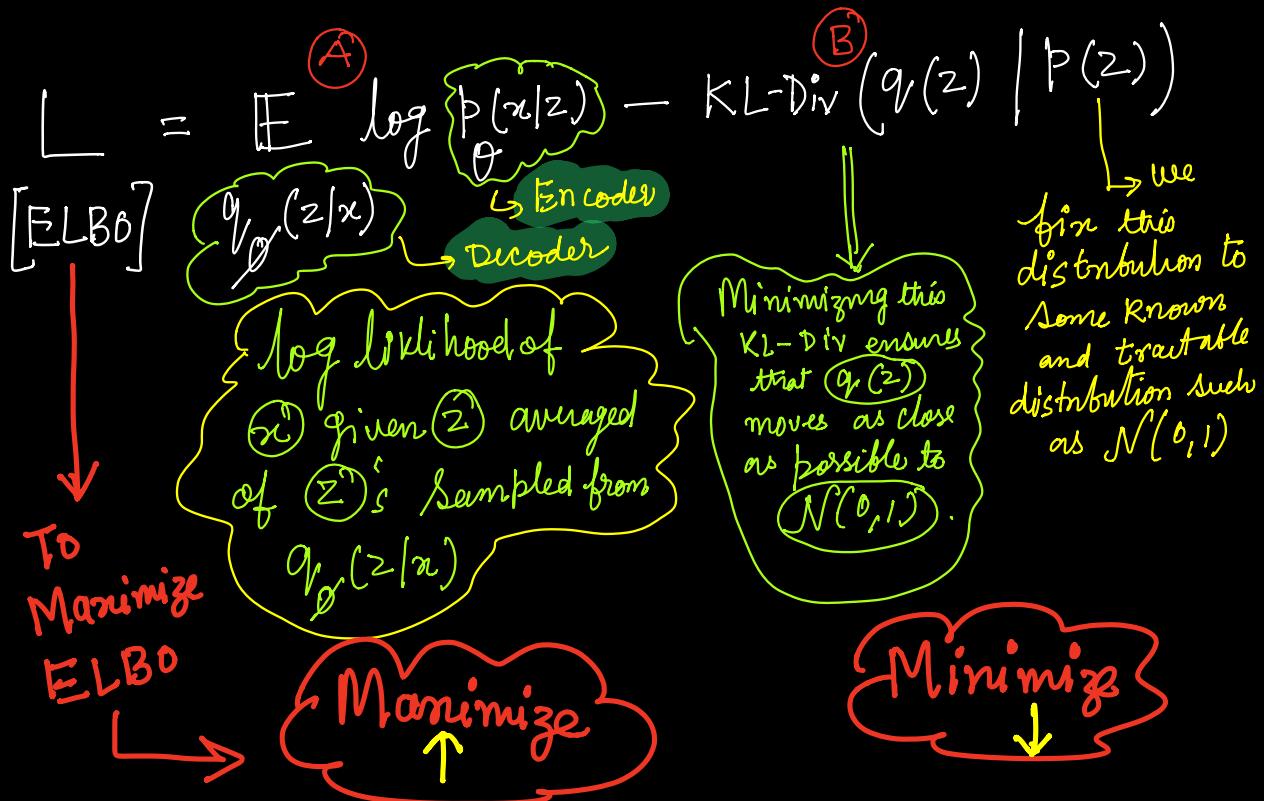
Minimizing (C) is equivalent to minimizing (B)

⑥ ELBO minimization is equivalent to

$\text{KL-Div}(q_\theta(z) \mid\mid p(z))$ minimization and

$\mathbb{E}_{q_{\theta}(z|x)} [\log p_{\theta}(x|z)]$ maximization

⑦



17

⑦ $\mathbb{E} [\log p_\theta(x|z)]$ maximization is equivalent to minimize $\|x - \hat{x}\|_2$.

$$\begin{aligned}
 & \textcircled{A} \quad \log p(x|z) \xrightarrow{\text{Maximizing}} \log p(\hat{x}|z) \xrightarrow{\text{Minimizing}} -\|x - \hat{x}\|_2^2 \\
 & \textcircled{B} \quad \text{is equivalent to} \\
 & \quad \min_{\hat{x}} \|x - \hat{x}\|_2^2
 \end{aligned}$$

⑧ Finally the loss function of VAE is

$$L(\phi, \theta; x, z) = \min \left\| x - \hat{x} \right\|_2 + KL(q_\phi(z|x) \left\| p_\theta(z) \right\|_{\mathcal{N}(0, I)}$$

↓
 Averaged over z 's
 Sampled from the distribution
 Predicted by $q_\phi(z|x)$ [ENCODER]
 Fed back to $p_\theta(x|z)$ [DECODER]
 to get an \hat{x} .
 ↓
 Distribution of z 's sampled from $q_\phi(z_M, z_O)$

Final Story :

(18)

- (A) Let us assume that our observation is MNIST ab.
- (B) It is a huge and varying dataset hence the images(X) must be following a very complex distribution ($p(X)$)
- (C) Using Encoder we try to reduce the dimensionality of (X) (to \mathbb{Z} latent variables), So that reduced dimensional latent vectors \mathbb{Z} follow some known tractable distribution.
- (D) Subsequently we also train a Decoder network that can sample latent vectors \mathbb{Z} from Gaussian distribution and generate images similar to MNIST.
- (E) Finally after training Encoder is discarded. we just need to sample \mathbb{Z} 's from $\mathcal{N}(0, 1)$ and feed it to Decoder that will generate new / new images similar to MNIST but not already present in it.