

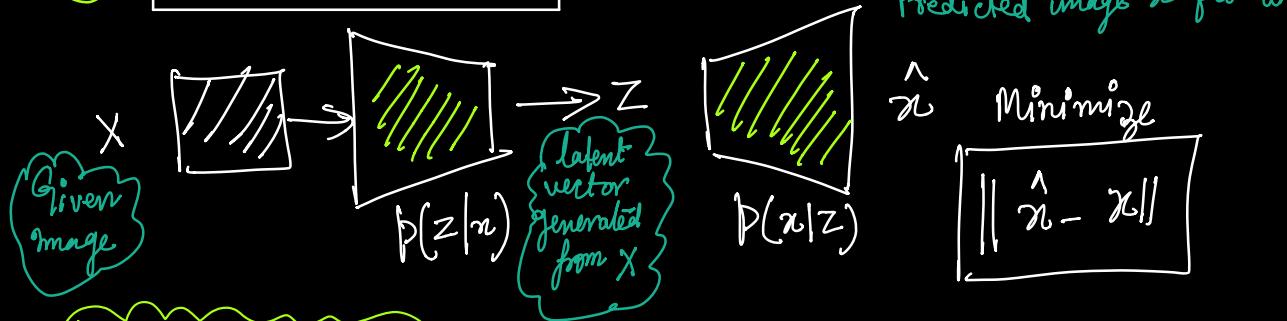
probability of z given x

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

①

) Can be converted into ② [Bottleneck]

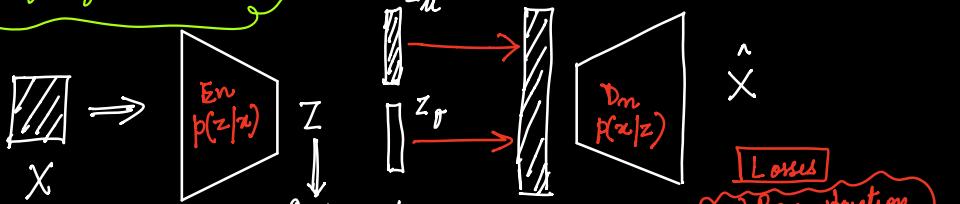
A Basic Autoencoder



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 z_u, z_o

z_u, z_o

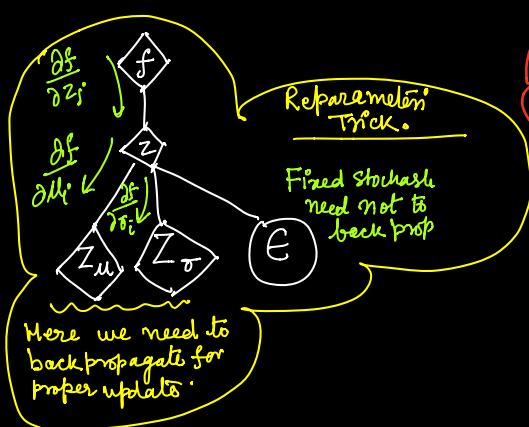
Sampling layer

sample a point from $G(z_u, z_o)$

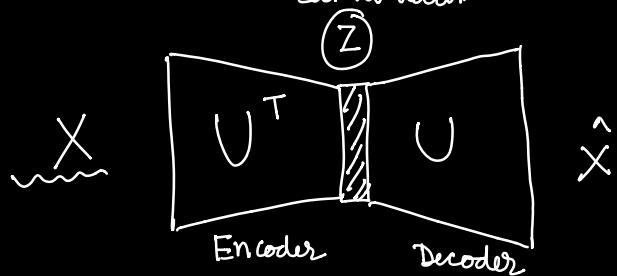
$(n \times 1)$

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

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



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

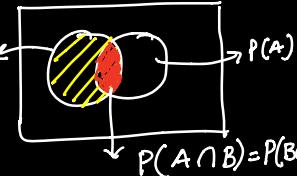
(3)

(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.)