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

Probability of "$Z$" given "$x$"

(1)

$(x)$ Can be converted into (2) [Bottleneck]

(A) **Basic Auto encoder**

$X$ — Given image

$p(z|x)$

$\rightarrow Z$ — latent vector generated from $X$

$p(x|z)$

Predicted image $\hat{x}$ from $Z$

$\hat{x}$   Minimize
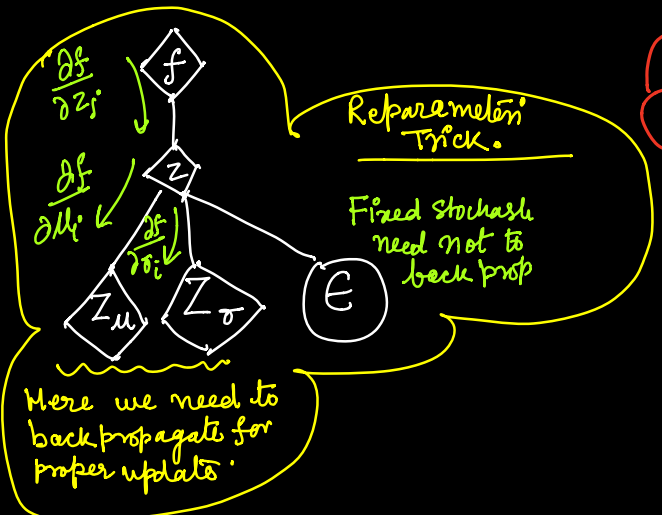
$$\| \hat{x} - x \|$$

Here our main question is (a) How to make sampling back propagable using reparametrization technique

(b) Why we want to train encoder so as to produce $Z's$ that follow a specific fixed distribution

(B) **Variational Auto encoders**

$X$ $\Rightarrow$ En $p(z|x)$   $Z$

Instead of the $Z$

$Z_u$

$Z_\sigma$

Dn $p(x|z)$   $\hat{X}$

we want our encoder to get $u, \sigma$   $Z_u, Z_\sigma$

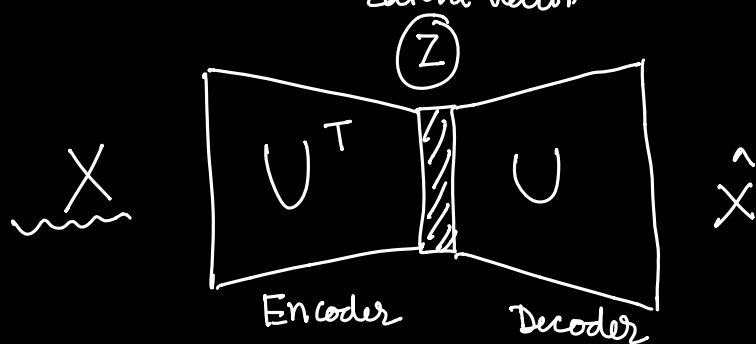Sampling layer sample a point from $G(Z_u, Z_\sigma)$

Losses
a) Reconstruction loss
min $\| X - \hat{X} \|$

b) $KL\text{-}Div(G(Z_u, Z_\sigma), N(0,1))$

$(n \times 1)$ $(n \times n)$ $\Rightarrow$ take only diagonal (variance)

$\frac{\partial f}{\partial z_i}$

$f$

$\frac{\partial f}{\partial u_i}$

$z$

$\frac{\partial f}{\partial \sigma_i}$

$Z_u$   $Z_\sigma$   $\epsilon$

Here we need to back propagate for proper update.

Reparameter Trick.

Fixed Stochash need not to back prop

$$KL\text{-}Div(G(Z_u, Z_\sigma), N(0,1))$$
$$= \frac{1}{2} \sum_{i=1}^{n} \left( u_i^2 + \sigma_i^2 - \log(1e\text{-}8 + \sigma^2) \right) - 1$$

In basic auto encoder data goes into the bottleneck and reconstructed

Latent vector

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

(reconstruction) error.



$\hat{X}$

Encoder          Decoder

(*) 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 do not ensure that Such A.E and PCA both learns the identical basis but may span the similar space

## Encoder (E)

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

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

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

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

## Decoder (D)

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

applying ① $\quad \hat{X} = U U^T X$

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

$\therefore \quad \min \| X - U U^T X \|$

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

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

In A.E $\boxed{U}$ may not be learned as orthonormal.

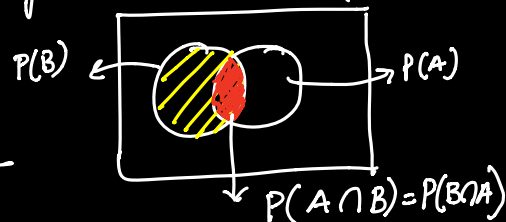One can train deep auto encoders with non-linearity in order to learn better representation.

# Important Concepts

## (A) Bayes Theorem:

① Conditional probability ÷ events are (A) & (B)

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

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

P(B) ← → P(A)

$$P(A \cap B) = P(B \cap A)$$

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

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

Hypothesis → Evidence → Prior

Assuming 1 student in the class of ㉑ has flu.
Event A: Student is ABC ⇒ $\boxed{P(A) = 1/20}$
Prior Probability (w/o any Knowledge)

Evidence B: 5 Girls & 15 Boys
Now given this new evidence what is the probability that ABC has flu.
$$P(A|B) = \frac{1}{5} \text{ (goes up) if girl} \Rightarrow 0 \text{ (if student is a Boy)}$$

New evidences are going to influence the Hypothesis

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

Here we can have 3 things:

Information — Since $P(x) \in (0,1)$ hence $(-\log)$.

| X — Event | P(X) — Probability | $-\log(P(X))$ |
|---|---|---|
| ① Virat Scored a century. | ↑ (highly probable event) | ↓ (less information) |
| ② Kenya wins Cricket world Cup | ↓ (rare event) | ↑ (high information) |
| ③ Tomorrow it rain or don't | 1 (Certain event) | 0 [no information] |

So basically rare events carry more information

# Ⓒ Average of Information is Entropy (H):

The expected value of information w.r.t any event Ⓧ averaged over all values Ⓧ Can attain is Entropy Ⓗ.

This is the expected value of $[\log p(x)]$ wrt $p(x)$

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

Summation over all x's

Probability of that Ⓧ to happen

Information Content in any Ⓧ

# Ⓓ KL-Divergence (KL-Div): In order to compute the

similarity between two distributions say Ⓟ and Ⓠ

KL-Div ( P || q ) Can be used defined as the

KL-Div of Ⓠ distribution wrt Ⓟ.

(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 wrt $p(x)$ as KL-Div is wrt $p(x)$

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

⇓ This is the Cross entropy between (p) and (q) distributions.

⇓ This is (-ve) entropy of (p) distribution

Now both expectations are wrt p(x)

Hence,

KL-Div $\left(p(x) | q(x)\right)$ Can be formally defined as the difference between average information of q(x) wrt p(x) and that of p(x) wrt p(x).

$$KL\text{-}Div\left(p(x) | q(x)\right) = -\sum p(x) \log q(x) + \sum p(x) \log p(x)$$

$$= \sum p(x) \log \frac{p(x)}{q(x)}$$

$$= -\sum p(x) \log \frac{q(x)}{p(x)}$$

⊛ KL-Div is not symmetric as $KL\text{-}Div(p|q) \neq KL\text{-}Div(q|p)$

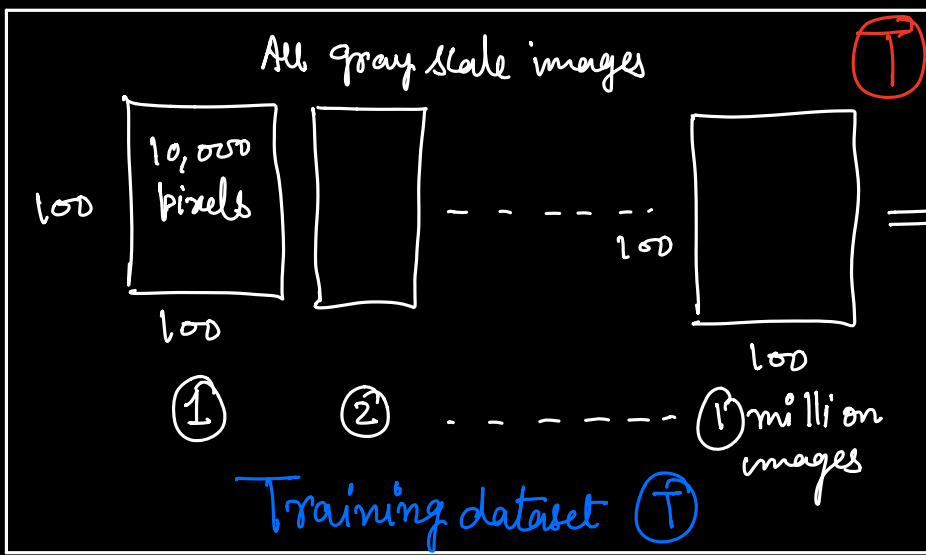↳ Hence it is a distance measure b/w a divergence.

⊛ KL-Div $\geq 0$ it is always +ve and b/w (0 & 1)

$$\therefore \quad KL\text{-}Div\left(q(z) \| P(z|x)\right) = -\sum q(z) \log \frac{p(z|x)}{q(z)}$$

(we will come back to this.)

All gray scale images ⑴

10,000 pixels

100

100

100

100

① ② - - - - - ①'million images

100

Training dataset ⑪

Let us assume that we have a very complex and huge training dataset ⑪ of 1 million images of size $100 \times 100 = 10,000$ pixels per image.

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.

But any $100 \times 100$ 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-occurences patterns.

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

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

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

Joint Probability distribution

Depending upon our training dataset ⑪, we wanted to estimate $P_\theta(x)$ where $P_\theta(x) =$ Probability of $(X \in T)$, with the distribution parameterized over $[\theta]$.

$$\theta^* = \arg\max_{\theta} \left[ P(x \in T) \right]$$

minimizing the rest.

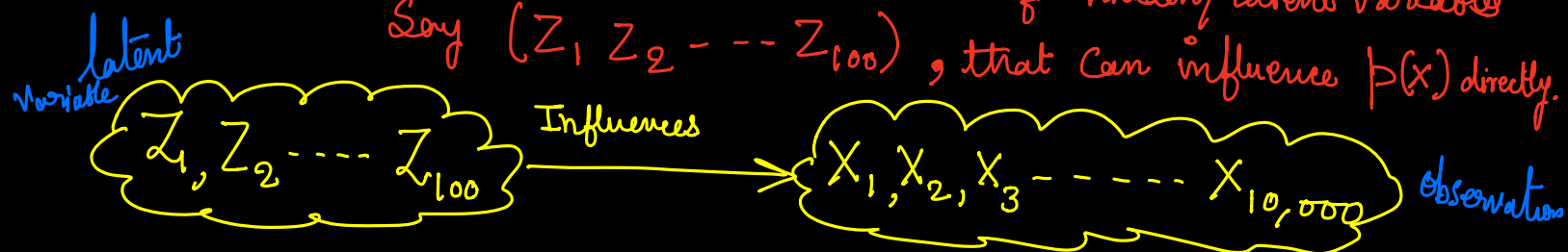But why we are interested to compute $P_\theta(X)$

→ (a) <u>Classification</u>: Helps us to discriminate b/w images that are coming from $T$ or rest.

→ (b) <u>Generative modeling</u>: 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, interpolation b/w 2 views/poses.

For this image sampling experiment $P(X) = P(X_1 X_2 ---- X_{10,000})$ 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.

$$P(X) = P(X_1 X_2 ---- X_{10,000}) = P(X_1) \cdot P(X_2/X_1) \cdot P(X_3/X_1 X_2) ----$$
$$---- \quad P(X_n|X_1 X_2 --- X_{n-1})$$

→ 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 $(Z_1, Z_2 --- Z_{100})$, that can influence $P(X)$ directly.

latent
variable
$Z_1, Z_2 ---- Z_{100}$ ——Influences——→ $X_1, X_2, X_3 ------ X_{10,000}$ observations

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 - - -.

external Parameters

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

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