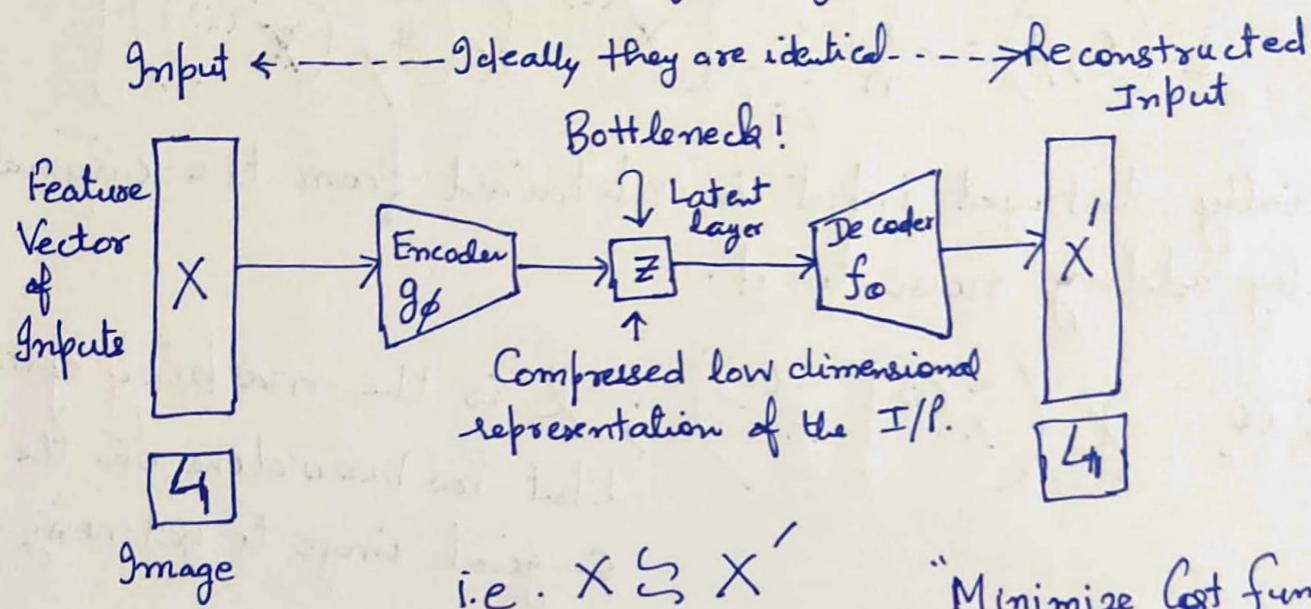


Autoencoders

①

① Let's take the example of 'Image Reconstruction'.

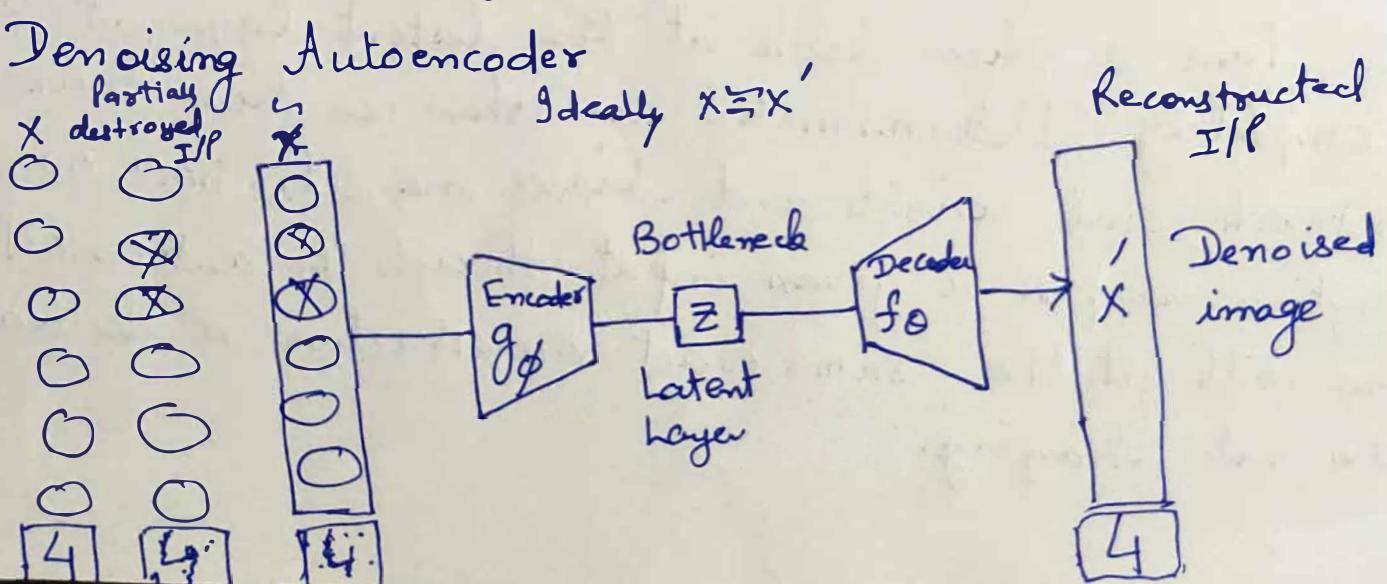


Stacked Autoencoder

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

θ and ϕ are the parameters which define the encoder and decoder. These are basically the weights and biases of the encoder and decoder modules. as these modules are basically the neural networks. i is an index which defines the features of I/P and it varies from 1 to n . $f_\theta(g_\phi(X^{(i)}))$ represents the data flow from the encoder to the decoder.

③ Denoising Autoencoder



①

$$\hat{x}^{(i)} = \chi(\hat{x}^{(i)} | x^{(i)})$$

②

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

② Partially destroyed input is obtained from the original image by adding noise to it.

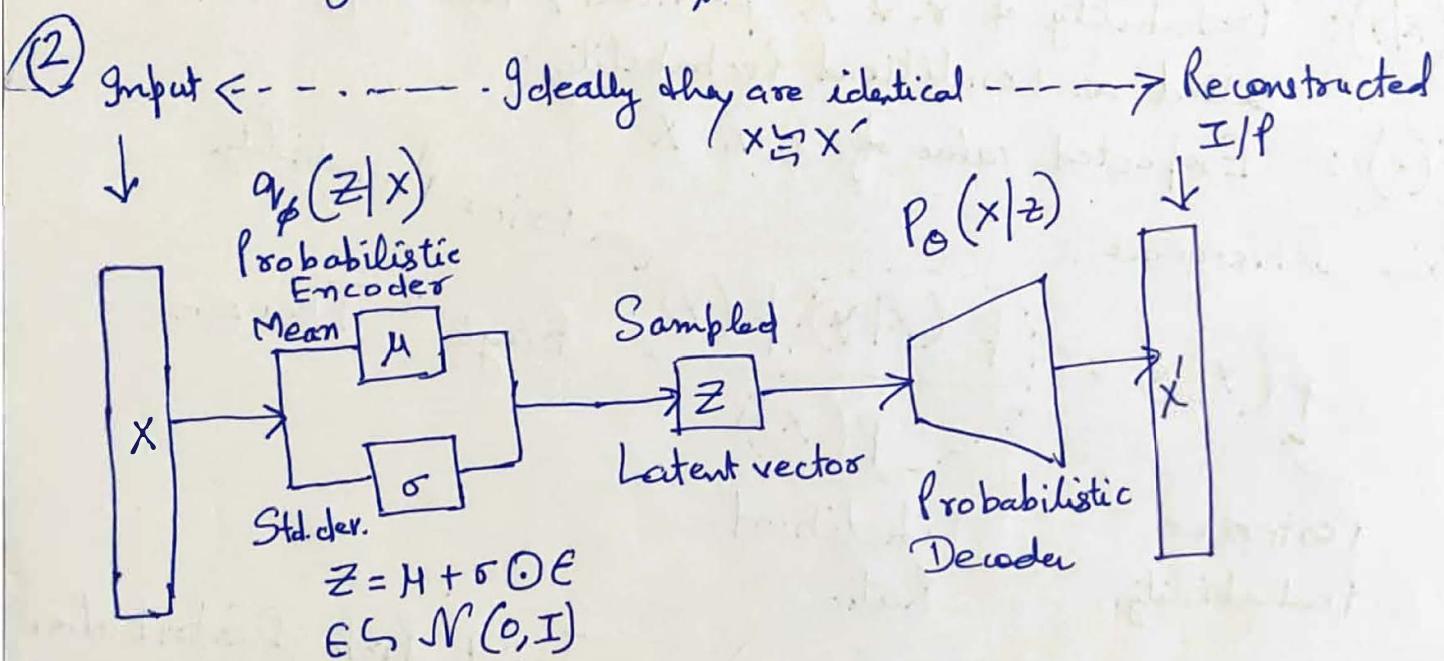
$$\hat{x}^{(i)} = \chi(\hat{x}^{(i)} | x^{(i)}) ; \chi \text{ is the masking operator that has been done on the original image to get noisy I/O}$$

③ Feed $\hat{x}^{(i)}$ into the autoencoder. Since z is the compressed version of the original image, and it will represent the original image by small number of features in the vector. If we can reconstruct the image back, then lot of BW^(Bandwidth) can be saved while transmitting an image. Noise will be added to the original during training process.

Variational Autoencoders :-

④ Take a closer look at the latent layer, it is completely Deterministic i.e. when we have trained the network and weights and biases are set, now on anytime we pass a given input through the autoencoder, we will get the same exact reconstruction as the weights are not changing.

① In contrast, the variational autoencoders (VAEs) ③ introduce an element of probabilistic twist on this idea of autoencoding. This will allow us to not only perform reconstructions, but generate new images similar to the input instances, similar to the input data, but importantly not forced^{strictly} to the reconstructions only.



○ → Circled dot operator → XNOR gate

Hadamard product
 element-wise \times of
 matrices of same size

Variational Autoencoder (VAE)

③ Replaced single deterministic layer Z with a random sampling operation. Now, we define a mean and standard deviation that captures a PD over the latent variable Z . We have converted a single vector space into a vector of means of latent variables and a vector of SDs of these latent variables that parametrize the PDFs over latent variable.

$$\textcircled{1} \quad \text{Loss function is given by} \quad \textcircled{4}$$

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

NOTE:-

$p(x)$: Probability of r.v. X .

$p(x|y)$: Probability of r.v. X provided Y has happened, also called as Conditional Probability.

$E(\cdot)$: Expected value of r.v. X

KL Divergence

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

↑ Prior Probability
↑ Likelihood Ratio
Posterior Probability

$$= \frac{p(x,y)}{p(x)} \rightarrow \text{Joint Prob. Distribution} \quad \textcircled{1}$$

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,

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

$$P(x) = \sum_{i=1}^5 p(x, y_i)$$

$$= \sum_{i=1}^4 p(x|y_i)p(y_i)$$

Expectation of r.v. X :- $E(x)$

Expected value of random variable is a weighted average of the possible values of x it can take; each value being weighted according to the probability of that event and is defined as

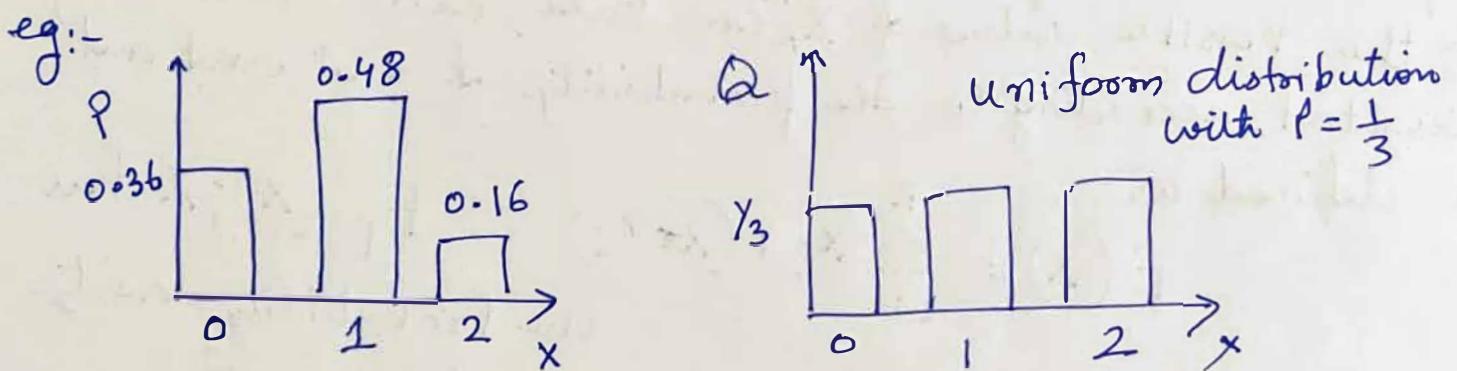
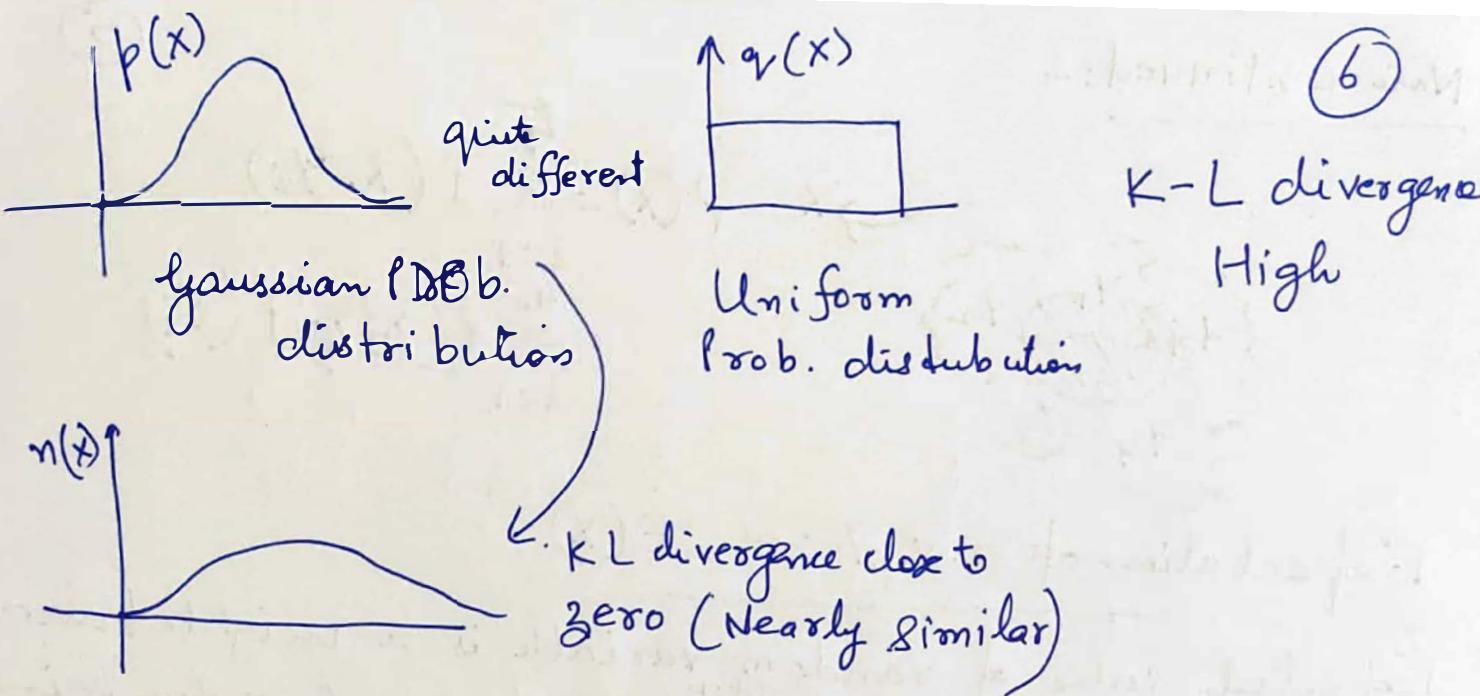
$$E(x) = \sum_{i=1}^K x_i p(x=x_i) \Leftarrow E_p(x); X \text{ has the probability density function as } P.$$

K-L Divergence:-

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

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

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



$$\begin{aligned}
 D_{KL}(Q || P) &= D_{KL}(Q || P) \\
 &= \sum_x Q(x) \log \left[\frac{Q(x)}{P(x)} \right] \\
 &= \frac{1}{3} \ln \left\{ \frac{0.33}{0.36} \right\} + \frac{1}{3} \ln \left\{ \frac{0.333}{0.48} \right\} + \frac{1}{3} \ln \left\{ \frac{0.333}{0.16} \right\} \\
 &\approx 0.09637 \text{ nats} \quad (\text{very small})
 \end{aligned}$$

Properties of KL divergence :-

- 1) $D_{KL}(P || Q) \text{ and } D_{KL}(Q || P) \geq 0$
- 2) $D_{KL}(P || Q) \neq D_{KL}(Q || P)$ (Not Symmetric)

Note 2:-

① Let X be an $K \times 1$ random vector (A random vector also called "multivariate r.v." is an K -dimensional column vector $X \in \mathbb{R}^{K \times 1}$ whose entries are random variables). The X is said to be multivariate normally distributed with mean μ and covariance Σ :

$$X \sim N(\mu, \Sigma)$$

iff its PDF is given by

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

where μ is an $K \times 1$ real vector and Σ is an $K \times K$ positive definite matrix.

$$\mu = E[X] = (E[x_1], E[x_2], \dots, E[x_K])^T$$

K-dimensional Mean Vector

$$\Sigma_{i,j} = E[(x_i - \mu_i)(x_j - \mu_j)] = \text{Cov}[x_i, x_j]$$

K × K - Covariance Matrix

The inverse of the covariance matrix is called Precision Matrix denoted by Σ^{-1} .

Theorem: Let x be an $K \times 1$ random vector. Assume two multivariate normal distributions P and Q specifying the probability distribution of x as

$$P(x): x \sim N(\mu_1, \Sigma_1)$$

$$Q(x): x \sim N(\mu_2, \Sigma_2)$$

① Then the Kullback - Leibler divergence of P from Q⁽⁸⁾
is given by (P & Q have the same dimension k)

$$KL[P \parallel Q] = \frac{1}{2} \left[(\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) + \text{tr}(\Sigma_2^{-1} \Sigma_1) - \ln \frac{|\Sigma_1|}{|\Sigma_2|} - k \right]$$

(1)

Note:- Multivariate Normal Distribution is a generalization of the one-dimensional (univariate) normal distribution to higher dimensions. A random vector \mathbf{x} is said to be k-variate normally distributed if every linear combination of its k components has a univariate Normal distribution. The multi-variate distribution is used to describe any set of correlated real-valued random variables each of which clusters around a mean value.

Proof of eq.(1) :-

We know that

$$KL[P(x) \parallel Q(x)] = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \quad (2)$$

Also,

$$P(x) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_1|}} \exp \left(\frac{-(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}{2} \right) \quad (3)$$

Take log on both sides of equation (3)

$$\Rightarrow \log P(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)$$

— (4)

Similarly,

$$\log Q(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)$$

— (5)

① Eq(2) can be rewritten as ⑨

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) [\log P(x) - \log Q(x)] - (6)$$

② Substituting equations (4) and (5) in eq(6), We get

$$\begin{aligned} KL(P(x) \parallel Q(x)) &= \sum_x P(x) \left\{ -\frac{\kappa}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right. \\ &\quad \left. + \frac{\kappa}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\} \end{aligned}$$

③ On Simplification, We get

$$\begin{aligned} KL(P(x) \parallel Q(x)) &= \sum_x P(x) \left[\frac{1}{2} \log \left(\frac{\Sigma_2}{\Sigma_1} \right) + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right. \\ &\quad \left. - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right] - (7) \end{aligned}$$

④ From eq(7), let's consider expression part by part:

$$\sum_x P(x) \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = E_p \left[\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right] - (8)$$

⑤ Let's look at some of the identities which are required.

$$E(x^T A x) = E \left(\underset{@}{\text{tr}}(x^T A x) \right) = E \left(\text{tr}(A x x^T) \right) = \text{tr}(E(A x x^T))$$

⑥ Also let's have a trick for Trace and Expectation:
If x is a scalar, then $E(x) = E(\text{tr}(x))$, since trace of a scalar is scalar.

$$\text{tr}(AB) = \text{tr}(BA), \quad \text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

$$\text{tr}(ABC) \stackrel{(b)}{\neq} \text{tr}(ACB) \quad \text{tr}(ACB) \stackrel{(d)}{=} \text{tr}(BAC)$$

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

$$\frac{1}{2} \mathbb{E}_P \left[(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right] = \mathbb{E}_P \left[\text{tr} \left(\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right) \right] \quad (10)$$

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

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

↑ Covariance Matrix

$$= \text{tr} \left[\Sigma_1 \frac{1}{2} \Sigma_1^{-1} \right] = \frac{1}{2} \text{tr} \left[I_K \right] = \frac{K}{2}. \quad - (9)$$

① Consider the second part of eq (7) for equation to calculate $\text{KL}(P(x) || Q(x))$, which is

$$\sum_x P(x) \left\{ \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\}$$

let's say
 $(A + B)^T \Sigma_2^{-1} (A + B)$
 $(A^T + B^T) \Sigma_2^{-1} (A + B)$
 $A^T \Sigma_2^{-1} A + A^T \Sigma_2^{-1} B + B^T \Sigma_2^{-1} A$
 $+ B^T \Sigma_2^{-1} B$
 Basically same

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

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

$$= \mathbb{E}_P \left[\frac{1}{2} (x - \mu_1)^T \Sigma_2^{-1} (x - \mu_1) + (x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right. \\ \left. + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

① Rewriting the previous expression:

(11)

$$E_p \left[\frac{1}{2} (x - M_1)^T \Sigma_2^{-1} (x - M_1) \right] + E_p \left[(x - M_1)^T \Sigma_2^{-1} (M_1 - M_2) \right]$$

$$+ E_p \left[(M_1 - M_2)^T \Sigma_2^{-1} (M_1 - M_2) \right]$$

CONSTANT

② $\Rightarrow \frac{1}{2} + \text{tr}(\Sigma_2^{-1}) + (M_1 - M_2)^T \Sigma_2^{-1} (M_1 - M_2) + 0$ (10)

similar to earlier derivation

$\therefore E(\text{constant}) = \text{Constant}$

③ $E_p \left[(x - M_1)^T \Sigma_2^{-1} (M_1 - M_2) \right]$

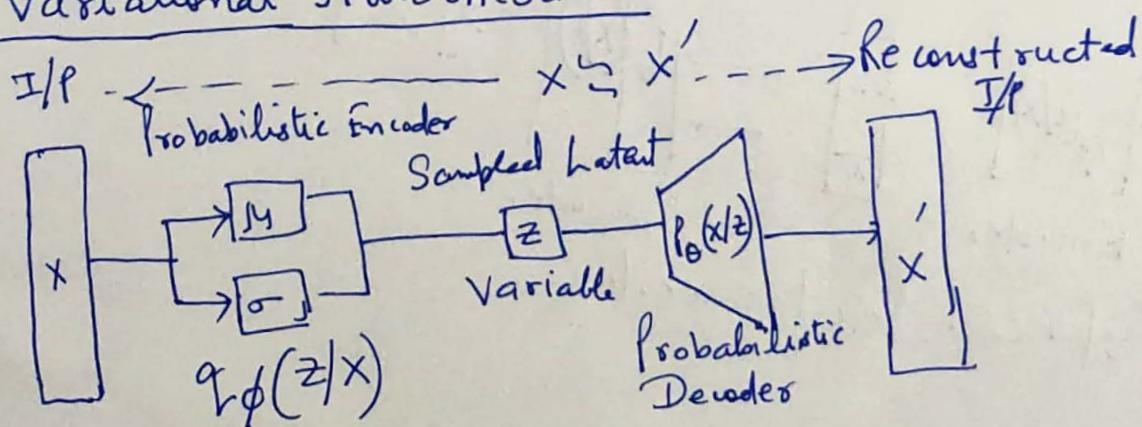
$$= \left[(E_p(x) - M_1)^T \Sigma_2^{-1} (M_1 - M_2) \right]$$

$$= (M_1 - M_1)^T \Sigma_2^{-1} (M_1 - M_2) = 0$$

④ Substitute eq(9) and eq(10) in eq(7) to obtain final and simplified KL-Divergence.

$$\text{KL}(P(x) || Q(x)) = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr}(\Sigma_2^{-1}) + (M_1 - M_2)^T \Sigma_2^{-1} (M_1 - M_2) \right]$$

Variational Autoencoder:



$$z = \mu + \sigma \odot \epsilon$$

$$\epsilon \sim N(0, I)$$

(12)

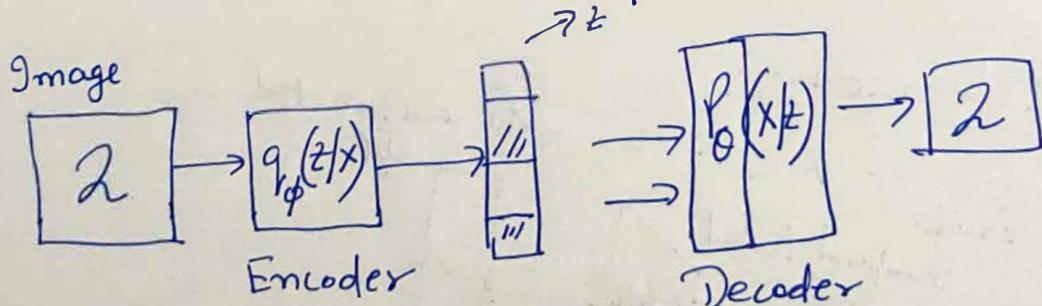
① Loss function in a VAE is given by

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

② Expectation operation which is carried on to the Prob. distribution $q_\phi(z|x)$ ∵ We can write $E[x] = E_p$; Here $P = q_\phi(z|x)$. $P_\theta(x|z)$ is any r.v. Y.

The goal of VAE :-

The goal of VAE is to find a distribution $q_\phi(z|x)$ of some latent variable z from which we can sample $z \sim q_\phi(z|x)$ to generate new samples $x \sim P_\theta(x|z)$ i.e. We want to find $q_\phi(z|x)$ from which we shall be sampling z from $q_\phi(z|x)$. . . z will be our sampled latent vector. Once it is done, then we will feed samples of z to the decoder and then we will generate new samples x' . The new samples mean that there are not present in x .



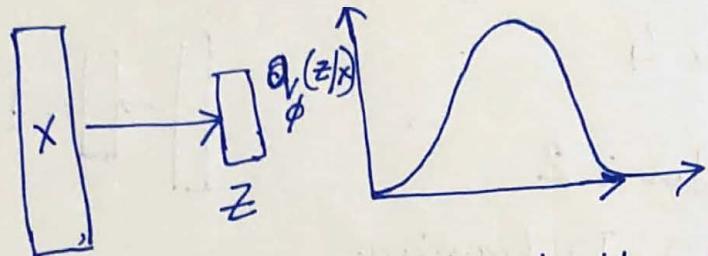
Typical Autoencoder.

① Traditional autoencoder data flow with one hidden unit is given by (13)

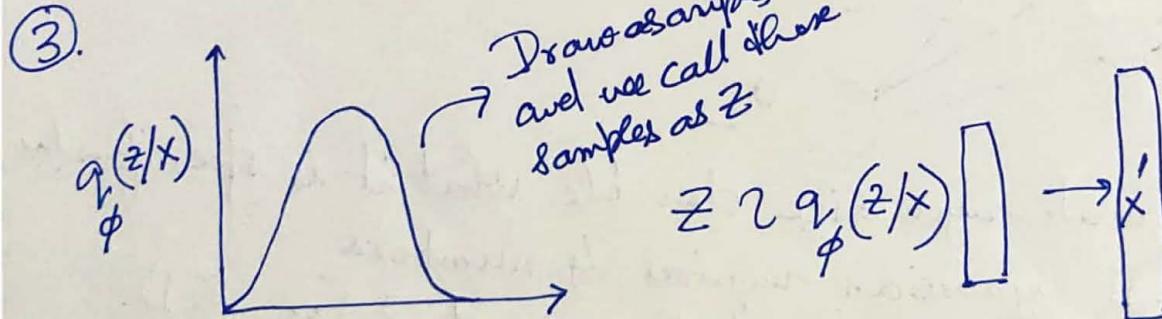
$$z = f(w_1^T x + b_1)$$

$$\hat{x} = g(w_2^T z + b_2)$$

② In the case of Variational Autoencoder:



At the end of the encoder, we don't get a value but we get a distribution: $q_\phi(z|x)$. Or more precisely we get the parameters of a distribution. In particular, encoder outputs a mean and variance which represent a Gaussian. Bayesian Machine learning is about learning distributions and not learning point estimates. \therefore Instead of actually find z , we are finding $q_\phi(z|x)$ that tells us the PDF of z given x with ϕ as the neural network (encoder.NN) parameters.



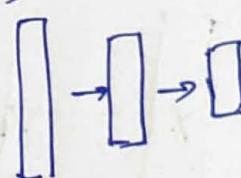
We trace a distribution $q_\phi(z|x)$ and we need actual numbers to pass in through the rest of the NN. We sample the gaussian distribution to obtain z which is sample vector and then we can pass z from the decoder. Follow the flow of information through decoder which is again a NN.

① The output of the decoder is again a distribution 14
 $p_\phi(x|z)$ from this distribution, we generate samples

② eg:- Parameterizing a Gaussian:-

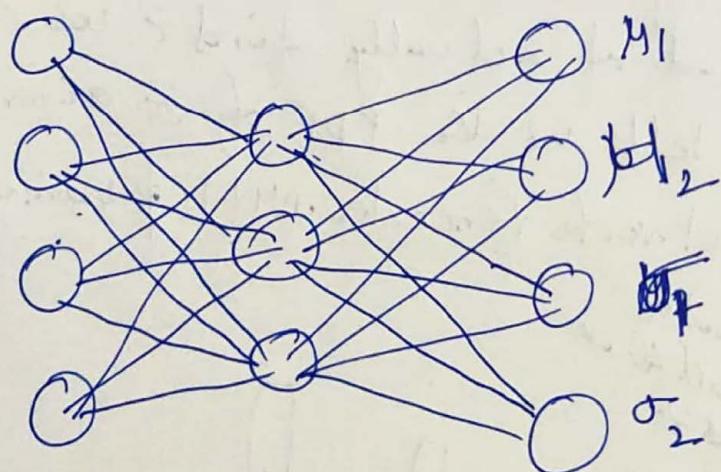
Suppose our encoder has the sizes $(4, 3, 2)$

- Input dimensionality is 4
- HL size is 3
- Size of latent vector is 2



i.e. $q_\phi(z|x)$ is a 2-D Gaussian.

We will use an axis-aligned Gaussian (one variance param per dim) rather than a full covariance Gaussian (requires full $D \times D$ matrix)



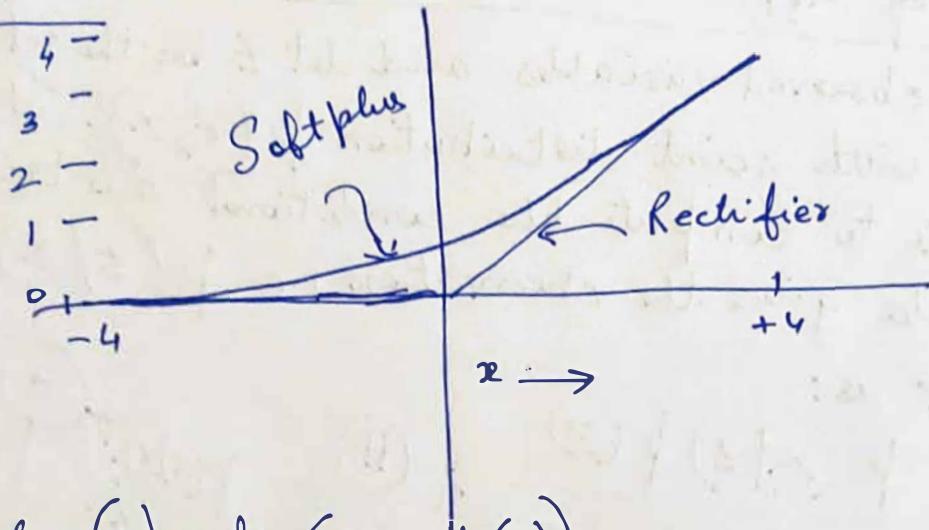
③ Make the final layer size double what it is specified

e.g. A 2-D Gaussian requires 4 numbers
Use first 2 for mean and the next 2 for SD.
NN can output any number. How can we fix this?

① Softplus :-

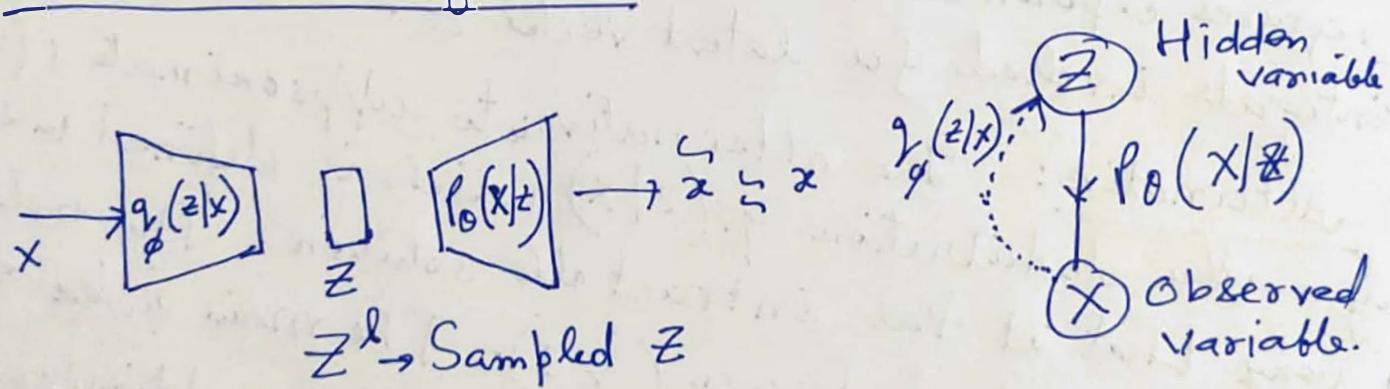
Nonlinearities

15



- $\text{Softplus}(a) = \log(1 + \exp(a))$
- Smooth version of ReLU
- Approximately linear when $|x|$ is large. and looks more like exponential when $|x|$ is small.
- Smooth / continuous / Differentiable.
- and always > 0

② Recall the VAE goal :-



Learn $q_\phi(z|x)$

Sample $q_\phi(z|x)$ to produce z^l

Feed z^l to decoder

Generate \hat{x} , new samples, but not present in input vector

③ Let's now look at variational autoencoders in a more analytical and mathematical framework:

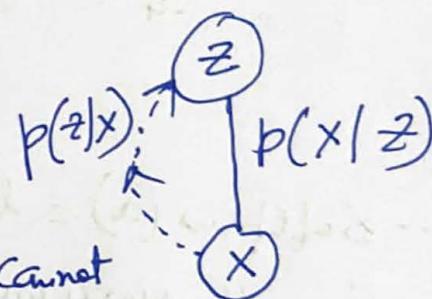
① The problem of Approximate Inference:-

(16)

Let X be a set of observed variables and 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 variable given the observations i.e. $p(z|x)$.

We can write it as:

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} \rightarrow \textcircled{A}$$



② Evaluating \textcircled{A} is difficult because $p(x)$ cannot be solved.

Reason: $p(x) = \int p(x|z)p(z) dz = \int p(x, z) dz$

The integral is not available in closed form or is intractable (requires exponential time to compute) due to multiple integrals involved for latent vector Z .

③ 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 using variational Inference (VI). The main idea of VI is to pose the inference problem as an optimization problem. By modeling $p(z|x)$ using $Q(z|x)$ where $Q(z|x)$ has a simple distribution such as Gaussian.

As discussed earlier let's calculate KL divergence between $p(z|x)$ and $Q(z|x)$.

$$\text{① } D_{KL} \left(Q_\phi(z|x) || P_\theta(z|x) \right) = \sum_z Q_\phi(z|x) \log \left(\frac{Q_\phi(z|x)}{P_\theta(z|x)} \right) \quad (17)$$

We will calculate the KL-divergence and will try to minimize the divergence and then $P_\theta(z|x)$ will be approximated.

$$\begin{aligned} \text{② } D_{KL} \left(Q_\phi(z|x) || P_\theta(z|x) \right) &= \underset{z \sim Q_\phi(z|x)}{\mathbb{E}} \left[\log \left(\frac{Q_\phi(z|x)}{P_\theta(z|x)} \right) \right] \\ &= \underset{z \sim Q_\phi(z|x)}{\mathbb{E}} \left[\log(Q_\phi(z|x)) - \log(P_\theta(z|x)) \right] \end{aligned} \quad (B)$$

Substituting ① in ②.

$$\begin{aligned} \text{② } &= \underset{z}{\mathbb{E}} \left[\log(Q_\phi(z|x)) - \log \left\{ \frac{P_\theta(x|z)P_\theta(z)}{P_\theta(x)} \right\} \right] \\ &\quad \text{② } z \sim Q_\phi(z|x) \end{aligned}$$

$$\begin{aligned} \text{② } &= \underset{z}{\mathbb{E}} \left[\log(Q_\phi(z|x)) - \log P_\theta(x|z) - \log P_\theta(z) + \log P_\theta(x) \right] \end{aligned}$$

Since the expectation is over z and $P_\theta(x)$ does not involve z , it can be moved out.

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

Rewriting the above equation,

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

$$\begin{aligned}
 \log P_\theta(x) - D_{KL}[Q_\phi(z|x) || P_\theta(z|x)] & \\
 = E[\log(P_\theta(x|z))] - E[\log Q_\phi(z|x) - \log P_\theta(z)] & \quad (18) \\
 = E[\log(P_\theta(x|z))] - D_{KL}[Q_\phi(z|x) || P_\theta(z)] & \quad (2) \\
 \end{aligned}$$

This is VAE objective function where the first term represents the reconstruction likelihood and the second term ensures that our learned distribution Q is similar to the prior distribution P .

$$\text{Loss} = -\text{objective function} \quad (\text{Reconstruction loss}) + D_{KL}[Q_\phi(z|x) || P_\theta(z)] \quad (\text{Regularization term})$$

Hence Prooved !!

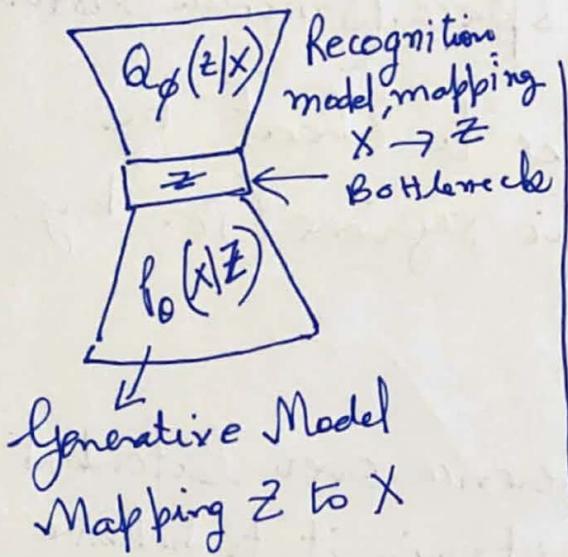
$$\begin{aligned}
 \text{Also } \log P_\theta(x) - D_{KL}[Q_\phi(z|x) || P_\theta(z|x)] & \\
 = -L(\theta, \phi) & \\
 \therefore \text{Our target is to find optimal } \theta, \phi \text{ such that} &
 \end{aligned}$$

$$\theta^*, \phi^* = \arg \min_{\theta, \phi} L(\theta, \phi)$$

① More Intuition about Loss function :-

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

Log likelihood → ← Regularizer →



So when we take log of gaussian, we get a squared error between the data sample \hat{x} and mean of the gaussian distribution if deviate from the gaussian distribution. $N(0, 1) = P_\theta(z)$

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

$$\log P_\theta(x|z) \propto (x - \mu_\theta(z))^T \Sigma_\theta^{-1}(z) (x - \mu_\theta(z))$$

← Squared Reconstruction error →

It basically tells us that how far is x from z

This squared reconstruction also is known as Data Fidelity term in inverse problems.

$$D_{KL}[Q_\phi(z|x) || P_\theta(z)]$$

↑
Inferred Latent Distribution
↓
Prior on Latent Distribution

① Common Choice of Prior:-

(20)

Normal Distribution

$$P_0(z) = \mathcal{N}(H=0, \sigma^2=1)$$

- It encourages encodings to distribute encodings evenly around the center of the latent space. (Smooth Distribution)
 - Penalize the n/w when it tries to cheat by clustering points in specific regions (i.e. by memorizing data)
- ② Therefore, we take the inference over latent variable distribution and constrain it to behave nicely by placing a prior on latent distribution, which is an initial hypothesis or guess about what the latent variable will look like.
- ③ This helps the n/w to enforce a latent space that roughly follows the prior latent distribution.

Problem: We cannot backpropagate gradients through sampling layers!

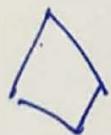
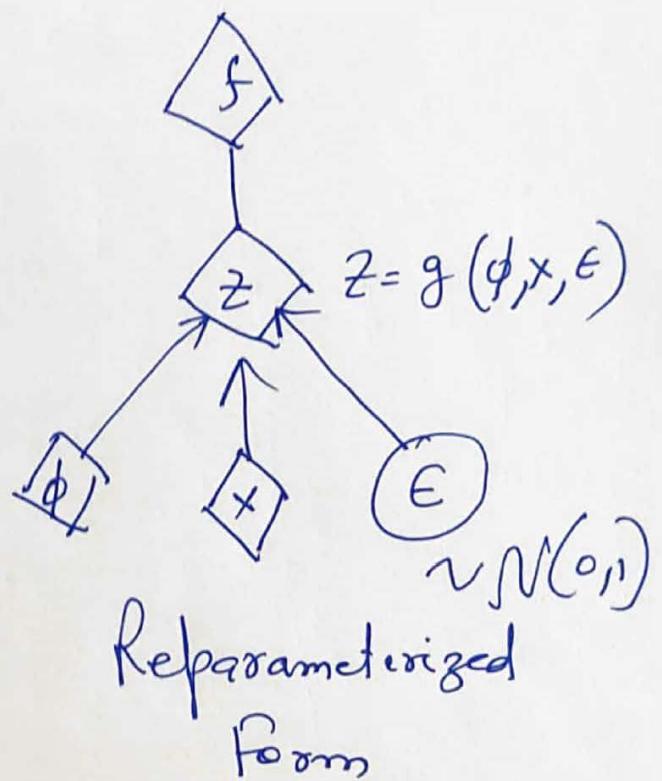
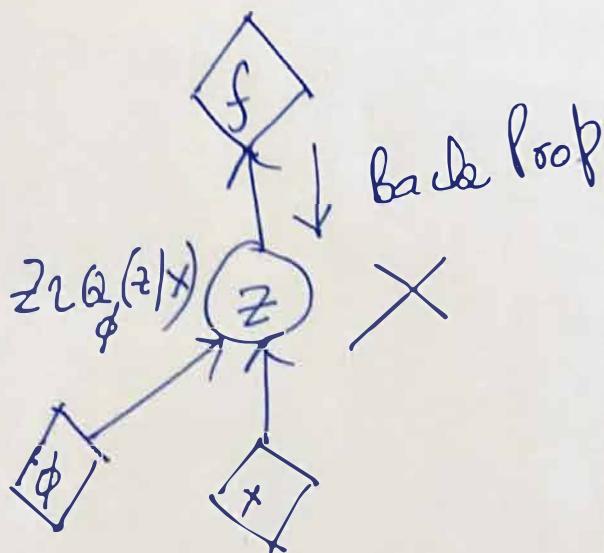
$$z \sim \mathcal{N}(H, \sigma^2)$$

∴ Consider the sampled latent vector z as a sum of

- a fixed vector H ; - and a fixed vector σ , scaled by random constants drawn from prior distribution

$$\Rightarrow z = H + \sigma \cdot \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, 1)$$

Reparameterize the Sampling layer:-



Deterministic Node



Stochastic Node.