



# AUTO-ENCODER → VAE → β-VAE

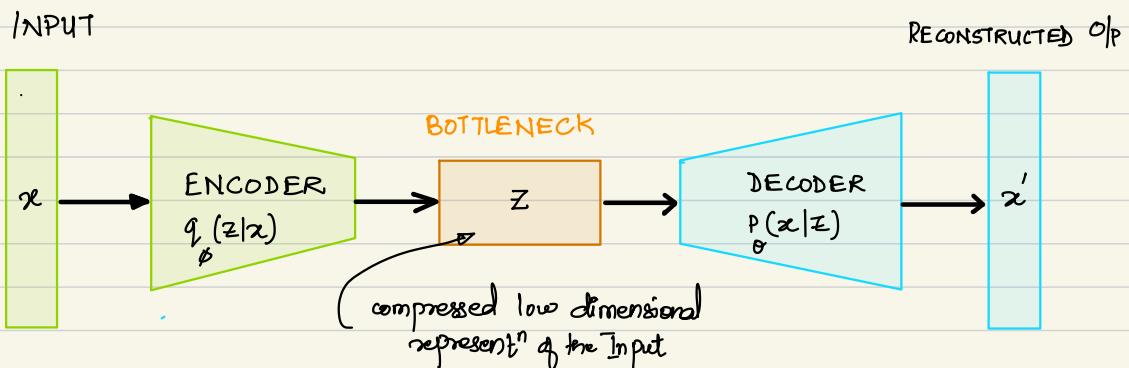
## Contents:

- ① Auto Encoders
- ② De-noising Auto Encoders
- ③ Sparse Auto-Encoders
- ④ Contractive Auto-Encoders
- ⑤ Variational Auto Encoders
  - (a) Loss function : Evidence Lower Bound
  - (b) Reparameterization Trick
- ⑥ Beta - VAE
- ⑦ VB-VAE and VB-VAE 2
- ⑧ TD-VAE
- ⑨ References

① AUTO ENCODER : is a neural network designed to learn an identity function in an unsupervised way to reconstruct the original input while compressing the data in the process so as to discover a more efficient & compressed representation. The original idea dates back to 1980s and was promoted by Geoffrey Hinton & Salakhutdinov in 2006.

It consists of 2 networks:

1. Encoder : which translates the original high dimension input into latent low dimensional code
2. Decoder : receives the code as input & recovers the data from the code with inverse transform.



The encoder essentially accomplishes Dimensionality Reduction, like PCA or Matrix Factorization. In addition, the auto encoder is explicitly optimized for data reconstruction from the code. A good intermediate representation not only can capture latent variables, but also benefits a full de-compression process.

The model contains an encoder  $q_\phi(z|x)$  parametrized by  $\phi$  & a decoder function  $p_\theta(x|z)$  parametrized by  $\theta$ . The low-dimensional code learned from the input  $x$  in the bottleneck layer is  $z = q_\phi(z|x)$  and the reconstructed input is  $x' = p_\theta(x|z)$ .

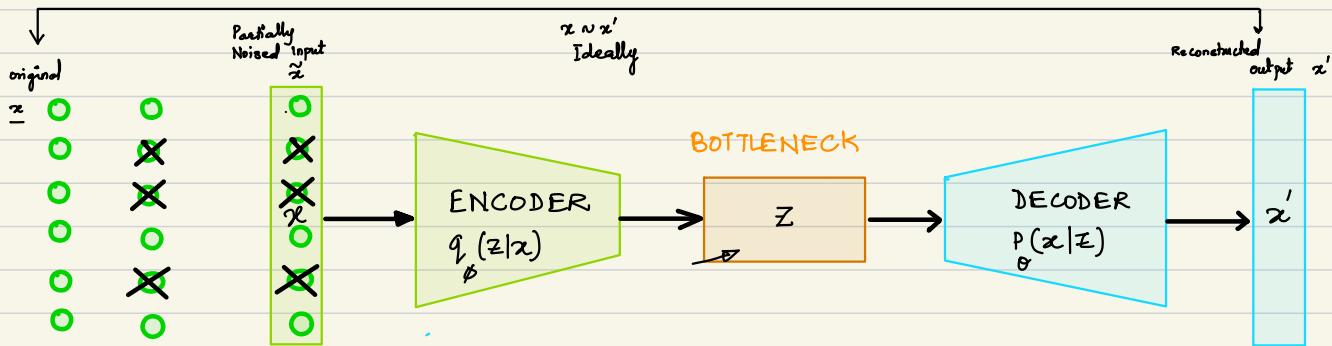
The parameters  $(\phi, \theta)$  are LEARNED together to o/p a reconstructed data sample, same as the original input, or in other words to learn the identity function. There are various METRICS to quantify the difference between the two vectors, such as CROSS ENTROPY when the activation function is SIGMOID or simple MEAN SQUARE ERROR loss.

$$\text{Loss}_{AE}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n \left( x^{(i)} - q_\phi(p_\theta(x^{(i)})) \right)^2$$

② **DENOISING AUTO-ENCODER** : when  $q_\phi$  &  $p_\theta$  are Deep Neural Networks, with more parameters than the data points, then Auto-encoders face the risk of OVERFITTING. To avoid overfitting and Improve Robustness, Denoising Auto-encoders, 2008, proposed a modification. The input is partially corrupted either by adding NOISE or by masking some of the values by a stochastic process.  $\tilde{x} \sim M_\phi(\tilde{x}|x)$ . The model is trained to recover the original input  $x$  (not the corrupt one  $\tilde{x}$ )

$\tilde{x} \sim M_\phi(\tilde{x}|x)$  = defines the mapping from true data samples to corrupt.

$$\text{Loss}_{\text{DAE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - p_\theta(q_\phi(\tilde{x}^{(i)})))^2$$



This design is motivated by the fact that humans can easily recognize an object or a scene even if the view is partially occluded or corrupted. To repair the partially destroyed input, the denoising AE has to discover & capture the relationship between dimensions of the input in order to infer the missing pieces. The noise is controlled by stochastic mapping  $M_\phi(\tilde{x}|x)$  and it is not specific to a particular type of corruption process (masking noise, Gaussian noise, salt & pepper noise etc). Naturally, the corruption process can be equipped with PRIOR knowledge.

In the original experiment of DENOISING, a fixed proportion of the input dimensions are selected at random and their values are forced to zero. This is similar to dropout.

③ **SPARSE AUTOENCODERS** : applies a "SPARSE" constraint on the hidden unit activation to avoid OVERFITTING & improve ROBUSTNESS. It forces the model to only have a small no. of hidden units being active at the same time. A neuron is activated when the value is close to 1 & inactive when the value is close to zero.

Let us assume that there are  $s_1$  neurons in the  $l^{th}$  hidden layer & the activation for the  $j^{th}$  neuron in this layer is labelled as  $a_j^{(l)}(x)$ ,  $j=1, 2, \dots, s_1$ . The fraction of the activation of this neuron  $\hat{s}_j$  is expected to be a small number  $s$ , known as the SPARSITY PARAMETER (normally  $s=0.05$ )

$$\hat{s}_j^{(l)} = \frac{1}{n} \sum_{i=1}^n [a_j^{(l)}(x^{(i)})] \approx s$$

This constraint is achieved by adding a penalty term into the loss function. The KL-divergence term measures the difference between two Bernoulli distributions, one with mean  $s$  & the other with the mean  $\hat{s}_j^{(l)}$ . The hyperparameter ' $\beta$ ' controls how strong the penalty we want to impose on the sparsity loss.

$$\begin{aligned} L_{\text{SAE}}(\theta) &= L(\theta) + \beta \sum_{l=1}^L \sum_{j=1}^{s_l} D_{\text{KL}}(s || \hat{s}_j^{(l)}) \\ &= L(\theta) + \beta \sum_{l=1}^L \sum_{j=1}^{s_l} s \log \frac{s}{\hat{s}_j^{(l)}} + (1-s) \log \left[ \frac{(1-s)}{1-\hat{s}_j^{(l)}} \right] \end{aligned}$$

**K-SPARSE ENCODER**:- The sparsity is enforced by only keeping the top k-highest activations in the bottleneck layer with Linear activation function. First we run a feed forward through the encoder to get the compressed code:  $z = q_\phi(x)$ . Set the values in the code vector only the k-largest are kept if all the other neurons are set to 0. This can be done by a ReLU with an adjustable threshold. This operation leads to a SPARSIFIED CODE  $z' = \text{sparsify}(z)$ . Compute the output and the LOSS from the sparsified( $z'$ ). Note that the BACKPROJECTION only goes through the top 'k' activated hidden units!

- ④ **CONTRACTIVE AUTOENCODER** : encourages the learned representation to stay in a contractive space for better robustness. It adds a term in the Loss Function to penalize the representation being too sensitive to the input, and thus improve the robustness to small perturbations around the training data points. The sensitivity is measured by the FROBENIUS norm of the JACOBIAN matrix of the encoder activations w.r.t the inputs :

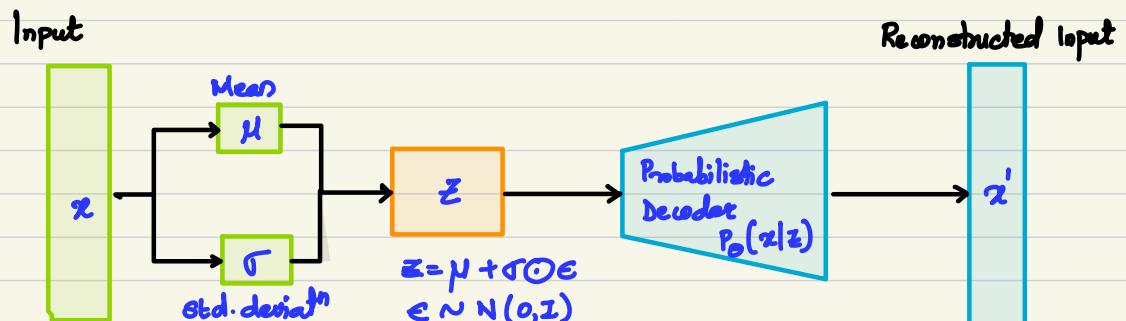
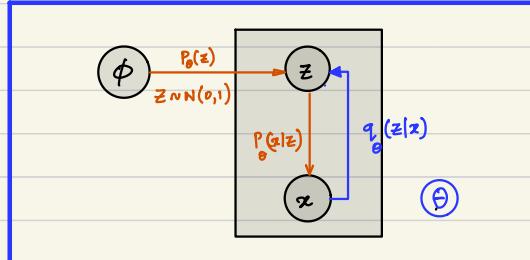
$$\left\| J_f(x) \right\|_F^2 = \sum_{ij} \left( \frac{\partial h_j(x)}{\partial x_i} \right)^2$$

where  $h_j$  is one unit output in the compressed code  $z = f(x)$ . This penalty term is the sum of all the squares of all partial derivatives of the learned encodings w.r.t the input dimensions.

- ⑤ **VAE: VARIATIONAL AUTOENCODERS** (Kingma & Welling, 2014) is deeply rooted in the methods of Variational Bayesian graphical models. Instead of mapping the input into a FIXED vector, we want to map it into a distribution  $p_\theta(z|x)$  parametrized by  $\theta$ .

- Prior  $q_\phi(z)$
- Likelihood  $p_\theta(x|z)$
- Posterior :  $q_\phi(z|x)$

### GRAPHICAL REPRESENTATION



$$\text{Loss}(\theta, \phi) = - \mathbb{E}_{z \sim q_\phi(z|x)} \left[ P_\theta(x|z) \right] + D_{KL} (q_\phi(z|x) || P_\theta(x))$$

⑥ BETA-VAE: If each variable in the inferred latent representation  $\mathbf{z}$  is only sensitive to one single generative factor of relative invariant to other factors, we will say this representation is **DISENTANGLED** or **FACTORIZED**. One benefit that often comes with disentangled representation is good interpretability.

$\beta$ -VAE (HIGGINS, 2017) is a modification of Variational AUTOENCODERS with a special emphasis to discover disentangled latent factors. As in VAE, we want to maximize the probability of generating real data, while keeping the distance between the real and the estimated posterior distributions small.

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

$$\text{subject to } D_{\text{KL}}(q_{\theta}(z|x) || p_{\theta}(z)) < \delta \quad \text{where } \delta \text{ is a small constant.}$$



## PART - I : Basics of KULLBACK - LEIBLER DIVERGENCE

$$1. D_{\text{KL}}(P || Q) = \sum_x P(x) \cdot \log \left[ \frac{P(x)}{Q(x)} \right] = \sum_x P(x) \cdot \log \left[ \frac{P(x)}{Q(x)} \right]$$

2. NOTE:  $D_{\text{KL}}(P || Q) \neq D_{\text{KL}}(Q || P)$  ie KL divergence is NOT symmetric

3. Also, KL divergence is always positive.

4. Let  $p(x)$  &  $q(x)$  be MULTI-VARIATE Gaussians :

$$p(x) \sim N(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \quad \boldsymbol{\mu}_1 \text{ & } \boldsymbol{\mu}_2 \text{ are means}$$

$$q(x) \sim N(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \quad \boldsymbol{\Sigma}_1 \text{ & } \boldsymbol{\Sigma}_2 \text{ are covariances.}$$

By defn of a multi-variate Gaussian the p.d.f is given by

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

5. let  $p(x)$  and  $q(x)$  have the same dimensions ' $k$ ' then prove that

$$D_{\text{KL}}(P(x) || Q(x)) = \frac{1}{2} \left[ \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} - k + \text{trace}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \right]$$

GOAL

PROOF :

$$D_{\text{KL}}(P(x) || Q(x)) = \sum_x P(x) \cdot \log \left[ \frac{P(x)}{Q(x)} \right] \quad \text{--- ①}$$

where,

$$P(x) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}_1|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right\}$$

$$\Rightarrow \log p(x) = -\frac{\kappa}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \quad \textcircled{2}$$

similarly  $\log q(x) = -\frac{\kappa}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \quad \textcircled{3}$

rearranging eqn ① as  $KL(p(x) || q(x)) = \sum_x p(x) [\log p(x) - \log q(x)] \quad \textcircled{1}$

substituting ② & ③ in ①  $\Rightarrow$  we get

$$KL(p(x) || q(x)) = \sum_x p(x) \left\{ -\frac{\kappa}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{\kappa}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right\}$$

(B)

$$KL(p(x) || q(x)) = \sum_x p(x) \left\{ \frac{1}{2} \log \left( \frac{\Sigma_2}{\Sigma_1} \right) + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right\} \quad \textcircled{4}$$

(A)

Consider the term (A)

$$\underbrace{\sum_x p(x) \cdot \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}_{(A)} = \mathbb{E}_{p(x)} \left[ \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right]$$

standard formulae:

- if  $X$  is a scalar, then  $\mathbb{E}[X] = \mathbb{E}[\text{trace}(X)]$

$$\text{trace}(AB) = \text{trace}(BA)$$

$$\text{trace}(ABC) = \text{trace}(BCA) = \text{trace}(CAB)$$

$$\text{trace}(ABC) \neq \text{trace}(ACB)$$

$$\mathbb{E}[\text{trace}(X)] = \text{trace} \mathbb{E}[X]$$

- Now consider,  $\mathbb{E}[X^T A X] = \mathbb{E}[\text{trace}(X^T A X)]$   
 $= \mathbb{E}[\text{trace}(A X X^T)]$   
 $= \text{trace} \mathbb{E}[A X X^T]$

$\therefore$  the term (A)

$$\sum_x p(x) \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = \frac{1}{2} \mathbb{E}_{p(x)} \left[ (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right]$$

(A)

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

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

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

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

*being constant it  
can be taken outside the  $\mathbb{E}_{p(x)}$*

*By defn  $= \Sigma_1$*

$\therefore \textcircled{A} = \text{trace} \left\{ \Sigma_1 \cdot \frac{1}{2} \Sigma_1^{-1} \right\} = \text{trace} \{ I_K \} = K$

Now consider the term ③,

$$\sum_{\mathbf{x}} P(\mathbf{x}) \cdot \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) =$$

Add and subtract  $\boldsymbol{\mu}_1$  in both the brackets

$$\begin{aligned} &= \sum_{\mathbf{x}} P(\mathbf{x}) \left\{ \frac{1}{2} \left[ (\mathbf{x} - \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right]^T \sum_2^{-1} \left[ (\mathbf{x} - \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right] \right\} \\ &= \sum_{\mathbf{x}} P(\mathbf{x}) \left\{ \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \cancel{\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)} + \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right\} \end{aligned}$$

$$= \mathbb{E}_{P(\mathbf{x})} \left[ \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right]$$

$$\begin{aligned} \textcircled{B} &= \mathbb{E}_{P(\mathbf{x})} \left[ \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right] + \mathbb{E}_{P(\mathbf{x})} \left[ (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right] + \mathbb{E}_{P(\mathbf{x})} \left[ \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right] \\ &\quad \textcircled{I} \qquad \textcircled{II} \qquad \textcircled{III} \\ &\text{Comparing with } \textcircled{A} \text{ above} \qquad \text{ZERO} \\ &= \text{trace} \left\{ \sum_2^{-1} \sum_1 \right\} + 0 + \mathbb{E}_{P(\mathbf{x})} \left[ \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right] \end{aligned}$$

$$\begin{aligned} \textcircled{II} \Rightarrow \mathbb{E}_{P(\mathbf{x})} \left[ (\mathbf{x} - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right] &= \mathbb{E}_{P(\mathbf{x})} \left[ (\mathbf{x} - \boldsymbol{\mu}_1)^T \right] \cdot \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= \left\{ \mathbb{E}_{P(\mathbf{x})} \cancel{[\mathbf{x}]} - \boldsymbol{\mu}_1 \right\}^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_1)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= 0 = \text{ZERO} \end{aligned}$$

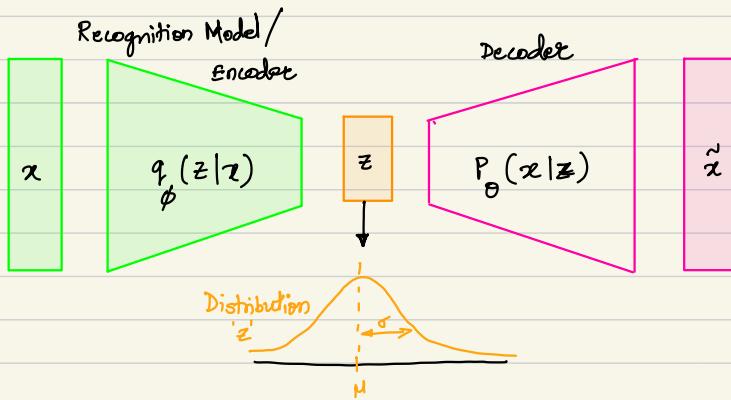
$$\therefore \textcircled{B} = \text{trace} \left[ \frac{\sum_2^{-1} \sum_1}{2} \right] + \frac{1}{2} \left[ (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right]$$

substituting ① & ③ in eqn ④ we get

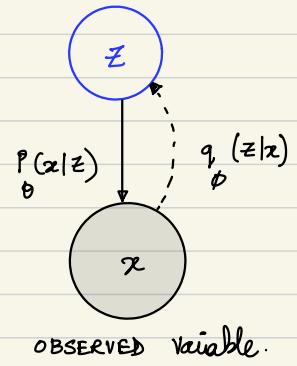
$$\begin{aligned} \text{KL}(P(\mathbf{x}) \parallel Q(\mathbf{x})) &= \sum_{\mathbf{x}} P(\mathbf{x}) \cdot \frac{1}{2} \log \left( \frac{\sum_2}{\sum_1} \right) + \text{trace} \left( \frac{\sum_2^{-1} \sum_1}{2} \right) + \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{k}{2} \end{aligned}$$

$$\therefore \text{KL}(P(\mathbf{x}) \parallel Q(\mathbf{x})) = \frac{1}{2} \left[ \log \left( \frac{\sum_2}{\sum_1} \right) - k + \text{tr} \left( \sum_2^{-1} \sum_1 \right) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \sum_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right]$$

### PART-III : OVERVIEW of VAE.



Latent variable



The goal of VAE is to LEARN the parameters of  $\phi$  of Decoder  $q_\phi(z|x)$ . This way we can sample from  $q_\phi(z|x)$  various values of  $\hat{z}$ . We will also LEARN  $\theta$  in  $P_\theta(x|z)$ . The derivation of the LOSS  $F^*$  of VAE can be solved through VARIATIONAL INFERENCE.

$$\text{Consider } p(z|x) = \frac{p(x|z) \cdot p(z)}{p(x)} \quad \textcircled{1} \quad \text{Bayes' Theorem}$$

where  $p(x) = \text{Denominator of RHS} = \int f(x|z) \cdot p(z) dz$  This integral is often INTRACTABLE.  
 $= \int p(x, z) dz$  can be written as a joint distribution but is not available in the CLOSED form to solve.

Here, the alternative is to approximate the intractable  $f(x|z)$ . POSTERIOR by another distribution  $q_\phi(z|x)$  which has a tractable solution. Variational Inference converts the problem into that of optimisation. Thus we try to find that  $q_\phi(z|x)$  using K-L divergence and minimize it.

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

By definition of expectation w.r.t  $Q_\phi(z|x)$

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

$$= \mathbb{E}_{z \sim Q_\phi} \left[ \log (Q_\phi(z|x)) - \underbrace{\log (P_\theta(z|x))}_{\text{lit. circled 1}} \right] \quad \textcircled{2}$$

$$= \mathbb{E}_{z \sim Q_\phi} \left[ \log (Q_\phi(z|x)) - \log P_\theta(z|x) - \log P_\theta(z) + \underbrace{\log P_\theta(z)}_{\text{constant wrt } z} \right]$$

$$= \mathbb{E}_{z \sim Q_\phi} \left[ \log (Q_\phi(z|x)) - \log P_\theta(z|x) - \log P_\theta(z) \right] + \log (P_\theta(z)) \quad \because \mathbb{E}[\text{const}] = \text{const.}$$

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

$$\underbrace{\log P_\theta(z) - D_{KL}(Q_\phi(z|x) || P_\theta(z|x))}_{-\text{Loss function}} = \mathbb{E}_{z \sim Q_\phi} [\log P_\theta(z|x)] - \mathbb{E}_{z \sim Q_\phi} [\log Q_\phi(z|x) - \log P_\theta(z)]$$

Note: KL Divergence Form

$$\log P_{\theta}(x) - D_{KL} \left( Q_{\phi}(z|x) \parallel P_{\theta}(z|x) \right) = \underbrace{\mathbb{E}_{z \sim Q_{\phi}} \left[ \log P_{\theta}(x|z) \right]}_{\text{Term (A)}} - \underbrace{D_{KL} \left( Q_{\phi}(z|x) \parallel P_{\theta}(z) \right)}_{\text{Term (B)}}$$

This is the  
OBJECTIVE function  
of VAE.

Term A =  $E_{\mathcal{E}} \left[ \log p_{\theta}(\mathcal{E}|x) \right]$  is the reconstruction likelihood.

Term (B) =  $D_{KL} \left[ Q_\phi(z|x) \parallel P_\theta(z) \right]$  ensures that our LEARNED distribution  $Q$  is similar to Prior  $P$ .

$$\therefore \text{LOSS fn} = -\mathbb{E}_{z \sim Q_\phi} [\log (P_\theta(z|x))] + D_{KL}[Q_\phi(z|x) || P_\theta(z)]$$

Reconstruction loss.
Regularizer

So our target is to LEARN the parameters ' $\phi$ ' & ' $\theta$ '

## PART IV: Loss f" TERMS in depth.

The regularizer  $D_{KL} \left[ Q_\phi(z|x) \parallel P_\theta(z) \right]$  is the divergence between two distributions. We assign  $P_\theta(z) \sim N(0, I)$ .  
 $\therefore$  we compare the term of the Loss fn to the k-l divergence b/w two Gaussian's done previously.  $\therefore N_2=0; I_2=I$

$$\text{Regularized } D_{KL} \left[ Q_\phi(z|x) || P_\theta(z) \right] = \frac{1}{2} \left[ \log \frac{1}{|\Sigma_\phi(x)|} - K + \underbrace{\text{trace} \left[ \Sigma_\phi(x) \right]}_{\substack{\text{summing over } k \text{ elements}}} + \underbrace{M_\phi(x)^T \cdot M_\phi(x)}_{\substack{\text{scalar quantity can be squared}}} \right]$$

- where 'k' is the dimension of the Gaussian.
  - also  $|\Sigma_g(x)|$  is the determinant of the diagonal matrix = product of its diagonal elements.

$$= \frac{1}{2} \left[ \sum_k \sum_{\phi} (\chi) + \sum_k N_{\phi}^2 (\chi) - \sum_k 1 - \log \prod_{k, \phi} \Sigma (\chi) \right]$$

$$= \frac{1}{2} \left[ \sum_k \sum_p (z_k) + \sum_k \mu_p^2(z_k) - \sum_k 1 - \sum_k \log (\sum_p (z_k)) \right]$$

$$\text{D}_{\text{KL}} \left[ N(\boldsymbol{\mu}_\phi(\mathbf{x}), \Sigma_\phi(\mathbf{x})) \parallel N(0, I) \right] = \frac{1}{2} \sum_k \left[ \Sigma_\phi(\mathbf{x}) + \boldsymbol{\mu}_\phi^2(\mathbf{x}) - 1 - \log \Sigma_\phi(\mathbf{x}) \right]$$

for numerical stability  $\sum_{\phi}(x)$  is replaced by  $\exp[\sum(x)]$

$$= \frac{1}{2} \sum_i \left[ \exp(\Sigma(x)) + N_{\phi}^2(x) - 1 - \sum_{\phi}(x) \right]$$

## OPTIMIZATION of THE LOSS FUNCTION

$-L(\theta, \phi) = \log P_\theta(x) - D_{KL} [Q_\phi(z|x) || P_\theta(z|x)]$  is referred to as VARIATIONAL LOWER BOUND or EVIDENCE LOWER BOUND (ELBO).

$$\theta^*, \phi^* = \underset{\theta, \phi}{\operatorname{argmax}} \{ L(\theta, \phi) \}$$

We know that  $D_{KL} [Q_\phi(z|x) || P_\theta(z|x)]$  is always positive  $\geq 0 \therefore L(\theta, \phi) \leq \log P_\theta(x)$

Thus minimizing the Loss fn we are actually maximizing the probability of generating real data samples.

$$L(\theta, \phi) = - \mathbb{E}_{z \sim Q_\phi(z|x)} [\log P_\theta(z|x)] + \frac{1}{2} \sum_k \left[ \exp(\sum_\theta z_k) + \frac{\mu_\phi^2(x)}{\phi} - 1 - \sum_\theta z_k \right]$$

During optimization, we apply the alternate optimization technique. i.e.

$$\begin{aligned} \text{Iteration #1: } & \left\{ \begin{array}{l} \theta^* = \nabla_\theta \{ L(\theta, \phi) \} \\ \phi^* = \nabla_\phi \{ L(\theta^*, \phi) \} \end{array} \right. \quad \phi = \text{const.} \\ & \phi^* \text{ is kept constant.} \end{aligned}$$

$$\begin{aligned} \text{Iteration #2: } & \left\{ \begin{array}{l} \theta^{**} = \nabla_\theta \{ L(\theta^*, \phi^*) \} \\ \phi^{**} = \nabla_\phi \{ L(\theta^{**}, \phi^*) \} \end{array} \right. \quad \theta^* \text{ is kept constant.} \\ & \theta^{**} \text{ is kept const.} \end{aligned}$$

NOTE: It is observed that differentiating wrt  $\theta$  is straight forward, but derivative wrt  $\phi$  is not possible. To overcome this issue REPARAMETRIZATION trick is used.

Optimization is carried out wrt both,  $\theta$  &  $\phi$ , to LEARN  $Q_\phi(z|x)$  &  $P_\theta(z|x)$  simultaneously

Derivative wrt  $\theta$ : ie  $\frac{\partial L}{\partial \theta_i} = \nabla_\theta \{ L(\theta, \phi) \}$

$$= \nabla_\theta \left\{ - \mathbb{E}_{z \sim Q_\phi(z|x)} [\log P_\theta(z|x)] + \frac{1}{2} \sum_k \underbrace{\left[ \exp(\sum_\theta z_k) + \frac{\mu_\phi^2(x)}{\phi} - 1 - \sum_\theta z_k \right]} \right\}$$

Term being constant wrt  $\theta$ ; its derivative = ZERO

$$= -\nabla_\theta \mathbb{E}_{z \sim Q_\phi(z|x)} [\log P_\theta(z|x)] \quad \text{we can easily move } \nabla_\theta \text{ inside } \mathbb{E}_z \text{ sign. Also } z \sim Q_\phi(z|x)$$

$$= \frac{1}{L} \sum_{l=1}^L \nabla_\theta \left[ \log P_\theta(z_l|x) \right] \quad \text{NOTE: We replace the Expectation by the average value according to the MONTE CARLO concept.}$$

Derivative wrt  $\phi$ : Here we will consider the optimal value of  $\theta^*$  calculated above

$$\therefore \hat{\phi}_i = \nabla_\phi \{ L(\theta^*, \phi) \}$$

$$= \nabla_\phi \left\{ - \mathbb{E}_{z \sim Q_\phi(z|x)} [\log P_\theta(z|x)] + \frac{1}{2} \sum_k \left[ \exp(\sum_\theta z_k) + \frac{\mu_\phi^2(x)}{\phi} - 1 - \sum_\theta z_k \right] \right\}$$

This derivation  $\nabla_\phi$  is difficult to estimate because ' $\phi$ ' appears in the distribution wrt which the expectation is taken

NOTE:

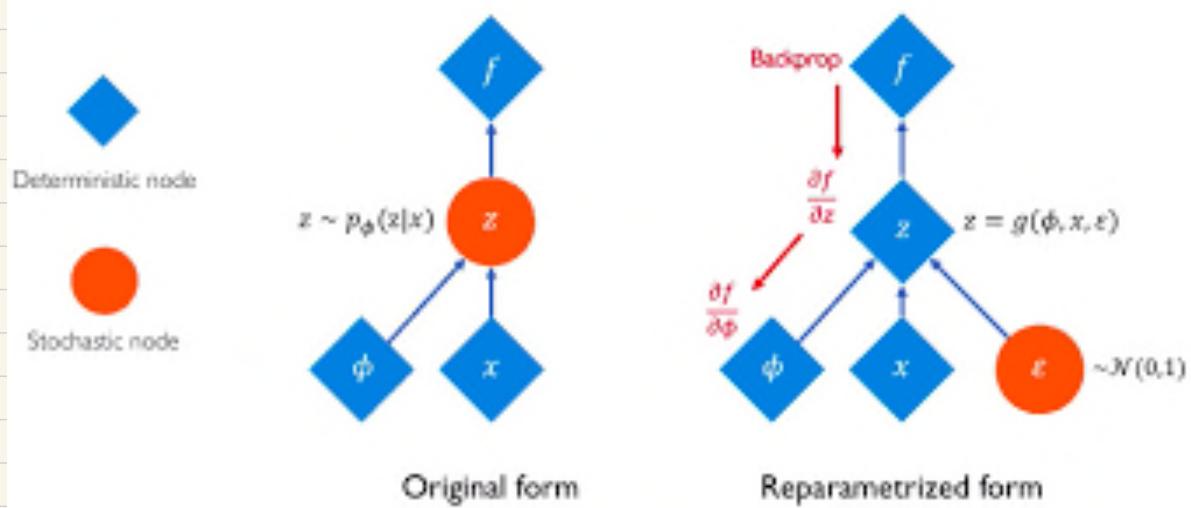
$$\nabla_\phi \mathbb{E}_{Q_\phi(z|x)} [f(z)] \neq \mathbb{E}_{Q_\phi(z|x)} [\nabla_\phi f(z)]$$

If we can somehow write this expectation in such a way that  $\phi$  appears inside the expectation, then we can push the gradient INSIDE the Expectation; ie if we can write

$$\mathbb{E}_{\theta_{\phi}(z|x)}[f(z)] = \mathbb{E}_{P(\epsilon)}[f(g_{\phi}(\epsilon, z))] \text{ such that } z = g_{\phi}(\epsilon, x) \xrightarrow{\text{any LINEAR TRANSFORMATION}} \text{with } \epsilon \sim N(0, 1)$$

$$\text{In our case, } g_{\phi}(\epsilon, x) = \mu(x) + \epsilon \Theta \sum_{\phi}^k (x)$$

$$= z \sim N(\mu(x), \Sigma(x))$$



$$\therefore \hat{\phi} = \nabla_{\phi} \{ L(\theta^*, \phi) \}$$

$$= \nabla_{\phi} \left\{ -\mathbb{E}_{\substack{z \sim p(z|x) \\ \epsilon \sim N(0,1)}} \left[ \log P_{\theta}(x|z^{(1)}) \right] + \frac{1}{2} \sum_k \left[ \exp \left( \sum_{\phi}^k (x) \right) + \mu_{\phi}^2(x) - 1 - \sum_{\phi}^k (x) \right] \right\}$$

NOTE:  $\phi$  is not present here  
 $\therefore \nabla_{\phi}$  can be now moved inside

$$= -\mathbb{E}_{\substack{z \sim p(z|x) \\ \epsilon \sim N(0,1)}} \left[ \nabla_{\phi} (\log P_{\theta}(x|z^{(1)})) \right] + \nabla_{\phi} \left[ \frac{1}{2} \sum_k \left( \exp \left( \sum_{\phi}^k (x) \right) + \mu_{\phi}^2(x) - 1 - \sum_{\phi}^k (x) \right) \right]$$

Replace  $E$  by mean by MONTE CARLO

$$= -\frac{1}{S} \sum_{s=1}^S \nabla \left[ \log P_{\theta}(x|z^{(s)}) \right] + \nabla_{\phi} \left[ \frac{1}{2} \sum_k \left( \exp \left( \sum_{\phi}^k (x) \right) + \mu_{\phi}^2(x) - 1 - \sum_{\phi}^k (x) \right) \right]$$