

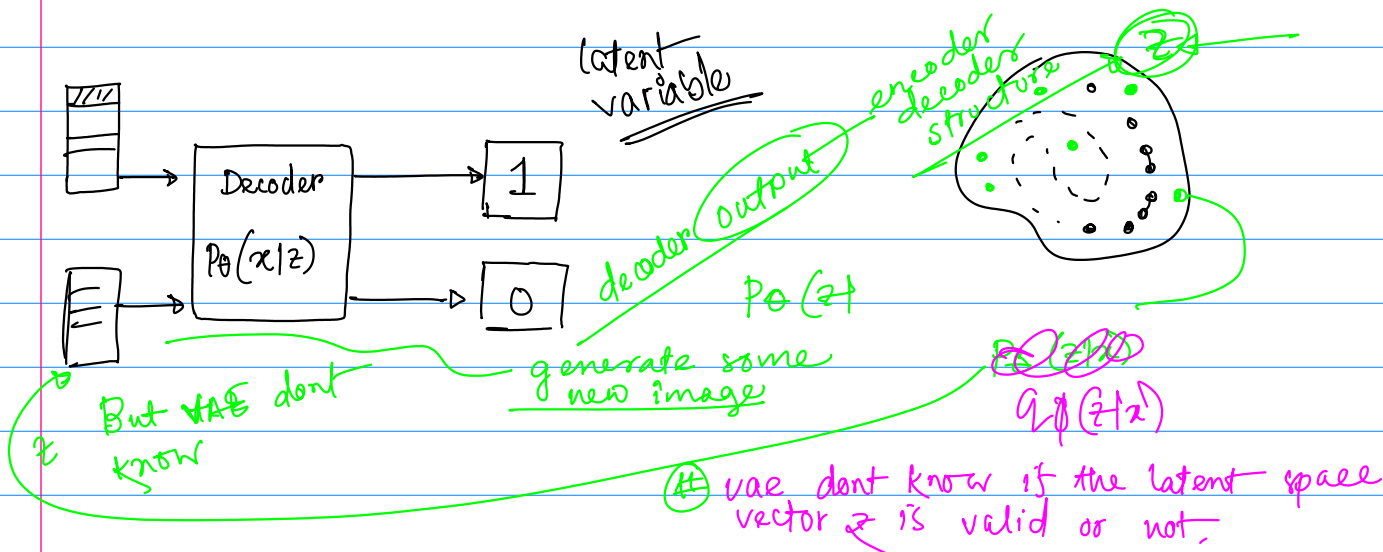
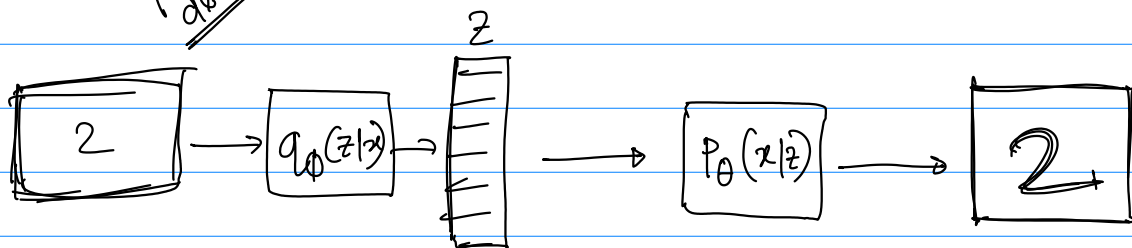
Loss function: $L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} [p_\theta(x|z)] + D_{KL}(q_\phi(z|x) || p_\theta(z))$

The KL divergence term is shown as $D_{KL}(q_\phi || p_\theta)$.

Goal of Variational Autoencoder:

To find $q_\phi(z|x)$ of some latent variable which we can sample or generate from $q_\phi(z|x)$ in order to generate new samples $x' \sim p_\theta(x|z)$.

not present in the input



it forces the encoder to produce a probability distribution function and then try to make them further

Probability Concepts Required:

$p(x)$: prob of a random var x

$p(x|y)$: prob of x given y has happened.

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

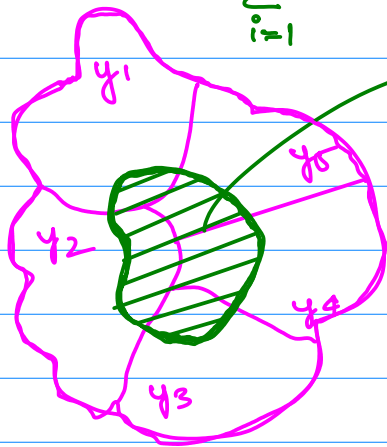
$$\text{posterior probability} = \frac{p(x, y)}{p(x)}$$

$\frac{p(y)}{p(x)}$ → likelihood ratio
 $p(y)$ → prior
 $p(x, y)$ → joint prob dist

∴ Total probability theorem: $y_1, y_2, y_3, \dots, y_n$ (no overlap, mutually exclusive ($y_i \cap y_j = \emptyset$))

Event x , union of multiple, mutually exclusive events:

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



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

$$= \sum_{i=1}^5 P(x|y_i) P(y_i) \quad \text{--- (2)}$$

$$p(y|x) = \frac{P(x|y) P(y)}{\left[\sum_{i=1}^N P(x|y_i) P(y_i) \right] \rightarrow P(x)}$$

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

$$= \frac{P(x|y) P(y)}{\sum_{i=1}^N P(x|y_i) P(y_i)}$$

Expectation of a Random Variable:

$$E(x) = \sum_{i=1}^K x_i \cdot P(x=x_i) = E(x)$$

③ $P(3) = \frac{1}{6}$
 $P(3) = \frac{1}{6}$

$p(x=3 | y \text{ is odd})$

$$P(x|y) = \frac{P(y|x) P(y)}{P(x)}$$

For fair die: $\{1, 2, 3, 4, 5, 6\}$
 $P = \frac{1}{6}$

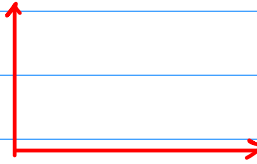
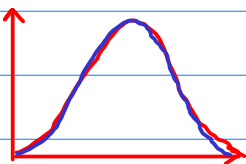
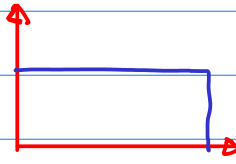
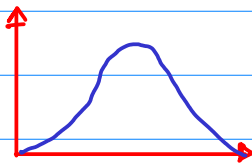
$$E(x) = \sum_{i=1}^6 x_i \cdot P(x_i)$$

$$= \frac{1}{6} + \frac{2}{6} + \frac{3}{6} + \frac{4}{6} + \frac{5}{6} + \frac{6}{6}$$

high value $= \left(\frac{21}{6}\right) =$
 $= \frac{21}{6} =$

KL-Divergence:

Kullback Divergence: 2022



low value

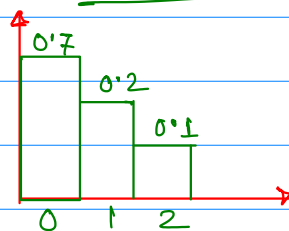
Measure of how one probability is different from the other: for two probability distribution P and Q

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

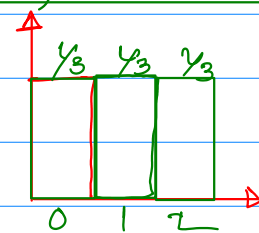
$$= \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad D_{KL}$$

Example:

P distribution



Q distribution:



uniform distribution

Here $D_{KL}(Q||P) = \sum_{x=0}^2 Q(x) \log \frac{Q(x)}{P(x)}$

$$= \frac{1}{3} \log\left(\frac{\frac{1}{3}}{0.7}\right) + \frac{1}{3} \log\left(\frac{\frac{1}{3}}{0.2}\right) + \frac{1}{3} \log\left(\frac{\frac{1}{3}}{0.1}\right)$$

$$= 0.09673 \text{ nats.}$$

Properties: ① $KL(P||Q) \text{ or } KL(Q||P) \geq 0$

② $KL(P||Q) \neq KL(Q||P)$

Non symmetric

Consider $P(x) = \mathcal{N}(x; \mu_1, \Sigma_1)$ Σ_1 and Σ_2
 $Q(x) = \mathcal{N}(x; \mu_2, \Sigma_2)$ covariance matrix

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

Equation of normal distribution:

then $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

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

Proof:

$$KL(P(x) \parallel Q(x)) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$$= \sum_x P(x) \cdot [\log P(x) - \log Q(x)] \quad \text{--- (1)}$$

and we put the equation of the distribution

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

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

$$= -\log(2\pi)^{K/2} - \log(|\Sigma_1|)^{1/2} - \frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)$$

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

Similarly $\log Q(x) = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{1}{2} (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) \quad \text{--- (3)}$

Putting (2) and (3) in (1) we get

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

$$\sum P(x) \left[-\frac{K}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_1|) - \frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) + \frac{K}{2} \log(2\pi) + \frac{1}{2} \log(|\Sigma_2|) + \frac{1}{2} (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) \right]$$

$$= \frac{1}{2} \sum_x P(x) \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} + (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) - (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) \right] \quad \text{--- (4)}$$

Considering partially $\sum P(x) = E(x)$ Expected value of x in P

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

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

Trace: for a square matrix; trace is the sum of the diagonal elements.

From the properties of trace we know

$$\begin{aligned} E(X^T A X) &= E(\text{tr}(X^T A X)) \\ &= E(\text{tr}(A X X^T)) \\ &= \text{tr}(E(A X X^T)) \end{aligned}$$

Therefore following this we can rewrite this like the following

$$\begin{aligned} &\frac{1}{2} E_p \left[(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right] \\ &= E_p \left[\text{tr} \left(\frac{1}{2} (x - \mu_1) (x - \mu_1)^T \Sigma_1^{-1} \right) \right] \\ &= E_p \left[\text{tr} \left(\frac{1}{2} (x - \mu_1) (x - \mu_1)^T \Sigma_1^{-1} \right) \right] \\ &\quad \text{trace of a vector/matrix} \\ &= \text{tr} \left[E_p \left[\frac{1}{2} (x - \mu_1) (x - \mu_1)^T \right] \Sigma_1^{-1} \right] \\ &\quad \text{Covariance matrix} \end{aligned}$$

$$= \text{tr} \left[\Sigma_1 \frac{1}{2} \Sigma_1^{-1} \right]$$

$$= \text{tr} \left[\frac{1}{2} \Sigma_1 \Sigma_1^{-1} \right]$$

$$= \text{tr} \left[\frac{1}{2} I_K \right]$$

$$= \frac{K}{2}$$

Considering the second part from equation (04)

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

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

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

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

Trace and Expectation trick:

$$a) \Rightarrow E(x) = E(\text{tr}(x))$$

Since trace of x is scalar

$$b) \Rightarrow \text{tr}(AB) = \text{tr}(BA)$$

$$c) \Rightarrow \text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

But it follows in cyclic order

$$c) \Rightarrow \text{tr}(ABC) \neq \text{tr}(ACB)$$

$$d) \Rightarrow E(\text{tr}(x)) = \text{tr}(E(x))$$

Expanding the equation we get the following

$$= \underbrace{E_p \left[\frac{1}{2} (\mathbf{x} - \mu_1)^T \Sigma_2^{-1} (\mathbf{x} - \mu_1) \right]}_{\text{tr} \left\{ \frac{1}{2} \Sigma_2^{-1} \Sigma_1 \right\}} + \underbrace{E_p \left[(\mathbf{x} - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]}_{\text{this part is constant due to the absence of } \mathbf{x}} + E_p \left[(\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

$\rightarrow 0$
zero

$$\begin{aligned} & E_p \left[(\mathbf{x} - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \\ &= \left[(E_p(\mathbf{x}) - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \\ &= \left[\underbrace{(\mu_1 - \mu_1)}_0^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \end{aligned}$$

Therefore summing everything up we get

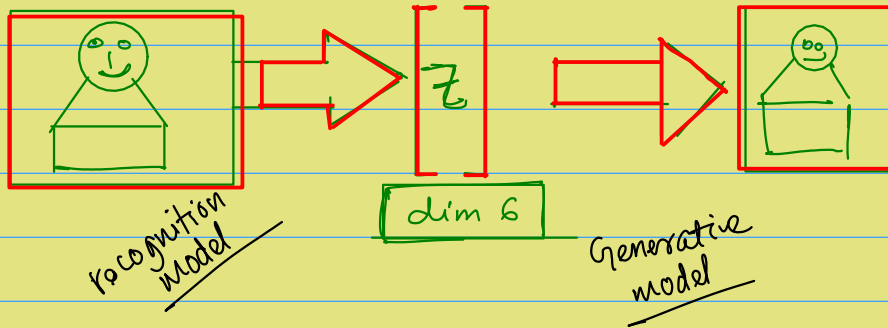
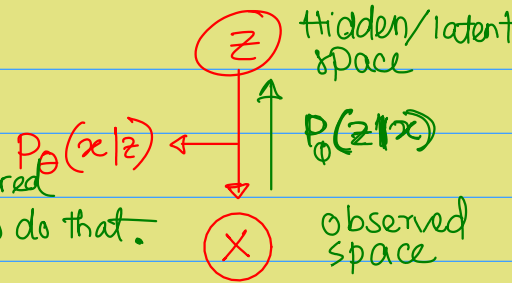
$$KL(P(\mathbf{x}) \parallel Q(\mathbf{x})) = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right]$$

(Proved)

LATENT VARIABLE

$$x \longrightarrow z \longrightarrow x'$$

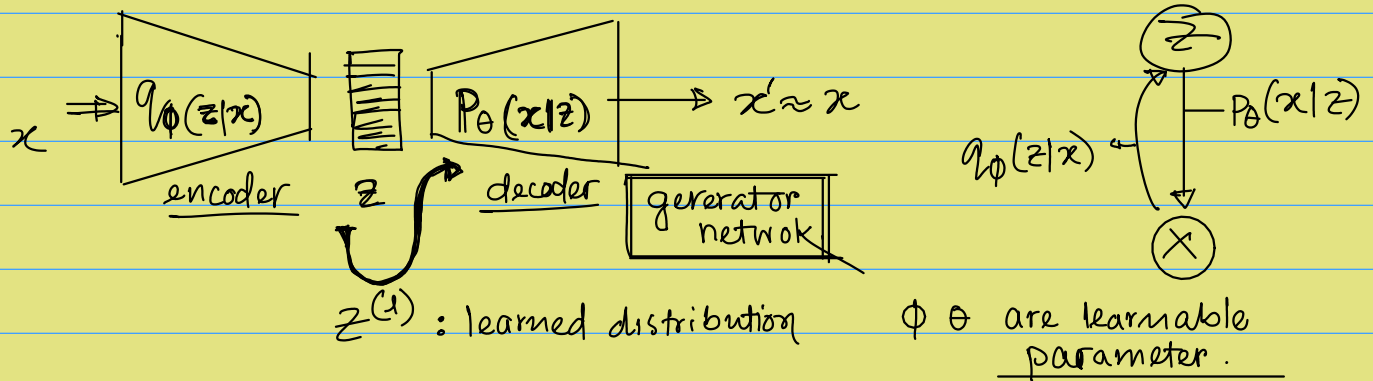
latent variable: some feature that are not measured yet. or technology is not here to do that.



AE: Fixed values
VAE: Gives out probability

DERIVATION OF LOSS FUNCTION

to find the distribution of $q_\phi(z|x)$ of some latent variables which can be sampled from $z \sim q_\phi(z|x)$. to generate the new samples x' from $p_\theta(x|z)$



The problem of approximate inference:

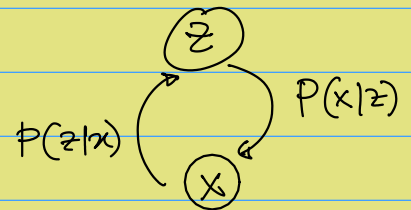
x : set of observed var

z : set of latent variables with joint prob distribution $p(z|x)$.

Aim is to learn the conditional distribution of z given the input set x . ie $p(z|x)$ need to learn :

$$p(z|x) = \frac{p(x|z) p(z)}{p(x)} \quad \text{--- (A) --- equation.}$$

evaluating $p(z|x)$ is difficult because $p(x)$ can not be calculated/solved.



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

Not available in the closed loop form.

Alternative:

approximate $p(z|x)$ by another dist $q_\phi(z|x)$ which is defined in suitable way done by variational inference. (vi) :

Pose inference problem into optimization problem.

By modeling $p(z|x)$ using $q(z|x)$ where $q(z|x)$ has a simple distribution like gaussian.

calculate D_{KL} between $p(z|x)$ and $q(z|x)$

$$D_{KL}(q(z|x) \parallel p(z|x)) = \sum_z q(z|x) \log \frac{q(z|x)}{p(z|x)}$$

$$= E_{z \sim q(z|x)} \left[\log \frac{q(z|x)}{p(z|x)} \right]$$

$$= E_{z \sim q(z|x)} \left[\log(q(z|x)) - \log(p(z|x)) \right]$$

mentor
experience
mentor
(B)

By substituting equation (A) in (B) we get

$$D_{KL}(q(z|x) \parallel p(z|x)) = E_z \left[\log(q(z|x)) - \log \frac{p(x|z) \cdot p(z)}{p(x)} \right]$$

where $z = z \sim q_\phi(z|x)$

$$= E_z \left[\log(q_\phi(z|x)) - \log(p_\theta(x|z)) - \log(p_\theta(z)) + \log(p_\theta(x)) \right]$$

Here since expectation is over z , we can bring out $p(x)$ that does not involve z .

z is being sampled from $q_\phi(z|x)$

$$\therefore D_{KL}(q_\phi(z|x) \parallel p_\theta(z|x)) - \log(p_\theta(x)) = E_z \left[\log q_\phi(z|x) - \log p_\theta(x|z) - \log p_\theta(z) \right]$$

Rearranging the equation again, we get

$$\log p_\theta(x) - D_{KL}[q_\phi(z|x) \parallel p_\theta(z|x)]$$

$$= E_z \left[\log p_\theta(x|z) \right] - E_z \left[\log q_\phi(z|x) - \log p_\theta(z) \right]$$

ক্লসে নাই কেন কথা
ক্লসে নাই কেন কথা

$$= E_z \left[\log p_\theta(x|z) \right] - D_{KL}[q_\phi(z|x) \parallel p_\theta(z)]$$

learned distribution

prior distribution

reconstruction likelihood

VAC objective function:

And the loss function here is simply $\text{loss} = - \text{objective function}$.

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

Therefore the target of the whole process is to find out the optimal θ and ϕ such that

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

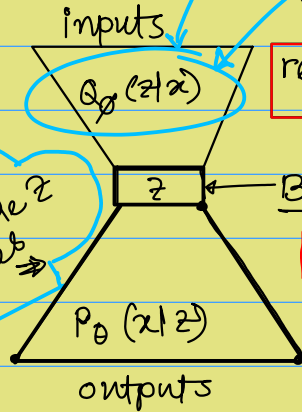
network parameters to learn

$$L(\theta, \phi) = \underbrace{-\mathbb{E}_{z \sim Q_{\phi}(z|x)} [\log P_{\theta}(x|z)]}_{\text{reconstruction loss}} - \underbrace{D_{KL}[Q_{\phi}(z|x) \parallel P(z)]}_{\text{regularizer}}$$

auto regularizer

Log likelihood

Regularizer



recognition model

Bottleneck

generative model

$$-\mathbb{E}_{z \sim Q_{\phi}(z|x)} [\log(P_{\theta}(x|z))] \quad \text{log likelihood recognition model.}$$

$$P_{\theta}(x|z) = \mathcal{N}(M_{\theta}(z), \Sigma_{\theta}(z))$$

We get a square error between the mean of the gaussian distribution and the data sample. !

Here

$$P_{\theta}(x|z) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_{\theta}(z)|}} \exp\left(-\frac{(x - M_{\theta}(z))^T \Sigma_{\theta}(z)^{-1} (x - M_{\theta}(z))}{2}\right)$$

data point

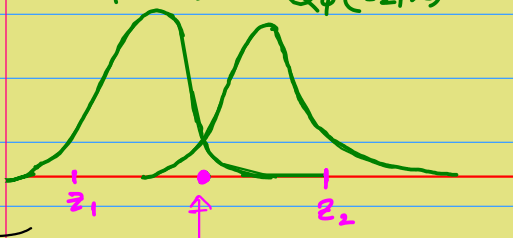
From the above we get

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

New data $\hat{z} = \frac{(z_1 + z_2)}{2}$ which is generated.

Now look at the term:

$Q_{\phi}(z_1|x)$ $Q_{\phi}(z_2|x)$



weighted combination of the variables z_1 and z_2

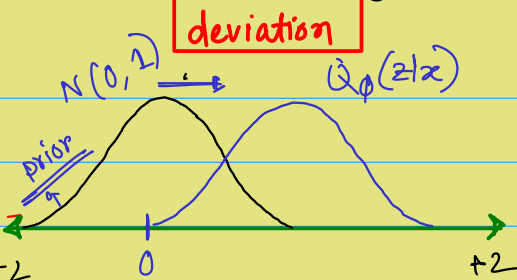
$$D_{KL}[Q_{\phi}(z|x) \parallel P(z)]$$

distribution of z given x

here we are assuming that distribution of the $P_{\theta}(z)$ is a normal distribution with $\mathcal{N}(0, 1)$

We want our data (\hat{z}) to be closer to a normal distribution. That is why divergence with $P_{\theta}(z) = \mathcal{N}(0, 1)$ is also calculated which is also working as a regularizer.

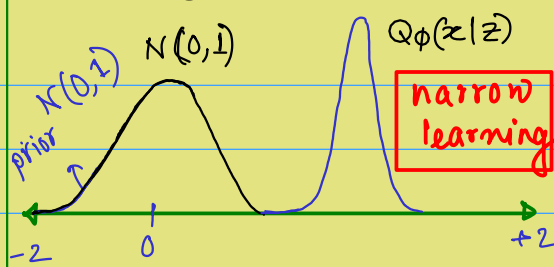
Data Fidelity
with No KL divergence



deviates from the normal distribution. Will be of normal shape but there is deviation.

⊕ No term is there to bring it back to the normal distribution.

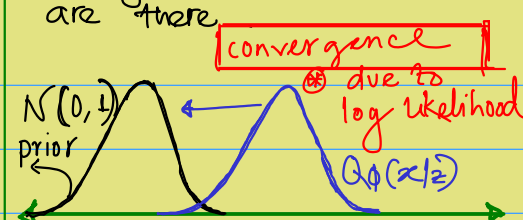
No data Fidelity
But regularization.



⊕ Network will cheat by learning very narrow distribution. / No large distribution.

Normal auto encoder !

Data Fidelity &
Regularization both are there



⊕ shape equal to a $N(0, 1)$ due to the regularized term KL divergence.

Aim and Objective: $D_{KL}(Q_\phi(z|x) \| P_\theta(z)) \rightarrow$ Regularization Term.

easy choice of z with $N(0, 1)$. We want $P_\theta(z)$ as close as possible to the $Q_\phi(z|x)$ so that they can be sampled easily.

Originally derived:

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

Now we put $\mu_2 = 0$ and $\Sigma_2 = I$

and $\mu_1 = \mu_\phi(z)$ and $\Sigma_\phi(z) = \Sigma_1$ we get

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

$$= D_{KL}(N(\mu_\phi(z), \Sigma_\phi(z)) \| N(0, I)) =$$

$$\frac{1}{2} \left[\text{tr}(\Sigma_\phi(z)) + \mu_\phi(z)^T \mu_\phi(z) - K - \log |\Sigma_\phi(z)| \right]$$

= Here K is the dimension of the Gaussian

→ $\text{tr}(\Sigma_\phi(z))$ is the trace function, which is the sum of the diagonals $\Sigma_\phi(z)$.

Then determinant of a diagonal matrix is product of its diagonals.

above equation becomes:

$$D_{KL}(N(\mu_\phi(x), \Sigma_\phi(x)) \| N(0, I))$$

$$= \frac{1}{2} \left(\sum_K \Sigma_\phi(x) + \sum_K \mu_\phi^2(x) - \sum_K 1 - \log \prod_K \Sigma_\phi(x) \right)$$

$$= \frac{1}{2} \left(\sum_K \Sigma_\phi(x) + \sum_K \mu_\phi^2(x) - \sum_K 1 - \sum_K (\log \Sigma_\phi(x)) \right)$$

$$= \frac{1}{2} \sum \left(\Sigma_\phi(x) + \mu_\phi^2(x) - 1 - \log \Sigma_\phi(x) \right)$$

For numerical stability we do the following ~~at~~ this

$$D_{KL}(N(\mu_\phi(x), \Sigma_\phi(x)) \| N(0, I))$$

Model

$$\Sigma_\phi(x) = e^{\hat{\Sigma}_\phi(x)}$$

$$= \frac{1}{2} \sum \left(\hat{\Sigma}_\phi(x) + \mu_\phi^2(x) - 1 - \hat{\Sigma}_\phi(x) \right)_{\text{if}}$$