

(c) Means of metrics from 7 runs of the training are as following.

NELBO : 101.687

K L Divergence : 19.261

Reconstruction loss : 82.425

Std. dev. for the same experiment are:

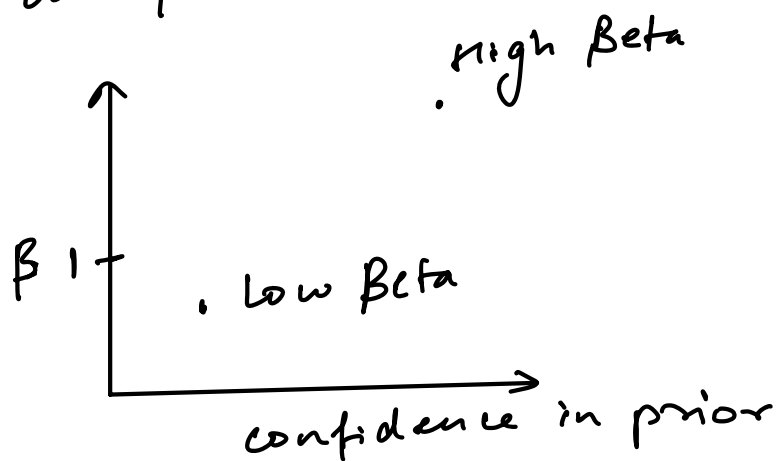
NELBO : 0.674

KL Divergence: 0.143

Reconstruction loss: 0.781

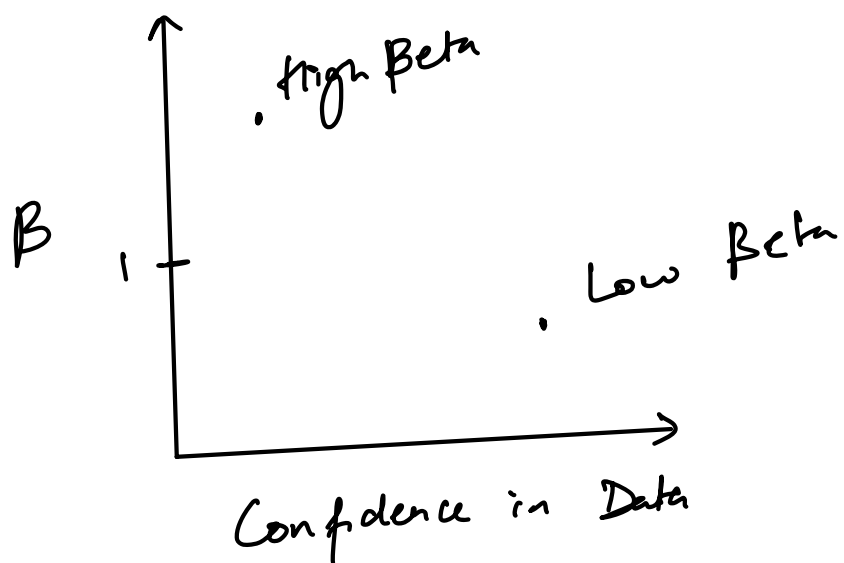
1(d) ' β ' in β -VAE seems to be acting like a weightage to KL Divergence. By tuning this, we are trying to convey how much we care for posterior $q_\phi(z|x)$ to be close to prior $p(z)$ - which is gaussian in our case.

- When we set $\beta \geq 1$, then the ELBO remains same as VAE.
- When we set $\beta < 1$, we are giving less weightage to KL Divergence as compared to Reconstruction Loss. This may make intuitive sense to have since we really don't have any knowledge about the prior and it's just a guess.
- When we set $\beta > 1$, we are giving more weightage to KL Divergence. We would probably want to do only when we want our posterior $p_\phi(z|x)$ to follow a certain shape which we would only want to do if we have some knowledge about our prior ($p(z)$).



The above perspective is from the perspective of prior. We can see this from the perspective of data also i.e. how much we want to learn from data. If we $\beta < 1$, what we are saying is that $p(z|n)$ is allowed to diverge from our prior $p(z)$ i.e. we want to learn more from data. This may result in overfitting.

For $\beta > 1$, we want to stick closer to prior. Essentially, we want to learn less from data and we are putting our bias in terms of selection of prior. This may lead to underfitting.



3 a) In this problem, we need to prove

$$\log p_{\theta}(n) \geq L_m(n) \geq L_1(n)$$

where $L_m(n) = E_{z^{(1)} \dots z^{(m)} \sim q_{\phi}(z|n)} \left(\log \frac{1}{m} \sum_{i=1}^m \frac{p_{\theta}(n, z^{(i)})}{q_{\phi}(z^{(i)}|n)} \right)$

$$L_1(n) = E_{z \sim q_{\phi}(z|n)} \left(\log \frac{p_{\theta}(n, z)}{q_{\phi}(z|n)} \right)$$

First, we'll talk about

$$\log p_{\theta}(n) \geq L_m(n)$$

$$p_{\theta}(n) = \sum_z p_{\theta}(n, z)$$

$$= \sum_{z \in Z} \frac{q(z)}{q(z)} p_{\theta}(n, z)$$

$$= E_{z \sim q(z)} \left[\frac{p_{\theta}(n, z)}{q(z)} \right]$$

Using Monte Carlo which samples 'k' elements at a time.

$$p_{\theta}(n) \approx E_{z^{(1)} \sim q(z)} \left[\frac{1}{k} \sum_{j=1}^k \frac{p_{\theta}(n, z^{(j)})}{q(z^{(j)})} \right]$$

Taking log on both sides

$$\log(p_\theta(x)) \approx \log\left(\mathbb{E}_{z^{(i)} \sim q(z)} \left[\frac{1}{k} \sum_{j=1}^k \frac{p_\theta(x, z^{(i)})}{q(z^{(i)})} \right]\right)$$

Using Jensen's Inequality, we get

$$\geq \mathbb{E}_{z^{(i)} \sim q(z)} \left[\log \left(\frac{1}{k} \sum_{j=1}^k \frac{p_\theta(x, z^{(i)})}{q(z^{(i)})} \right) \right]$$

For better estimates of $z^{(i)}$ we want $z^{(i)}$ to be picked from $q_\phi(z|x)$. So the above equation becomes.

$$\geq \mathbb{E}_{z^{(i)} \sim q_\phi(z|x)} \left[\log \left(\frac{1}{k} \sum_{j=1}^k \frac{p_\theta(x, z^{(i)})}{q_\phi(z^{(i)}|x)} \right) \right]$$

which is RHS. Hence,

$$\boxed{\log p_\theta(x) \geq L_m(x)}$$

Now coming to the second part,

$$L_m(n) \geq L_1(n),$$

$$\log \frac{1}{m} \sum_{i=1}^m \left(\frac{p_\theta(n, z^{(i)})}{q_\phi(z^{(i)}|n)} \right)$$

$$= \log E_{z \sim q_\phi(z^{(i)}|n)} \left(\frac{p_\theta(n, z^{(i)})}{q_\phi(z^{(i)}|n)} \right)$$

Jensen's Inequality states that,

$$\log E[x] \geq E[\log(x)]$$

$$\therefore \geq E_{z \sim q_\phi(z^{(i)}|n)} \left[\log \frac{p_\theta(n, z^{(i)})}{q_\phi(z^{(i)}|n)} \right]$$

$$\geq L_1(n)$$

Hence, $\boxed{L_m(n) \geq L_1(n)}$

There are two intuitions as well for this:

① The more samples we draw from $q_\phi(z|x)$, the better estimator it is for the actual $q_\phi(z|x)$ and hence while training we can bring it closer to the actual $q_\phi(z)$.

② If we could draw all the samples from z i.e. $m \rightarrow \infty$, we could compute the actual $P_\theta(x)$ and hence the $L_m(x) \rightarrow 0$. So higher the value of m , lesser is $L_m(x)$.