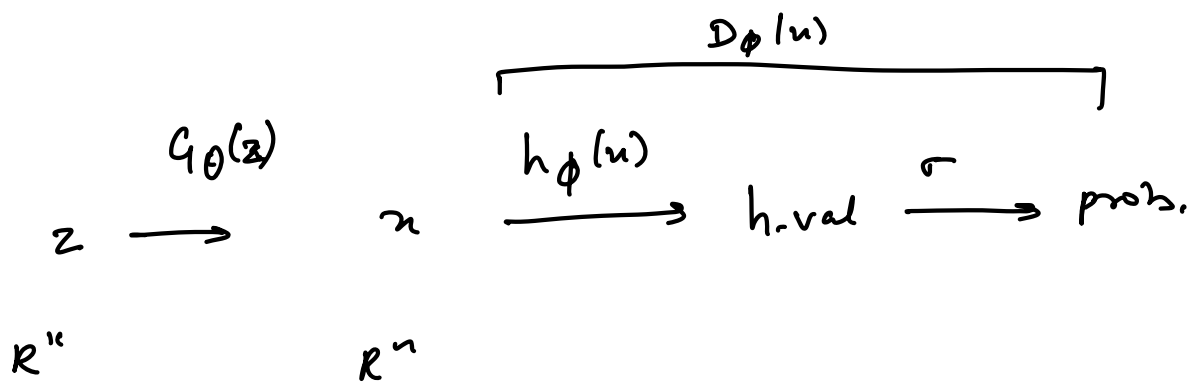


Setup



$$2(a) \quad L_G^{\text{minimax}}(\theta, \phi) = E_{z \sim N(0,1)} [\log (1 - D_\phi(G_\theta(z)))]$$

substituting $D_\phi(x)$ with $\sigma(h_\phi(x))$, we get

$$L_G^{\text{minimax}}(\theta, \phi) = E_{z \sim N(0,1)} [\log (1 - \sigma(h_\phi(G_\theta(z))))]$$

$$\frac{\partial L_G^{\text{minimax}}}{\partial \theta} = E_{z \sim N(0,1)} \left[- \frac{\sigma'(h_\phi(G_\theta(z)))}{1 - \sigma(h_\phi(G_\theta(z)))} \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right]$$

we know that,

$$\sigma'(x) = \sigma(x) (1 - \sigma(x))$$

$$\therefore \frac{\partial L_G^{\text{minimax}}}{\partial \theta} = E_{z \sim N(0,1)} \left[- \frac{\sigma(h_\phi(G_\theta(z))) (1 - \sigma(h_\phi(G_\theta(z))))}{1 - \sigma(h_\phi(G_\theta(z)))} \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right]$$

$$= E_{z \sim N(0,1)} \left[- \sigma(h_\phi(G_\theta(z))) \frac{\partial}{\partial \theta} h_\phi(G_\theta(z)) \right]$$

\uparrow
 ≈ 0 when $D(G_\theta(z)) \approx 0$

$$\therefore \frac{\partial L_G^{\text{minimax}}}{\partial \theta} \approx 0$$

because we are multiplying $\frac{\partial}{\partial \theta} h_D(G_\theta(z))$ with a term ≈ 0 . So, overall value ≈ 0 .

- Why is this a problem for generator?

Ans: It is problematic because if the gradients are very small (≈ 0), it would be hard to train the generator as the updates to the weights θ would be negligible. The training would be essentially stalled.

$$3 a) L_D(\phi, \theta) = -E_{n \sim p_{data}(n)} [\log D_\phi(n)] - E_{n \sim p_\theta(n)} [\log (1 - D_\phi(n))]$$

$$= - \int p_{data}(n) \log D_\phi(n) dn$$

$$- \int p_\theta(n) \log (1 - D_\phi(n)) dn$$

$$= - \int (p_{data}(n) \log D_\phi(n) + p_\theta(n) \log (1 - D_\phi(n))) dn$$

$$= \int f(D_\phi(n)) dn$$

where $f(t) = -p_{data}(n) \log t - p_\theta(n) \log (1-t)$

and $t = D_\phi(n)$ ——— ①

The point where $f(t)$ minimizes (and in turn $L_D(\phi, \theta)$) is where

$$f'(t) = 0$$

$$\Rightarrow \frac{-p_{data}(n)}{t} + \frac{p_\theta(n)}{1-t} = 0$$

$$\Rightarrow \frac{P_0(n)}{1-t} = \frac{P_{data}(n)}{t}$$

$$\Rightarrow \frac{P_0(n)}{P_{data}(n)} = \frac{1-t}{t}$$

$$\Rightarrow \frac{P_0(n)}{P_{data}(n)} = \frac{1}{t} - 1$$

$$\Rightarrow \frac{P_0(n)}{P_{data}(n)} + 1 = \frac{1}{t}$$

$$\Rightarrow \frac{P_0(n) + P_{data}(n)}{P_{data}(n)} = \frac{1}{t}$$

$$\Rightarrow t = \frac{P_{data}(n)}{P_0(n) + P_{data}(n)}$$

Substitution based on ①

\Rightarrow

$$D_{\phi}(n) = \frac{P_{data}(n)}{P_{\theta}(n) + P_{data}(n)}$$

$$\therefore D^*(n) = \frac{P_{data}(n)}{P_{\theta}(n) + P_{data}(n)}$$

\nearrow
point where $f(D_{\phi}(n))$
minimizes

3(b) Given,

$$D_{\phi}(n) = \sigma(h_{\phi}(n))$$

$$= \frac{1}{1 + e^{-h_{\phi}(n)}} \quad \text{--- (1)}$$

If $D_{\phi} = D^*$, then

$$D_{\phi}(n) = \frac{P_{data}(n)}{P_{\theta}(n) + P_{data}(n)} \quad \text{--- from (3a)}$$

From (1), (3a)

$$\frac{1}{1 + e^{-h_{\phi}(n)}} = \frac{P_{data}(n)}{P_{\theta}(n) + P_{data}(n)}$$

$$1 + e^{-h_{\phi}(n)} = \frac{P_{\theta}(n) + P_{data}(n)}{P_{data}(n)}$$

$$1 + e^{-h_{\phi}(n)} = \frac{p_{\theta}(n)}{p_{data}(n)} + 1$$

cancelling "1" on each side and taking
log

$$-h_{\phi}(n) = \log \frac{p_{\theta}(n)}{p_{data}(n)}$$

$$h_{\phi}(n) = - \log \frac{p_{\theta}(n)}{p_{data}(n)}$$

$$= \log \frac{p_{data}(n)}{p_{\theta}(n)}$$

3(c) Given ,

$$L_G(\theta, \phi) = E_{x \sim p_\theta(x)} [\log(1 - D_\phi(x))] \\ - E_{x \sim p_\theta(x)} [\log D_\phi(x)]$$

$$= E_{x \sim p_\theta(x)} \left[\log \frac{1 - D_\phi(x)}{D_\phi(x)} \right]$$

if $D_\phi = D^*$, then

$$L_G(\theta, \phi) = E_{x \sim p_\theta(x)} \left[\log \frac{1 - \frac{p_{data}(x)}{p_{data}(x) + p_\theta(x)}}{\frac{p_{data}(x)}{p_{data}(x) + p_\theta(x)}} \right]$$

$$= E_{x \sim p_\theta(x)} \left[\log \frac{\frac{p_\theta(x)}{p_{data}(x) + p_\theta(x)}}{\frac{p_{data}(x)}{p_{data}(x) + p_\theta(x)}} \right]$$

$$= E_{x \sim p_{\theta}(x)} \left[\log \frac{p_{\theta}(x)}{p_{data}(x)} \right]$$

$$= KL(p_{\theta}(x) \parallel p_{data}(x))$$

$$3(d) \quad n \mathcal{L} = -E_{x \sim p_{data}(x)} [\log p_{\theta}(x)]$$

$$= -E_{x \sim p_{data}(x)} \left[\log \frac{p_{\theta}(x)}{p_{data}(x)} \right]$$

$$= -E_{x \sim p_{data}(x)} \left[\log \frac{p_{\theta}(x)}{p_{data}(x)} \right]$$

$$= -E_{x \sim p_{data}(x)} [\log p_{data}(x)]$$

↑
A term which is independent of ' θ ' and hence can be ignored for the purpose of training the generator as that is tuning ' θ '.

$$= -E_{x \sim p_{data}(x)} \left[\log \frac{p_{\theta}(x)}{p_{data}(x)} \right]$$

$$= E_{x \sim p_{data}(x)} \left[\log \frac{p_{data}(x)}{p_{\theta}(x)} \right]$$

$$= KL(P_{data}(x) \parallel P_{\theta}(x))$$

So, VAE is in essence trying to minimize $KL(P_{data}(x), P_{\theta}(x))$ and L_G is trying to minimize $KL(P_{\theta}(x) \parallel P_{data}(x))$.

Since, KL divergence is not symmetric, I don't think we can say that these objectives are same.

But I would note that these are similar in nature and are varying because of the mechanics of the learning process. In VAE, we start with the actual samples (P_{data}) & we try to see whether model (P_{θ}) can generate similar sample. In GAN, the process starts with generation i.e. P_{θ} . I feel this is cause of different KL divergences.

$$4(a) \quad \ln_{\Phi}(x, y) = \log \frac{P_{\text{data}}(x, y)}{P_{\theta}(x, y)}$$

$$= \log \frac{P_{\text{data}}(x|y)}{P_{\theta}(x|y)} + \log \frac{P_{\text{data}}(y)}{P_{\theta}(y)}$$

$$= \log \frac{P_{\text{data}}(x|y)}{P_{\theta}(x|y)}$$

$$= \log \frac{N(\varphi(x) | \mu_y, I)}{N(\varphi(x) | \hat{\mu}_y, I)}$$

$$= \log \frac{\frac{1}{\sqrt{2\pi}} e^{-1/2 (\varphi(x) - \mu_y)^2}}{\frac{1}{\sqrt{2\pi}} e^{-1/2 (\varphi(x) - \hat{\mu}_y)^2}}$$

$$= \log \left(e^{-\frac{1}{2} \left((\varphi(n) - \mu_y)^2 - (\varphi(n) - \hat{\mu}_y)^2 \right)} \right)$$

$$= -\frac{1}{2} \left((\varphi(n) - \mu_y)^2 - (\varphi(n) - \hat{\mu}_y)^2 \right)$$

$$= -\frac{1}{2} \left(\mu_y^2 - 2\varphi(n)\mu_y - \hat{\mu}_y^2 + 2\varphi(n)\hat{\mu}_y \right)$$

$$= -\frac{1}{2} \left(-2\varphi(n)(\mu_y - \hat{\mu}_y) + \mu_y^2 - \hat{\mu}_y^2 \right)$$

$$= \varphi(n)(\mu_y - \hat{\mu}_y) - \frac{\mu_y^2 - \hat{\mu}_y^2}{2}$$

<To Be Completed>

S(a)

$$KL(P_{\theta}^{(n)} || P_{data}^{(n)}) = E_{n \sim N(\theta, \epsilon^2)} \left[\log \frac{\exp\left(-\frac{1}{2\epsilon^2} (n - \theta)^2\right)}{\exp\left(-\frac{1}{2\epsilon^2} (n - \theta_0)^2\right)} \right]$$

$$= E_{n \sim N(\theta, \epsilon^2)} \left[\log \left(\exp \left(-\frac{1}{2\epsilon^2} \left((n - \theta)^2 - (n - \theta_0)^2 \right) \right) \right) \right]$$

$$= E_{n \sim N(\theta, \epsilon^2)} \left[-\frac{1}{2\epsilon^2} \left((n - \theta)^2 - (n - \theta_0)^2 \right) \right]$$

$$= E_{n \sim N(\theta, \epsilon^2)} \left[\frac{1}{2\epsilon^2} \left((n - \theta_0)^2 - (n - \theta)^2 \right) \right]$$

$$= \frac{1}{2\epsilon^2} E_{n \sim N(\theta, \epsilon^2)} \left[-2n\theta_0 + \theta_0^2 + 2n\theta - \theta^2 \right]$$

$$= \left(\frac{\theta_0^2 - \theta^2}{2\epsilon^2} \right) + \frac{1}{2\epsilon^2} E_{n \sim N(\theta, \epsilon^2)} \left[-2n\theta_0 + 2n\theta \right]$$

$$= \left(\frac{\theta_0^2 - \theta^2}{2\epsilon^2} \right) + \left(\frac{2\theta - 2\theta_0}{2\epsilon^2} \right) E_{n \sim N(\theta, \epsilon^2)} [n]$$

$$= \left(\frac{\theta_0^2 - \theta^2}{2\epsilon^2} \right) + \left[\frac{2\theta - 2\theta_0}{2\epsilon^2} \right] \theta$$

$$= \frac{\theta_0^2 - \theta^2 + 2\theta^2 - 2\theta_0\theta}{2\epsilon^2}$$

$$= \frac{\theta_0^2 + \theta^2 - 2\theta_0\theta}{2\epsilon^2}$$

$$= \frac{(\theta - \theta_0)^2}{2\epsilon^2}$$

hence proved that

$$KL(P_\theta^{(n)} || P_{data}^{(n)}) = \frac{(\theta - \theta_0)^2}{2\epsilon^2}$$

(b) In the above expression, if $\epsilon \rightarrow 0$
the $KL \rightarrow \infty$ i.e. KL divergence
would be a very high number.

$$\frac{\partial KL}{\partial \theta} = \frac{1}{2\epsilon^2} (2(\theta - \theta_0))$$

$$= \frac{\theta - \theta_0}{\epsilon^2}$$

This derivative value also be very high if
 $\theta \neq \theta_0$. The issue that it will cause during
trainings would be the updates to weights
are very large and the training
will not be stable. The weights will
keep moving in different directions
without converging in a stable manner.

(c) All the loss functions we have discussed so far have either a very high value or zero. Essentially, the gradient descent process which navigates the loss curve becomes a process of random search which is randomly looking for the point where $loss = 0$. There is no smooth path which leads to that point. So, we need a function where the loss curve is smooth and there is a path to minima from all the points.