

1.

a.

$$p_{\theta}(x|z) = \prod_{j=1}^G \hat{y}_j^{x_j} (1 - \hat{y}_j)^{1-x_j}$$

b. The output dimension of the encoder should be 2, since  $z \in \mathbb{R}^2$ .

c.

$$\begin{aligned} \log p_{\theta}(x) &= \log \sum_z q_{\phi}(z|x) \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \geq \sum_z q_{\phi}(z|x) \log \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \\ &= \sum_z q_{\phi}(z|x) \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \\ &= \sum_z q_{\phi}(z|x) \log p_{\theta}(x|z) + \sum_z q_{\phi}(z|x) \log \frac{p(z)}{q_{\phi}(z|x)} \\ &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(x|z)] - D_{KL}(q_{\phi}(z|x), p(z)) \end{aligned}$$

d. The KL divergence is non-symmetric and non-negative. Also  $D_{KL}(q, p) = 0$  when  $q$  and  $p$  are the same distribution.

e. No, they will not be the same because the Equation 2 is more computationally expensive than Equation 1 but will be more accurate in the end.

f. No, we build the encoder because it allows us to compute a better Gaussian distribution for our data and slowly improve its accuracy in representing the data. This allows us to reduce our lower bound giving us an overall more accurate representation of our data's distribution.

g. The KL-divergence  $D_{KL}(q_{\phi}(z|x), q_{\phi}(z|x)) = 0$  because  $q_{\phi}(z|x) = q_{\phi}(z|x)$ , This is a basic property of the KL-divergence. It makes sense intuitively because if the distributions are the same then we have a perfect representation of our data so the lower-bound between them would be 0.

h.

$$\begin{aligned} KL(q_{\phi}(z|x), p(z)) &= \sum_z q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{p(z)} \\ &= \sum_z q_{\phi}(z|x) \log \frac{\frac{1}{\sqrt{2\pi\sigma_{\phi}^2}} \exp\left(-\frac{1}{2\sigma_{\phi}^2}(z - \mu_{\phi})^2\right)}{\frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left(-\frac{1}{2\sigma_p^2}(z - \mu_p)^2\right)} \end{aligned}$$

$$\begin{aligned}
&= \sum_z q_\phi(z|x) \log \frac{\exp\left(-\frac{1}{2\sigma_\phi^2}(z - \mu_\phi)^2\right)}{\exp\left(-\frac{1}{2\sigma_p^2}(z - \mu_p)^2\right)} \\
&= \sum_z q_\phi(z|x) \left[ -\frac{1}{2\sigma_\phi^2}(z - \mu_\phi)^2 + \frac{1}{2\sigma_p^2}(z - \mu_p)^2 \right] \\
&= \sum_z q_\phi(z|x) \frac{1}{2\sigma_\phi^2} \left[ -(z - \mu_\phi)^2 + (z - \mu_p)^2 \right] \\
&= \sum_z q_\phi(z|x) \frac{1}{2\sigma_\phi^2} \left[ -z^2 + 2z\mu_\phi - \mu_\phi^2 + z^2 - 2z\mu_p + \mu_p^2 \right]
\end{aligned}$$

Since  $\text{var}(x) = E[x^2] - E[x]^2$

$$= 1 + \sum_z \frac{1}{2\sigma_\phi^2} [2z\mu_\phi - \mu_\phi^2 - 2z\mu_p + \mu_p^2]$$

I am stuck here trying to do it without integrals but eventually we reach that.

$$\begin{aligned}
&KL(q_\phi(z|x), p(z)) \\
&= \log \frac{\sigma}{\sigma} + \frac{\sigma^2 + (\mu_\phi - \mu_p)^2}{2\sigma^2} - \frac{1}{2} = \frac{1}{2} + \frac{(\mu_\phi - \mu_p)^2}{2\sigma^2} - \frac{1}{2} \\
&KL(q_\phi(z|x), p(z)) = \frac{(\mu_\phi - \mu_p)^2}{2\sigma^2}
\end{aligned}$$

i. The Lagrangian is:

$$\begin{aligned}
L(q_\phi(z|x), \lambda) &= \sum_z q_\phi(z|x) \log(p_\theta) - \sum_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p(z)} \\
&\quad + \lambda \left( \sum_z q_\phi(z|x) - 1 \right) \\
0 &= \frac{\partial L}{\partial q_\phi(z_i|x)} \Rightarrow \log q_\phi(z_i|x) = \lambda - 1 + \log p_\theta(x|z_i)p(z_i) \\
1 &= \sum_z q_\phi(z|x) = \sum_{z_i \in Z} p_\theta(x|z_i)p(z_i)e^{\lambda-1} = e^{\lambda-1} \sum_{z_i \in Z} p_\theta(x|z_i)p(z_i) \\
e^{\lambda-1} &= \frac{1}{\sum_{z_i \in Z} p_\theta(x|z_i)p(z_i)} \Rightarrow \lambda = 1 - \log \sum_{z_i \in Z} p_\theta(x|z_i)p(z_i)
\end{aligned}$$

If we substitute in  $\lambda$ :

$$\begin{aligned}
\log q_\phi(z_i|x) &= \lambda - 1 + \log p_\theta(x|z_i)p(z_i) \\
\log q_\phi(z_i|x) &= \log p_\theta(x|z_i)p(z_i) - \log \sum_{z_i \in Z} p_\theta(x|z_i)p(z_i)
\end{aligned}$$

Thus,

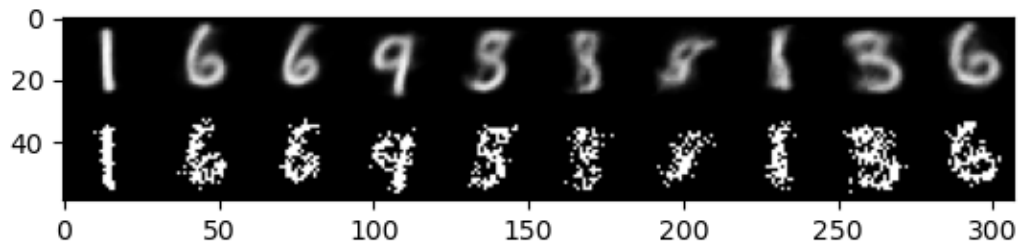
$$q_\phi(z_i|x) = \frac{p_\theta(x|z_i)p(z_i)}{\sum_{z_i \in Z} p_\theta(x|z_i)p(z_i)}$$

j.  $q_\phi(z_i|x)$  should be (3) because according to Bayes Rule:

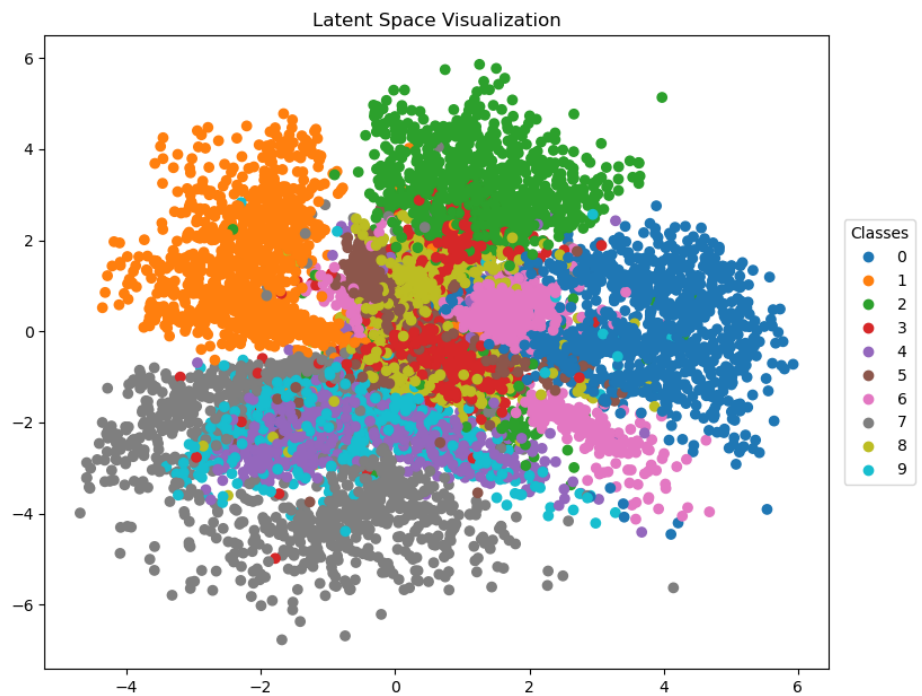
$$q_{\phi}(z_i|x) = \frac{p_{\theta}(x|z_i)p(z_i)}{\sum_{z_i \in \mathcal{Z}} p_{\theta}(x|z_i)p(z_i)} = p_{\theta}(z|x)$$

2.

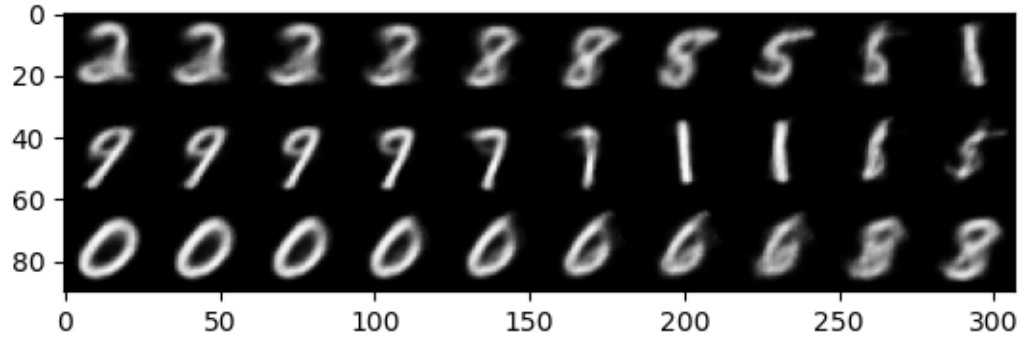
- a. Code
- b. Code
- c. Code
- d. Code
- e. Next page
  - i.



ii.



iii.



3.

a. The cost function for GANs is:

$$\max_{\theta} \min_{\omega} - \sum_x \log D_{\omega}(x) - \sum_z \log(1 - D_{\omega}(G_{\theta}(z)))$$

b. Source: lecture slides

$$\begin{aligned} \min_D - \int_x p_{data} \log D(x) dx - \int_z p_z(z) \log(1 - D(G_{\theta}(z))) dz = \\ \min_D - \int_x p_{data} \log D(x) + p_G(x) \log(1 - D(x)) dx \end{aligned}$$

c. Using Euler-Lagrange formalism:

$$S(D) = \int_x L(x, D, \dot{D}) dx$$

Since there are no  $\dot{D}$ :

$$\begin{aligned} 0 &= \frac{\partial L(x, D, \dot{D})}{\partial D} - \frac{d}{dx} \frac{\partial L(x, D, \dot{D})}{\partial \dot{D}} = \frac{\partial L(x, D, \dot{D})}{\partial D} \\ \frac{\partial L(x, D, \dot{D})}{\partial D} &= \frac{p_{data}(x)}{D(x)} + \frac{p_G(x)}{1 - D(x)} = 0 \end{aligned}$$

So,

$$D^* = \frac{p_{data}(x)}{p_{data}(x) + p_G(x)}$$

d.

$$\begin{aligned} - \int_x p_{data} \log D(x) + p_G(x) \log(1 - D(x)) dx = \\ - \int_x p_{data} \log \frac{p_{data}(x)}{p_{data}(x) + p_G(x)} + p_G(x) \log \frac{p_G(x)}{p_{data}(x) + p_G(x)} dx \end{aligned}$$

Using JSD,

$$\begin{aligned} - \int_x p_{data} \log \frac{p_{data}(x)}{p_{data}(x) + p_G(x)} + p_G(x) \log \frac{p_G(x)}{p_{data}(x) + p_G(x)} dx \\ = -2 JSD(p_{data}, p_G) + \log(4) \end{aligned}$$

and

$$JSD(p_{data}, p_G) = \frac{1}{2} KL(p_{data}, M) + \frac{1}{2} KL(p_G, M) \text{ with } M = \frac{1}{2} (p_{data} + p_G)$$

Thus, the optimal Generator  $G^*(x)$  generates the distribution  $p_G^*$ , such that  $p_G^* = p_{data}$ .

- e.  $KL(P_1, P_2) = Undefined$ ,  $KL(P_1, P_3) = Undefined$ . The KL divergence would be undefined for both since there are some values of  $x$  for which the probability in  $P_1$  is 0 or  $P_2$  is 0 and same for  $P_1$  and  $P_3$ , while the other is nonzero. This would lead to a division by 0 in the KL divergence equation.

$$W_1(P_1, P_2) = 0.5, \quad W_1(P_1, P_3) = 1$$

4.

- Code
- Code
- Code
- Code
- Code
- Epoch: 10



Epoch: 30



Epoch: 50



Epoch: 90

