

### Question 1

Describe steps needed to sample from such a model

$\Rightarrow$

We first obtain set of sample  $z^{(i)}$  where  $i = 1, \dots, L$

$\Rightarrow$  We allow expectation

be approximated by  $E[f] = \int f(z) p(z) dz$  to

$$\hat{f} = \frac{1}{L} \sum_{i=1}^L f(z^{(i)})$$

$\Rightarrow$  We let directed graph with some nodes instantiated with values given by

$$P(L=L^0, S=S^0) = \sum_{D, F, G} P(D) P(F) P(G)$$

$\Rightarrow$  Next we input the factor graph

$\Rightarrow$  We order all given variables from top to bottom so that the parent variables come before child variables

$\Rightarrow$  Next, for each variable in the parent variable, we retrieve their sampled values which exist already after the ordering

$\Rightarrow$  However, with the parent factor being deterministic, this will simply compute the child value

$\Rightarrow$  Lastly, we output the sampled values, within each variable in the graph

## Question 2

- $\Rightarrow$  This is because we cannot compute any given new point likelihood say  $y$  under the probability distribution efficiently, since it involves integrating over the hidden variable.
- $\Rightarrow$  Also the naive spherical Gaussian noise model produces noisy sample or rather overly smooth ones instead if noise is not added to it.
- $\Rightarrow$  Thus the system may converge into a local minimum for which any latent variable is completely ignored with the encoder predicting prior resulting into posterior collapse.

## Question 3

- $\Rightarrow$  Example resulting very small KL-divergence

$$\Rightarrow KL(p, q) = \log \frac{\sigma_1}{\sigma_2} + \frac{\sigma_2^2 - (\mu_1 + \mu_2)^2}{2\sigma_1^2} - \frac{1}{2}$$

- $\Rightarrow$  Example resulting very large KL-divergence

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_1^2} + \frac{1}{2}$$



#### Question 4

$\Rightarrow$  This will be because the log-probability will be an expectation under the complicated posterior distribution which is what we are to approximate, and we cannot usually evaluate it

$\Rightarrow$  However, with the lower-bound this encourages the fit to concentrate on the plausible parameters

#### Question 5

$\Rightarrow$  (1) This basically optimizes the likelihood using the approximate distribution thus maximizing the log-likelihood

$\Rightarrow$  (2) This forces the approximation to match the true posterior meaning that when variational approximate posterior is perfect then KL-Vanishes

#### Question 6

$\Rightarrow$  So with the optimization of the log-likelihood this results into ignoring the precision which only can be fine-tuned through hyper-parameter tuning hence will result into balancing problems for the relevant parameters  
~~just that~~

### Question 7

PAGE NO.:

DATE: / /

$\Rightarrow$  Given that  $p_\theta(z_n|x_n)$  we have that

$$p_\theta(z|x) = \frac{p_\theta(x|z)p_\theta(z)}{p_\theta(x)}$$

$$= \frac{p_\theta(x|z)p_\theta(z)}{\int p_\theta(x|z)p_\theta(z)dz}$$

$\Rightarrow$  This is equivalent of  $q_\phi(z|x)$

$\Rightarrow$  Thus, the log-marginal likelihood will be given as

$$\Rightarrow \log p_\theta(x) = \Delta_{KL}[q_\phi(z|x) || p_\theta(z|x)] + L(\theta, \phi|x)$$

$$\Rightarrow L(\theta, \phi|x) = E_{q_\phi(z|x)}(\log p_\theta(x, z) - \log q_\phi(z|x))$$

$$\Rightarrow \int [\log p_\theta(x|z) - \log q_\phi(z|x) + \log p_\theta(z)] q_\phi(z|x) dz$$

$\Rightarrow$  This simplifies into

$$\int \frac{p_\theta(z)}{q_\phi(z|x)} q_\phi(z|x) dz + \int \log p_\theta(x|z) q_\phi(z|x) dz$$

$$\Rightarrow E_{q_\phi(z|x)}[\log p_\theta(x|z) - \Delta_{KL}[q_\phi(z|x) || p_\theta(z)]]$$

$\Rightarrow$  Optimizing the ELBO, we have the mean lower bound over sample given as



$$\Rightarrow L(\theta, \phi(z)) = -\frac{1}{N} \sum_{n=1}^N -\Delta_{KL}(q_{\phi}(z|z_n) || p_{\theta}(z)) + E_{q_{\phi}}[\log p_{\theta}(z_n|z)]$$

$\Rightarrow$  let the reconstruction term be given as

$$\Rightarrow L_n^{\text{recon}} = -E_{q_{\phi}(z|z_n)}[\log p_{\theta}(z_n|z)]$$

with  $\text{reg}$

$$\Rightarrow L_n^{\text{reg}} = \Delta_{KL}(q_{\phi}(z|z_n) || p_{\theta}(z))$$

$\Rightarrow$  Thus the loss minimization function becomes

$$L(\theta, \phi(z)) = \frac{1}{N} \sum_{n=1}^N L_n^{\text{recon}} + L_n^{\text{reg}}$$



### Question 8

$\Rightarrow$  This normally happens because the network simply involves a sampling step and there is literally zero way to differentiate through this, implying that it will be impossible to make updates to the parameters that will occur earlier within the network

$\Rightarrow$  However, via reparameterization trick we can simply move the stochastic part into numerous branches of the network which instead draw a sample from the given norm to draw from intended Gaussian function

$\Rightarrow$  Under this approach we can compute the derivatives since there is no need for the backpropagation algorithm to pass down the stochastic branch

Question 9 8/10

See attached code files

Question 11

$$\Rightarrow \text{From } L(x; \theta_d, \phi_g) = \frac{1}{n} \sum_{i=1}^n [\log \Delta_\phi(x^{(i)}) + \log(1 - \Delta_\phi(G_\theta(z^{(i)})))]$$

$$\Rightarrow \nabla_{\theta_d} L(x; \theta_d, \phi_g) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_d} \log \Delta_\phi(x^{(i)}) + \nabla_{\theta_d} \log(1 - \Delta_\phi(G_\theta(z^{(i)}))]$$

$$= \nabla_{\theta_d} \sum_{i=1}^n \log \Delta_\phi(x^{(i)}) + \log(1 - \Delta_\phi(G_\theta(z^{(i)}))]$$

Question 12

$\Rightarrow$  From

$$L(\theta_d, \theta_g) = \frac{1}{n} \sum_{i=1}^n \log \Delta(x^{(i)}) + \log(1 - \Delta(G(z^{(i)}))]$$

$$\Rightarrow \frac{\partial L(\theta_d, \theta_g)}{\partial z_d} = \frac{1}{n} \sum_{i=1}^n 0 + \frac{(0 - \Delta(G(z^{(i-1)})))}{1 - \Delta(G(z^{(i)}))}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{-\Delta(G(z^{(i-1)}))}{1 - \Delta(G(z^{(i)}))}$$



### Question 13

$$\Rightarrow \frac{\partial L(\theta_d, \theta_g)}{\partial z_d^i} = \frac{1}{n} \sum_{i=1}^n \frac{-\Delta(G(z^{(i)}))}{1 - \Delta(G(z^{(i)}))}$$

for next value  $i+1$  we will have

$$= \frac{1}{n} \sum_{i=1}^n \frac{-\Delta(G(z^{(i+1)}))}{1 - \Delta(G(z^{(i+1)}))}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{-\Delta(G(z^{(i)}))}{1 - \Delta(G(z^{(i+1)}))}$$

with  $w_d^{i+1} = \Delta(G(z^{(i+1)}))$  and

$$g_d^i = \Delta(G(z^{(i)}))$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{-g_d^i}{1 - w_d^{i+1}}$$

### Question 14

$\Rightarrow$  with the output generator function

$$L(\theta_d, \theta_g) = \frac{1}{n} \sum_{i=1}^n \log(1 - \Delta(G(z^{(i)})))$$

$$\Rightarrow \frac{\partial L(\theta_d, \theta_g)}{\partial g(z_i, \theta_g)} = \frac{1}{n} \sum_{i=1}^n \frac{\Delta(z^{(i)})}{1 - \Delta(G(z^{(i)}))} \quad \text{with}$$

$$w_g^i = \Delta(z^{(i)})$$

and

$$\frac{\partial L(\theta_d, \theta_g)}{\partial z_d^i} = \frac{1}{1 - \Delta(G(z^{(i)}))}$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \frac{\Delta(z^{(i)})}{1 - \Delta(G(z^{(i)}))}$$

Question 15

$$\Rightarrow L(\theta_d, \theta_g) = \frac{1}{n} \sum_{i=1}^n \log(1 - \Delta(G(z^{(i)})))$$

$$\Rightarrow \nabla_{\theta_g} L(\theta_d, \theta_g) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_g} \log(1 - \Delta(G(z^{(i)})))$$

$$\Rightarrow \frac{\nabla_{\theta_g}}{n} \sum_{i=1}^n \frac{\Delta(z^{(i)})}{1 - \Delta(G(z^{(i)}))}$$

Question 16

$$\nabla_{\theta_d} L(\theta_d, \theta_g) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_d} \log(1 - \Delta(G(z^{(i)})))$$

$$= \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_d} \frac{0 - G(z^{(i)})}{1 - \Delta(G(z^{(i)}))}$$

$$= \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_d} \frac{-G(z^{(i)})}{1 - \Delta(G(z^{(i)}))}$$

$$\Rightarrow \frac{\nabla_{\theta_d}}{n} \sum_{i=1}^n \frac{-G(z^{(i)})}{1 - \Delta(G(z^{(i)}))}$$



Question 17

$$\Rightarrow \nabla_{\theta_g} L(\theta_d, \theta_g) = \nabla_{\theta_g} \sum_{i=1}^n \log(1 - \Delta(G(z^{(i)})))$$

with fixed learning rate

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \nabla_{\theta_g} \frac{\Delta(G(z^{(i)}))}{1 - \Delta(G(z^{(i)}))}$$

$$\Rightarrow \nabla_{\theta_g} \sum_{i=1}^n \frac{\Delta(G(z^{(i)}))}{1 - \Delta(G(z^{(i)}))}$$

Question 18

$\Rightarrow$  Deriving  $\frac{\partial x_{L+1}}{\partial x_L}$

$$\Rightarrow \text{From } \frac{\partial L}{\partial x^{L+1}} = [W^{(L)}]^T \cdot \frac{\partial L}{\partial x^L} \cdot g^{(L-1)} x^{(L-1)}$$

$$\Rightarrow \frac{\partial L}{\partial w^{(L)}} = \frac{\partial L}{\partial x^L} [a^{(L-1)}]^T$$

$$\Rightarrow \frac{\partial L}{\partial a^{(L)}} = \frac{\partial L}{\partial x^L}$$

$\Rightarrow$  For a mini-batch extension we have

$$a^{(L-1)} = \begin{pmatrix} a_1^{(L-1)} \\ a_2^{(L-1)} \\ \vdots \\ a_n^{(L-1)} \end{pmatrix}$$

⇒ This error can be propagated for layer  $L$  back to its previous layer  $L-1$  by

$$\Rightarrow \frac{\partial L}{\partial x^{(L-1)}} = [W^{(L)}]^T \cdot \frac{\partial L}{\partial x^{(L)}} \cdot g^{(L-1)'}(x^{(L-1)})$$

⇒ Implied that for all mini-batch samples

$$\frac{\partial L}{\partial w^{(L)}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L^{(i)}}{\partial w^{(L)}} \quad \text{and}$$

$$\frac{\partial L}{\partial a^{(L)}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L^{(i)}}{\partial a^{(L)}}$$

⇒ Applying chain rule, we obtain that

$$\frac{\partial x_{L+1}}{\partial x_L} = \frac{\partial x_{L+1}}{\partial z} \bigg|_{z=f(x)} \cdot \frac{\partial f(x)}{\partial x_L}$$

$$\Rightarrow \frac{\partial x_{L+1}}{\partial x_L}(x_1, x_2, \dots, x_n) = \frac{\partial x_{L+1}(z_1, \dots, z_n)}{\partial z_L} \cdot p$$

$z_1 = f_1(x_1, x_2)$   
 $z_2 = f_2(x_1, x_2)$   
 $\vdots$   
 $z_L = f_L(x_1, x_2)$

where  $p = \frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_L}$

⇒ For the one-pass and  $x_L$  for the top most layer we have



$$\frac{\partial x_{L+1}(x_1, x_2)}{\partial x_1} = \frac{\partial x_{L+1}(z_1, z_2)}{\partial z_1} \cdot \frac{\partial f(x_1, x_2)}{\partial x_1}$$

$z_1 = f_1(x_1, x_2)$   
 $z_2 = f_2(x_1, x_2)$

which simplifies into

$$\frac{\partial x_{L+1}(x_1, \dots, x_n)}{\partial x_i} = \sum_{j=1}^n \frac{\partial x_{L+1}(z_1, \dots, z_n)}{\partial z_j} \cdot \frac{\partial f_j(x_1, \dots, x_n)}{\partial x_i}$$

### Question 19

=> The consequence is that when training deep models the backpropagation learning rule tends to get stuck in the local minima

=> This error will be a function of all weights in a given multidimensional space. However, there is no proof that this global minimum error surface has been reached already

=> Thus, the derivative of the given transfer function will vanish to zero for very small summed inputs and similarly for very large summed inputs as well

### Question 20

See attached code file

THANK YOU!!!