

**Question 1.** Answer

1.1 The expected complete data log likelihood (ECLL) :

$$\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})] = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

we can rearrange

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) &= \log p_{\theta}(\mathbf{x}) - \mathbb{E}_q \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \\ &= \log p_{\theta}(\mathbf{x}) - \mathbb{E}_q \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{\frac{p_{\theta}(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})}} \\ &= \log p_{\theta}(\mathbf{x}) - \mathbb{E}_q [\log q_{\phi}(\mathbf{z}|\mathbf{x}) - \log p_{\theta}(\mathbf{z}, \mathbf{x}) + \log p_{\theta}(\mathbf{x})] \\ &= -\mathbb{E}_q [\log q_{\phi}(\mathbf{z}|\mathbf{x}) - \log p_{\theta}(\mathbf{z}, \mathbf{x})] \\ &= \mathbb{E}_q [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + H(q) \end{aligned}$$

Since the  $H(q)$  is non-negative and  $q(\mathbf{z}|\mathbf{x})$  is fixed, hence, maximizing  $ECLL = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}, \mathbf{z})] = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})]$  w.r.t  $\theta$  is equal to maximize  $\log p_{\theta}(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$ . It has the maximum value where the value of  $D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$  is 0. This means the maximizer of the ECLL coincides with that of the marginal likelihood only if  $q(\mathbf{z}|\mathbf{x})$  perfectly matches  $p(\mathbf{z}|\mathbf{x})$ .

1.2

$$\begin{aligned} \mathcal{L}(\theta, \phi; \mathbf{x}) &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x})||p(\mathbf{z}))] \\ &= \mathbb{E}_{q_{\phi}}[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p(\mathbf{z})} - \log q_{\phi}(\mathbf{z} | \mathbf{x}) + \log q(\mathbf{z})] \\ &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log p(\mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x}) + \log q(\mathbf{z})] \\ &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x})] \\ &= \log p_{\theta}(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \end{aligned}$$

Hence,

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) - \mathcal{L}(\theta, \phi; \mathbf{x}) &= D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &= D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x})||p(\mathbf{z})) - D_{\text{KL}}(q_i^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) \\ &\quad + D_{\text{KL}}(q_i^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) - D_{\text{KL}}(q_{\phi}^*(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) \end{aligned}$$

Since the  $q_{\phi}$  is parameterized by a neural network, it means there is an additional constraint imposed by requiring that all the variational parameters lie in the range of the network. Within a specific parametric family, nothing is more general than freely optimizing the variational parameters  $_{\phi}$ . This cost is known as the amortization gap (last row in above equation). For a network with infinite capacity, this gap goes away, and we can be (only) as good as the free-form optimization setting. However that is not the case in any practical implementations, as all networks have finite capacity. Therefore,  $D_{\text{KL}}(q_i^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) \leq D_{\text{KL}}(q_{\phi}^*(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$

- 1.3 (a) In terms of bias in estimating the marginal likelihood via the ELBO, in the best case scenario (i.e. when both approaches are optimal within the respective families), The  $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$  has less bias comparing with  $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$  because we model the  $q_i^*$  for each observation and optimize them jointly. While with  $q_{\phi^*}$ , the family chosen is simpler than the  $q_i^*$ , it will invoke higher bias.
- (b) In terms of the computational point of view (efficiency), we use  $q_{\phi^*}$  to optimize the parameters of the neural network instead of the individual parameters of each observation. Hence the parameters won't depending on the number of observations. So, with the presence of massive data points, the  $q_{\phi^*}$  approach is more efficient.
- (c) In terms of memory (storage of parameters), same argument as (b), we specify our parameters of  $q_{\phi^*}$  and do not rely on the instance numbers the instance numbers. So, we have could have fewer parameters when we have massive observations. Hence. having a maximizer  $\phi^*$  of the ELBO requires less memory than  $q_i$  for each  $x_i$  since the number of parameters for the latter need to optimize grows (at least) linearly with the number of observations.

**Question 2** (5-5-5-5). One way to enforce autoregressive conditioning is via masking the weight parameters.<sup>1</sup> Consider a two-hidden-layer convolutional neural network without kernel flipping, with kernel size  $3 \times 3$  and padding size 1 on each border (so that an input feature map of size  $5 \times 5$  is convolved into a  $5 \times 5$  output). Define mask of type A and mask of type B as

$$(\mathbf{M}^A)_{::ij} := \begin{cases} 1 & \text{if } i < 2 \\ 1 & \text{if } i = 2 \text{ and } j < 2 \\ 0 & \text{elsewhere} \end{cases} \quad (\mathbf{M}^B)_{::ij} := \begin{cases} 1 & \text{if } i < 2 \\ 1 & \text{if } i = 2 \text{ and } j \leq 2 \\ 0 & \text{elsewhere} \end{cases}$$

where the index starts from 1. Masking is achieved by multiplying the kernel with the binary mask (elementwise). Specify the receptive field of the output pixel that corresponds to the third row and the third column (index 33 of Figure 1 (Left)) in each of the following 4 cases :

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 1 – (Left)  $5 \times 5$  convolutional feature map. (Right) Template answer.

1. If we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^A$  for the second layer.
2. If we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^B$  for the second layer.
3. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^A$  for the second layer.
4. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^B$  for the second layer.

Your answer should look like Figure 1 (Right).

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

**Answer 1.** 1. 

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

 2. 

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

 3. 

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

 4. 

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

1. An example of this is the use of masking in the Transformer architecture.

**Question 3.** Answer

3.1 Based on the question, we have the following properties of the normalizing flows :

$$\begin{aligned}
 \mathbf{x} &= F(\mathbf{u}), \text{ where } \mathbf{u} \sim P_U(\mathbf{u}), \mathbf{u} = F^{-1}(\mathbf{x}) \\
 P_X(\mathbf{x}) &= P_U(\mathbf{u}) \left| \det \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \right| = P_U(\mathbf{u}) |\det J_{F^{-1}}(\mathbf{x})| \\
 P_U(\mathbf{u}) &= P_X(\mathbf{x}) \left| \det \frac{\partial \mathbf{x}}{\partial \mathbf{u}} \right| = P_X(\mathbf{x}) |\det J_F(\mathbf{u})| \\
 P_{F^{-1}(X)}(\mathbf{u}) &= P_X(F(\mathbf{u})) |\det J_F(\mathbf{u})| \\
 P_{F(U)}(\mathbf{x}) &= P_U(F^{-1}(\mathbf{x})) |\det J_{F^{-1}}(\mathbf{x})| = \frac{P_U(F^{-1}(\mathbf{x}))}{|\det J_F(\mathbf{u})|}
 \end{aligned}$$

We can first show that :

$$\begin{aligned}
 D_{KL}[P_X || P_{F(U)}] &= \int_{\mathbf{x}} P_X(\mathbf{x}) \log \frac{P_X(\mathbf{x})}{P_{F(U)}(\mathbf{x})} d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_X(\mathbf{x}) - \log P_{F(U)}(\mathbf{x})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_U(\mathbf{u}) + \log |\det J_{F^{-1}}(\mathbf{x})| - \log P_{F(U)}(\mathbf{x})] d\mathbf{x}
 \end{aligned}$$

Then we can show for the reverse KL divergence as :

$$\begin{aligned}
 D_{KL}[P_{F^{-1}(X)} || P_U] &= \int_{\mathbf{u}} P_{F^{-1}(X)}(\mathbf{u}) \log \frac{P_{F^{-1}(X)}(\mathbf{u})}{P_U(\mathbf{u})} d\mathbf{u} \\
 &= \int_{\mathbf{u}} P_{F^{-1}(X)}(\mathbf{u}) [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] d\mathbf{u} \\
 &= \int_{\mathbf{u}} P_{F^{-1}(X)}(\mathbf{u}) [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] dF^{-1}(\mathbf{x}) \\
 &= \int_{\mathbf{x}} P_{F^{-1}(X)}(\mathbf{u}) [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] |\det J_{F^{-1}}(\mathbf{x})| d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_{F^{-1}(X)}(\mathbf{u}) |\det J_{F^{-1}}(\mathbf{x})| [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(F(\mathbf{u})) |\det J_F(\mathbf{u})| |\det J_{F^{-1}}(\mathbf{x})| [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(F(\mathbf{u})) [\log P_{F^{-1}(X)}(\mathbf{u}) - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(F(\mathbf{u})) [P_X(F(\mathbf{u})) |\det J_F(\mathbf{u})| - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_X(\mathbf{x}) + \log |\det J_F(\mathbf{u})| - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_U(\mathbf{u}) + \log |\det J_{F^{-1}}(\mathbf{x})| + \log |\det J_F(\mathbf{u})| - \log P_U(\mathbf{u})] d\mathbf{x} \\
 &= \int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_U(\mathbf{u}) + \log |\det J_{F^{-1}}(\mathbf{x})| - \log \frac{P_U(\mathbf{u})}{|\det J_F(\mathbf{u})|}] d\mathbf{x}
 \end{aligned}$$

By replacing  $\frac{P_U(\mathbf{u})}{|\det J_F(\mathbf{u})|} = P_{F(U)}(\mathbf{x})$ , we have the above equation equal to

$$\int_{\mathbf{x}} P_X(\mathbf{x}) [\log P_U(\mathbf{u}) + \log |\det J_{F^{-1}}(\mathbf{x})| - \log P_{F(U)}(\mathbf{x})] d\mathbf{x}$$

which is the same as the final equation from the  $D_{KL}[P_X||P_{F(U)}]$

- 3.2 (a) For scenario 1), we would use the reverse KL divergence  $D_{KL}[P_{F(U)}||P_X]$  as the objective to optimize. We can write the objective function as :

$$\mathcal{L} = D_{KL}[P_{F(U)}||P_X] = -E_{P_{F(U)}}[\log P_{F(U)} - \log P_X] \quad (1)$$

$$= -E_{p_U(\mathbf{u})}[\log p_U(\mathbf{u}) - \log |\det J_F(\mathbf{u})| - \log P_X(F(\mathbf{u}))]. \quad (2)$$

In order to minimize the reverse KL divergence as described above, we need to sample from the base distribution  $P_U$  as well as compute and differentiate through the transformation  $F$  and its Jacobian determinant and evaluation of  $P_X$ . It is not necessarily sample from  $P_X$ .

- (b) For scenario 2), we would use the forward KL divergence  $D_{KL}[P_X||P_{F(U)}]$  as the objective to optimize. The objective function for the forward KL divergence  $D_{KL}[P_X||P_{F(U)}]$  can be written as follows :

$$\mathcal{L} = D_{KL}[P_X||P_{F(U)}] = -E_{P_X(\mathbf{x})}[\log P_{F(U)}] + const. \quad (3)$$

$$= -E_{P_X(\mathbf{x})}[\log p_U(\mathbf{u})(F^{-1}(\mathbf{x})) + \log |\det J_{F^{-1}}(\mathbf{x})|] + const. \quad (4)$$

Since we have samples from the data distribution  $P_X(\mathbf{x})$  and we assume a set of samples  $\{\mathbf{x}_n\}_{n=1}^N$  from  $P_X(\mathbf{x})$ , we can estimate  $E_{P_X}$  by Monte Carlo with following :

$$\mathcal{L} \approx -\frac{1}{N} \sum_{n=1}^N [\log p_U(\mathbf{u})(F^{-1}(\mathbf{x})) + \log |\det J_{F^{-1}}(\mathbf{x})|] + const.$$

and then minimize the above Monte Carlo approximation of the KL divergence with stochastic gradient-based methods. Hence, we don't need to evaluate  $P_X(\mathbf{x})$ .

**Question 4.** Answer

1. Maximizing  $V(d, g)$  w.r.t  $d$  implies :  $d_{k+1} = d_k + \alpha \nabla_{d_k} V(d, g)$  ; Minimizing  $V(d, g)$  w.r.t  $g$  implies :  $g_{k+1} = g_k - \alpha \nabla_{g_k} V(d, g)$ .

Since the  $\nabla_{d_k} V(d, g) = g_k$  and  $\nabla_{g_k} V(d, g) = d_k$ , hence

$$[d_{k+1}, g_{k+1}]^\top = \begin{pmatrix} 1 & \alpha \\ -\alpha & 1 \end{pmatrix} \begin{pmatrix} d_k \\ g_k \end{pmatrix}$$

Thus,

$$A = \begin{pmatrix} 1 & \alpha \\ -\alpha & 1 \end{pmatrix}$$

2. We know that a stationary point is a point on the surface of the graph (of the function) where all its partial derivatives are zero (equivalently, the gradient is zero).

In this case, by differentiating, we get :  $\begin{pmatrix} d_k \\ g_k \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ . Therefore the stationary points on this graph occur when  $d_k = g_k = 0$ , which is  $d = g = 0$ . Therefore the coordinates of the stationary point are (0,0).

3. Assume  $\exists \lambda, s.t. |\lambda| \geq 1$ ,  $\lambda + Sp(A)$ , then  $\exists x, s.t. Ax = \lambda x$ . we can see that

$$\left\| \begin{pmatrix} d_{k+1} \\ g_{k+1} \end{pmatrix} \right\| = \left\| A \begin{pmatrix} d_k \\ g_k \end{pmatrix} \right\| = |\lambda|^{k+1} \left\| \begin{pmatrix} d_0 \\ g_0 \end{pmatrix} \right\|$$

Since  $\lambda \geq 1$ , when  $k \rightarrow \infty$  the behaviour of  $d_k$  and  $g_k$  is diverging.

Also, since the

$$Trace(A) = 1 + 1 = 2$$

$$det(A) = 1 + \alpha^2 > 1$$

Hence, as we increasing the value of  $\alpha$ , it implies a faster divergence.

**Question 5.** Answer

5.1 Since,

$$\begin{aligned}
 & \|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2 \\
 = & \\
 & (W_p \mathbf{f}_1 - \mathbf{f}_2)^\top (W_p \mathbf{f}_1 - \mathbf{f}_2) \\
 = & \\
 & \mathbf{f}_1^\top W_p^\top W_p \mathbf{f}_1 - \mathbf{f}_2^\top W_p \mathbf{f}_1 - \mathbf{f}_1^\top W_p^\top \mathbf{f}_2 + \mathbf{f}_2^\top \mathbf{f}_2 \\
 = & \\
 & \text{tr}(W_p^\top W_p \mathbf{f}_1^\top \mathbf{f}_1) - \text{tr}(W_p \mathbf{f}_1 \mathbf{f}_2^\top) - \text{tr}(W_p^\top \mathbf{f}_2 \mathbf{f}_1^\top) + \text{tr}(\mathbf{f}_2^\top \mathbf{f}_2)
 \end{aligned}$$

Since  $F_1 = F_2 = \mathbb{E} [\mathbf{f}_1 \mathbf{f}_1^\top] = W(X + X')W^\top$  and  $F_{12} = F_{21}^\top = \mathbb{E} [\mathbf{f}_1 \mathbf{f}_2^\top] = WXW^\top$ .  
Hence, we have

$$\begin{aligned}
 & \text{tr}(W_p^\top W_p \mathbf{f}_1^\top \mathbf{f}_1) - \text{tr}(W_p \mathbf{f}_1 \mathbf{f}_2^\top) - \text{tr}(W_p^\top \mathbf{f}_2 \mathbf{f}_1^\top) + \text{tr}(\mathbf{f}_2^\top \mathbf{f}_2) \\
 = & \\
 & \frac{1}{2} \text{tr}(W_p^\top W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(W_p^\top F_{12}) + \text{tr}(F_2)
 \end{aligned}$$

Thus,

$$\begin{aligned}
 & J(W, W_p) = \frac{1}{2} \mathbb{E}_{x_1, x_2} [\|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2] \\
 = & \\
 & J(W, W_p) = \frac{1}{2} [\text{tr}(W_p^\top W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)]
 \end{aligned}$$

5.2

$$\dot{W}_p = -\frac{\partial J}{\partial W_p} = -\frac{1}{2}(2W_p F_1 - F_{12}^\top - F_{12}) = -W_p F_1 + F_{12}^\top$$

5.3 Since we don't have stop gradient and  $W$  is the weight for the target network, then we have additional terms. With  $\tilde{W}_p(t) = W_p(t) - I_{n_2}$  and we have :

$$\begin{aligned}
 & \dot{W} = -\frac{\partial J}{\partial W} \\
 = & \\
 & -W_p^\top W_p W(X + X') + (W_p^\top + W_p)WX - W(X + X') \\
 = & \\
 & -(W_p^\top W_p + I)WX' - (W_p^\top W_p - W_p^\top - W_p + I)WX \\
 = & \\
 & -(W_p^\top W_p + I)WX' - (W_p - I)^\top (W_p - I)WX \\
 = & \\
 & -(W_p^\top W_p + I)WX' - \tilde{W}_p^\top \tilde{W}_p WX
 \end{aligned}$$

with  $\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X)$  and we have :

$$\frac{d}{dt} \text{vec}(W) = -[X' \otimes (W_p^T W_p + I) + X \otimes \tilde{W}_p^T \tilde{W}_p] \text{vec}(W)$$

Hence

$$\frac{d}{dt} \text{vec}(W(t)) = -H(t) \text{vec}(W(t))$$

if  $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$ , then apply the proved property : " $\frac{d}{dt} \mathbf{w}(t) = -H(t) \mathbf{w}(t)$ ", satisfies the constraint  $\|\mathbf{w}(t)\|_2 = e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$ , implying that  $\mathbf{w}(t) \rightarrow 0$ ", we have  $\|\text{vec}(W(t))\|_2 = e^{-\lambda_0 t} \|\text{vec}(W(0))\|_2 \rightarrow 0$  and there is no chance for  $W$  to learn any meaningful features.

5.4 In the case when both the Stop-Grad **and** the predictor are removed ( $W_P = I$ ), the gradient of  $W(t)$  is :

$$\dot{W}(t) = -\frac{\partial J}{\partial W(t)} = -W_P^\top (W_p W (X + X') + W X) = -W X'$$

$X'$  is positive definite matrix with the same property as the previous question  $W(t) \rightarrow 0$  .

5.5 We can see from Q5.2–5.4 that by combining the stop-gradient and the predictor, the  $\dot{W}(t)$  won't become zero as time passes. As a result, the  $W$  can learn some useful features in addition to the training.