

Due Date: April 29th 23:59, 2020Instructions

- For all questions, show your work!
- Please use a document preparation system such as LaTeX, unless noted otherwise.
- Unless noted that questions are related, assume that notation and definitions for each question are self-contained and independent.
- Submit your answers electronically via Gradescope.
- **TAs for this assignment are Samuel Lavoie, Jae Hyun Lim, Sanae Lotfi.**

This assignment covers mathematical and algorithmic techniques underlying the four most popular families of deep generative models. Thus, we explore autoregressive models (Question 1), reparameterization trick (Question 2), variational autoencoders (VAEs, Questions 3-4), normalizing flows (Question 5), and generative adversarial networks (GANs, Question 6).

Question 1 (4-4-4-4). One way to enforce autoregressive conditioning is via masking the weight parameters.¹ Consider a two-hidden-layer convolutional neural network without kernel flipping, with kernel size 3×3 and padding size 1 on each border (so that an input feature map of size 5×5 is convolved into a 5×5 output). Define mask of type A and mask of type B as

$$(\mathbf{M}^A)_{::ij} := \begin{cases} 1 & \text{if } i = 2 \text{ and } j < 2 \\ 1 & \text{if } i = 3 \\ 0 & \text{elsewhere} \end{cases} \quad (\mathbf{M}^B)_{::ij} := \begin{cases} 1 & \text{if } i = 2 \text{ and } j \leq 2 \\ 1 & \text{if } i = 3 \\ 0 & \text{elsewhere} \end{cases}$$

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 fourth column (index 34 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) 5×5 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.

Using \mathbf{M}^A for the first and second layer, below is the receptive field of pixel 34,

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 2 – (Left) Receptive field of output pixel 34 when Mask-A is applied to both the layers.

1. An example of this is the use of masking in the Transformer architecture (Problem 3 of HW2 practical part).

2. If we use \mathbf{M}^A for the first layer and \mathbf{M}^B for the second layer. Using \mathbf{M}^A for the first layer and \mathbf{M}^B for second layer, below is the receptive field of pixel 34,

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 3 – Receptive field of output pixel-34 when Mask-A is applied to layer-1 and Mask-B is applied to Layer-2.

3. If we use \mathbf{M}^B for the first layer and \mathbf{M}^A for the second layer. Using \mathbf{M}^B for the first layer and \mathbf{M}^A for second layer, below is the receptive field of pixel 34,

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 4 – Receptive field of output pixel-34 when Mask-B is applied to layer-1 and Mask-A is applied to Layer-2.

4. If we use \mathbf{M}^B for the first layer and \mathbf{M}^B for the second layer.
Using \mathbf{M}^B for the first layer and second layer, below is the receptive field of pixel 34,

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 5 – Receptive field of output pixel-34 when Mask-B is applied to layer-1 and Layer-2.

Question 2 (6-3-6-3). Reparameterization trick is a standard technique that makes the samples of a random variable differentiable. The trick represents the random variable as a simple mapping from another random variable drawn from some simple distribution². If the reparameterization is a bijective function, the induced density of the resulting random variable can be computed using the change-of-variable density formula, whose computation requires evaluating the determinant of the Jacobian of the mapping.

Consider a random vector $Z \in \mathbb{R}^K$ with a density function $q(\mathbf{z}; \phi)$ and a random variable $Z_0 \in \mathbb{R}^K$ having a ϕ -independent density function $q(\mathbf{z}_0)$. We want to find a deterministic function $\mathbf{g} : \mathbb{R}^K \rightarrow \mathbb{R}^K$ that depends on ϕ , to transform Z_0 , such that the induced distribution of the transformation

2. More specifically, these mapping should be differentiable wrt the density function's parameters.

has the same density as Z . Recall the change of density for a bijective, differentiable \mathbf{g} :

$$q(\mathbf{g}(\mathbf{z}_0)) = q(\mathbf{z}_0) |\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)|^{-1} = q(\mathbf{z}_0) \left| \det \left(\frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} \quad (1)$$

1. Assume $q(\mathbf{z}_0) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and $\mathbf{g}(\mathbf{z}_0) = \mu + \sigma \odot \mathbf{z}_0$, where $\mu \in \mathbb{R}^K$ and $\sigma \in \mathbb{R}_{>0}^K$. Note that \odot is element-wise product.

$$\begin{aligned} \mathbf{g}(\mathbf{z}_0) = \mu + \sigma \odot \mathbf{z}_0 &\implies \left| \det \left(\frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} = |\text{diag}(\sigma)|^{-1} \\ &= \frac{1}{\sqrt{\text{diag}(\sigma^2)}} \end{aligned}$$

$$\begin{aligned} \mathbf{g}(\mathbf{z}_0) = \mu + \sigma \odot \mathbf{z}_0 &\implies \mathbf{z}_0 = \frac{\mathbf{g}(\mathbf{z}_0) - \mu}{\sigma} \quad (\text{division with } \sigma \text{ is element wise}) \\ q(\mathbf{z}_0) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \mathbf{z}_0^T \mathbf{z}_0\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\mathbf{g}(\mathbf{z}_0) - \mu}{\sigma}\right)^T \left(\frac{\mathbf{g}(\mathbf{z}_0) - \mu}{\sigma}\right)\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (\mathbf{g}(\mathbf{z}_0) - \mu)^T (\text{diag}(\sigma^2)^{-1}) (\mathbf{g}(\mathbf{z}_0) - \mu)\right) \\ \implies q(\mathbf{g}(\mathbf{z}_0)) &= q(\mathbf{z}_0) \left| \det \left(\frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (\mathbf{g}(\mathbf{z}_0) - \mu)^T (\text{diag}(\sigma^2)^{-1}) (\mathbf{g}(\mathbf{z}_0) - \mu)\right) \frac{1}{\sqrt{\text{diag}(\sigma^2)}} \\ \implies q(\mathbf{g}(\mathbf{z}_0)) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} (\mathbf{g}(\mathbf{z}_0) - \mu)^T (\text{diag}(\sigma^2)^{-1}) (\mathbf{g}(\mathbf{z}_0) - \mu)\right) = \mathcal{N}(\mu, \text{diag}(\sigma^2)) \end{aligned}$$

2. Compute the time complexity of evaluating $|\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)|$ when $\mathbf{g}(\mathbf{z}_0) = \mu + \sigma \odot \mathbf{z}_0$. Use the big \mathcal{O} notation and expressive the time complexity as a function of K .

In this case, $|\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)| = \sigma^K$, computing power of a number is of log complexity. Hence, $|\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)|$ can be calculated in $\mathcal{O}(\log K)$.

3. Assume $\mathbf{g}(\mathbf{z}_0) = \mu + \mathbf{S} \mathbf{z}_0$, where \mathbf{S} is a non-singular $K \times K$ matrix. Derive the density of $\mathbf{g}(\mathbf{z}_0)$ using Equation (1).

$$\mathbf{g}(\mathbf{z}_0) = \mu + \mathbf{S} \mathbf{z}_0 \implies \left| \det \left(\frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} \left| \det \left(\frac{\partial (\mu + \mathbf{S} \mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} = |\mathbf{S}|^{-1} = \frac{1}{|\mathbf{S}|}$$

$$\begin{aligned}
\mathbf{g}(\mathbf{z}_0) &= \mu + \mathbf{S}\mathbf{z}_0 \implies \mathbf{z}_0 = \mathbf{S}^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu) \\
q(\mathbf{z}_0) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\mathbf{z}_0^T \mathbf{z}_0\right) \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}((\mathbf{S}^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu)))^T (\mathbf{S}^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu))\right) \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}((\mathbf{g}(\mathbf{z}_0) - \mu)^T (\mathbf{S}^{-1})^T \mathbf{S}^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu))\right) \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}((\mathbf{g}(\mathbf{z}_0) - \mu)^T (\mathbf{S}\mathbf{S}^T)^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu))\right) \\
\implies q(\mathbf{g}(\mathbf{z}_0)) &= q(\mathbf{z}_0) \left| \det \left(\frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0} \right) \right|^{-1} \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\mathbf{g}(\mathbf{z}_0) - \mu)^T (\mathbf{S}\mathbf{S}^T)^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu)\right) \frac{1}{|\mathbf{S}|} \\
\implies q(\mathbf{g}(\mathbf{z}_0)) &= \frac{1}{\sqrt{2\pi|\mathbf{S}\mathbf{S}^T|}} \exp\left(-\frac{1}{2}(\mathbf{g}(\mathbf{z}_0) - \mu)^T (\mathbf{S}\mathbf{S}^T)^{-1}(\mathbf{g}(\mathbf{z}_0) - \mu)\right) = \mathcal{N}(\mu, \mathbf{S}\mathbf{S}^T)
\end{aligned}$$

4. The time complexity of the general Jacobian determinant is at least $\mathcal{O}(K^{2.373})^3$. Assume instead $\mathbf{g}(\mathbf{z}_0) = \mu + \mathbf{S}\mathbf{z}_0$ with \mathbf{S} being a $K \times K$ lower triangular matrix; i.e. $\mathbf{S}_{ij} = 0$ for $j > i$, and $\mathbf{S}_{ii} > 0$. What is the time complexity of evaluating $|\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)|$?

Here, \mathbf{S} is a lower-triangular matrix, hence \mathbf{S}^{-1} is a lower triangular matrix too. So $\mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)$ is a lower-triangular matrix too. Computing determinant of a lower triangular matrix is multiplication of all the diagonal elements. There are total K such elements, hence the time complexity of evaluating $|\det \mathbf{J}_{\mathbf{z}_0} \mathbf{g}(\mathbf{z}_0)|$ is $\mathcal{O}(K)$.

Question 3 (5-5-6). Consider a latent variable model $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$, where $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and $\mathbf{z} \in \mathbb{R}^K$. The encoder network (aka “recognition model”) of variational autoencoder, $q_\phi(\mathbf{z}|\mathbf{x})$, is used to produce an approximate (variational) posterior distribution over latent variables \mathbf{z} for any input datapoint \mathbf{x} .⁴ This distribution is trained to match the true posterior by maximizing the evidence lower bound (ELBO):

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

Let \mathcal{Q} be the family of variational distributions with a feasible set of parameters \mathcal{P} ; i.e. $\mathcal{Q} = \{q(\mathbf{z}; \pi) : \pi \in \mathcal{P}\}$; for example π can be mean and standard deviation of a normal distribution. We assume q_ϕ is parameterized by a neural network (with parameters ϕ) that outputs the parameters, $\pi_\phi(\mathbf{x})$, of the distribution $q \in \mathcal{Q}$, i.e. $q_\phi(\mathbf{z}|\mathbf{x}) := q(\mathbf{z}; \pi_\phi(\mathbf{x}))$.

1. Show that maximizing the expected complete data log likelihood (ECLL)

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

for a fixed $q(\mathbf{z}|\mathbf{x})$, wrt the model parameter θ , is equivalent to maximizing

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$$

3. https://en.wikipedia.org/wiki/Computational_complexity_of_mathematical_operations

4. Using a recognition model in this way is known as “amortized inference”; this can be contrasted with traditional variational inference approaches (see, e.g., Chapter 10 of Bishop’s *Pattern Recognition and Machine Learning*), which fit a variational posterior independently for each new datapoint.

This means the maximizer of the ECLL coincides with that of the marginal likelihood only if $q(\mathbf{z}|\mathbf{x})$ perfectly matches $p(\mathbf{z}|\mathbf{x})$.

$$\begin{aligned}
\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})] &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}, \mathbf{z})] \\
&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{z}|\mathbf{x})p(\mathbf{x})] \\
&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x})] \\
&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{z}|\mathbf{x})] + \log p_\theta(\mathbf{x}) \quad (\text{Because } \log p_\theta(\mathbf{x}) \text{ is not dependent on } \mathbf{z}) \\
&= \log p_\theta(\mathbf{x}) + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log \frac{p_\theta(\mathbf{z}|\mathbf{x})q(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})}] \\
&= \log p_\theta(\mathbf{x}) + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})}] + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log q(\mathbf{z}|\mathbf{x})] \\
&= \log p_\theta(\mathbf{x}) - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})}] + \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log q(\mathbf{z}|\mathbf{x})] \\
& \quad (\text{Since } \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log q(\mathbf{z}|\mathbf{x})] \text{ is constant with respect to } \theta, \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log q(\mathbf{z}|\mathbf{x})] = C) \\
&= \log p_\theta(\mathbf{x}) - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})}] + C \\
&= \log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) + C
\end{aligned}$$

Hence,

$$\boxed{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})] = \log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) + C}$$

From the above, we can say that maximizing ECLL with respect to θ is equivalent to maximizing $\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$.

2. Consider a finite training set $\{\mathbf{x}_i : i \in \{1, \dots, n\}\}$, n being the size the training data. Let ϕ^* be the maximizer $\arg \max_\phi \sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$ with θ fixed. In addition, for each \mathbf{x}_i let $q_i \in \mathcal{Q}$ be an “instance-dependent” variational distribution, and denote by q_i^* the maximizer of the corresponding ELBO. Compare $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$ and $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$. Which one is bigger?

Here, θ is fixed. Hence, maximizing ELBO \implies minimizing $D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$ For any instance \mathbf{x}_i , $q_i^*(\mathbf{z})$ would maximize the ELBO or minimize the $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$. Hence, $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$ is smaller or equal to $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$. In best case, q_{ϕ^*} will be same as q_i^* , that's when $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$ and $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$ will be equal. So the relation is as below,

$$\boxed{D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i)) \geq D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))}$$

3. Following the previous question, compare the two approaches in the second subquestion
 - (a) in terms of bias of estimating the marginal likelihood via the ELBO, in the best case scenario (i.e. when both approaches are optimal within the respective families) When we have optimal encoding mechanism for both the approaches, the bias for $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$ is expected to be lower than the bias for $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$ based on the relation in subquestion 2, the conclusion would also be same as subquestion 2.

$$\boxed{D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i)) \geq D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))}$$

- (b) from the computational point of view (efficiency):

In terms of computation, per iteration cost would be the same as we have same number of

parameters. However, required number of iterations will be more when we have instance-specific variational distribution. Ideally, we will have to train n times more parameters for instance specific variational distribution over instance-agnostic distribution(q_ϕ).

(c) in terms of memory (storage of parameters):

Instance specific method would require n times more memory than instance-agnostic method, as we are storing 1 distribution for each example and we have n training examples. So instance-specific variational optimization is expensive in terms of storage requirements as well as optimization.

Question 4 (8-8). Let $p(x, z)$ be the joint probability of a latent variable model where x and z denote the observed and unobserved variables, respectively. Let $q(z|x)$ be an auxiliary distribution which we call the *proposal*, and define⁵

$$\mathcal{L}_K[q(z|x)] = \int \cdots \int \left(q(z_1|x) \cdots q(z_K|x) \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} \right) dz_1 dz_2 \cdots dz_K$$

We've seen in class that this objective is a tighter lower bound on $\log p(x)$ than the evidence lower bound (ELBO), which is equal to \mathcal{L}_1 ; that is $\mathcal{L}_1[q(z|x)] \leq \mathcal{L}_K[q(z|x)] \leq \log p(x)$.

In fact, $\mathcal{L}_K[q(z|x)]$ can be interpreted as the ELBO with a refined proposal distribution. For z_j drawn i.i.d. from $q(z|x)$ with $2 \leq j \leq K$, define the *unnormalized* density

$$\tilde{q}(z|x, z_2, \dots, z_K) := \frac{p(x, z)}{\frac{1}{K} \left(\frac{p(x, z)}{q(z|x)} + \sum_{j=2}^K \frac{p(x, z_j)}{q(z_j|x)} \right)}$$

(Hint: in what follows, you might need to use the fact that if w_1, \dots, w_K are random variables that have the same distribution, then $K\mathbb{E}[w_1] = \sum_i \mathbb{E}[w_i] = \mathbb{E}[\sum_i w_i]$. You need to identify such w_i 's before applying this fact for each subquestion.)

1. Show that $\mathcal{L}_K[q(z|x)] = \mathbb{E}_{z_{2:K}}[\mathcal{L}_1[\tilde{q}(z|x, z_2, \dots, z_K)]]$; that is, the importance-weighted lower bound with K samples is equal to the average ELBO with the unnormalized density as a refined proposal.

$\mathbb{E}_{z_{2:K}}[\mathcal{L}_1[\tilde{q}(z|x, z_2, \dots, z_K)]]$ is,

$$= \int \cdots \int q(z_2|x) \cdots q(z_K|x) \left(\int \tilde{q}(z|x, z_2, \dots, z_K) \log \frac{p(x, z_1)}{\tilde{q}(z|x, z_2, \dots, z_K)} dz_1 \right) dz_2 dz_3 \cdots dz_K$$

Substituting $\tilde{q}(z|x, z_2, \dots, z_K) := \frac{p(x, z)}{\frac{1}{K} \left(\frac{p(x, z)}{q(z|x)} + \sum_{j=2}^K \frac{p(x, z_j)}{q(z_j|x)} \right)}$ and multiplying with $q(z_1|x)$,

$$\begin{aligned} &= \int \cdots \int q(z_1|x) q(z_2|x) \cdots q(z_K|x) \frac{\frac{p(x, z_1)}{q(z_1|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} dz_1 dz_2 \cdots dz_K \\ &= \mathbb{E}_{z_{1:K}} \left[\frac{\frac{p(x, z_1)}{q(z_1|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} \right] \end{aligned}$$

5. Note that $\mathcal{L}_K[\cdot]$ is a "functional" whose input argument is a "function" $q(\cdot|x)$.

Using $K\mathbb{E}[w_1] = \sum_i \mathbb{E}[w_i] = \mathbb{E}[\sum_i w_i]$ in the above equation,

$$= \frac{1}{K} \sum_{i=1}^K \mathbb{E}_{z_{1:K}} \left[\frac{\frac{p(x, z_i)}{q(z_i|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} \right] \text{ because } z_i\text{s are IID}$$

Let's take summation inside the expectation,

$$\begin{aligned} &= \mathbb{E}_{z_{1:K}} \left[\frac{\frac{1}{K} \sum_{i=1}^K \frac{p(x, z_i)}{q(z_i|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} \right] \\ &= \int \cdots \int \left(q(z_1|x) \cdots q(z_K|x) \log \frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)} \right) dz_1 dz_2 \cdots dz_K \\ &= \mathcal{L}_K[q(z|x)] \\ &\implies \boxed{\mathcal{L}_K[q(z|x)] = \mathbb{E}_{z_{2:K}} [\mathcal{L}_1[\tilde{q}(z|x, z_2, \dots, z_K)]]} \end{aligned}$$

2. Show that $q_K(z|x) := \mathbb{E}_{z_{2:K}} [\tilde{q}(z|x, z_2, \dots, z_K)]$ is in fact a probability density function. Also, show that $\mathcal{L}_1[q_K(z|x)]$ is an even tighter lower bound than $\mathcal{L}_K[q(z|x)]$. This implies $q_K(z|x)$ is closer to the true posterior $p(z|x)$ than $q(z|x)$ due to resampling, since $\mathcal{L}_K[q(z|x)] \geq \mathcal{L}_1[q(z|x)]$. (Hint: $f(x) := -x \log x$ is concave.)

$$\begin{aligned} &\mathbb{E}_{z_{2:K}} [\tilde{q}(z|x, z_2, \dots, z_K)] \text{ is,} \\ &= \int \cdots \int q(z_2|x) \cdots q(z_K|x) \tilde{q}(z|x, z_2, \dots, z_K) dz_2 dz_3 \cdots dz_K \\ &= \int \cdots \int q(z_2|x) \cdots q(z_K|x) \frac{p(x, z)}{\frac{1}{K} \left(\frac{p(x, z)}{q(z|x)} + \sum_{j=2}^K \frac{p(x, z_j)}{q(z_j|x)} \right)} dz_2 dz_3 \cdots dz_K \end{aligned}$$

To prove that $\mathbb{E}_{z_{2:K}} [\tilde{q}(z|x, z_2, \dots, z_K)]$ is a probability distribution in z , let's integrate it in z space.

$$\begin{aligned} &\int_{z_1} \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] dz_1 \\ &= \int_{z_1} \left(\int \cdots \int q(z_2|x) \cdots q(z_K|x) \frac{p(x, z_1)}{\frac{1}{K} \left(\frac{p(x, z_1)}{q(z_1|x)} + \sum_{j=2}^K \frac{p(x, z_j)}{q(z_j|x)} \right)} dz_2 dz_3 \cdots dz_K \right) dz_1 \\ &= \int \cdots \int q(z_1|x) \cdots q(z_K|x) \frac{\frac{p(x, z_1)}{q(z_1|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} dz_1 dz_2 \cdots dz_K \\ &= \mathbb{E}_{z_{1:K}} \left[\frac{\frac{p(x, z_1)}{q(z_1|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \right] \end{aligned}$$

Using $K\mathbb{E}[w_1] = \sum_i \mathbb{E}[w_i] = \mathbb{E}[\sum_i w_i]$ in the above equation,

$$= \frac{1}{K} \sum_{i=1}^K \mathbb{E}_{z_{1:K}} \left[\frac{\frac{p(x, z_i)}{q(z_i|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \right]$$

Taking summation inside expectation,

$$\begin{aligned}
 &= \mathbb{E}_{z_{1:K}} \left[\frac{\frac{1}{K} \sum_{i=1}^K \frac{p(x, z_i)}{q(z_i|x)}}{\frac{1}{K} \sum_{j=1}^K \frac{p(x, z_j)}{q(z_j|x)}} \right] = \mathbb{E}_{z_{1:K}} [1] = 1 \\
 &\implies \boxed{\int_{z_1} \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] dz_1 = \int_z q_K(z|x) dz = 1}
 \end{aligned}$$

Hence, we can say that $q_K(z|x)$ is a probability distribution.

Here, $f(x) = -x \log x$ is a concave function, from Jensen's inequality, it implies that $f(\mathbb{E}(x)) \geq \mathbb{E}[f(x)]$. We will use this property. From answer of the previous question, we know that $\mathcal{L}_K[q(z|x)] = \mathbb{E}_{z_{2:K}} [\mathcal{L}_1[\tilde{q}(z|x, z_2, \dots, z_K)]]$. $\mathcal{L}_K[q(z|x)] =$

$$\begin{aligned}
 &= \mathbb{E}_{z_{2:K}} [\mathcal{L}_1[\tilde{q}(z|x, z_2, \dots, z_K)]] \\
 &= \mathbb{E}_{z_{2:K}} \left[\int \tilde{q}(z_1|x, z_2, \dots, z_K) \log \frac{p(x, z_1)}{\tilde{q}(z_1|x, z_2, \dots, z_K)} dz_1 \right] \\
 &= \mathbb{E}_{z_{2:K}} \left[\int \tilde{q}(z_1|x, z_2, \dots, z_K) \log p(x, z_1) dz_1 \right] - \mathbb{E}_{z_{2:K}} \left[\int \tilde{q}(z_1|x, z_2, \dots, z_K) \log \tilde{q}(z_1|x, z_2, \dots, z_K) dz_1 \right] \\
 &\leq \int \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] \log p(x, z_1) dz_1 - \int \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] \log \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] dz_1 \\
 &\leq \int \mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)] \frac{\log p(x, z_1)}{\mathbb{E}_{z_{2:K}} [\tilde{q}(z_1|x, z_2, \dots, z_K)]} dz_1 \\
 &\leq \int q_K(z|x) \frac{\log p(x, z)}{q_K(z|x)} dz = \mathcal{L}_1[q_K(z|x)] \\
 &\implies \boxed{\mathcal{L}_K[q(z|x)] \leq \mathcal{L}_1[q_K(z|x)]}
 \end{aligned}$$

Thus, we can say that $\mathcal{L}_1[q_K(z|x)]$ is a tighter bound than $\mathcal{L}_K[q(z|x)]$.

Question 5 (5-5-5-6). Normalizing flows are expressive invertible transformations of probability distributions. In this exercise, we will see how to satisfy the invertibility constraint of some family of parameterizations. For the first 3 questions, we assume the function $g : \mathbb{R} \rightarrow \mathbb{R}$ maps from real space to real space.

1. Let $g(z) = af(bz + c)$ where f is the ReLU activation function $f(x) = \max(0, x)$. Show that g is non-invertible.

$$g(z) = \begin{cases} a(bz + c) & \text{if } z > \frac{-c}{b} \\ 0 & \text{otherwise} \end{cases} \quad \frac{dg(z)}{dz} = \begin{cases} ab & \text{if } z > \frac{-c}{b} \\ 0 & \text{otherwise} \end{cases}$$

Since $\frac{dg(z)}{dz}$ is 0 for $z \leq \frac{-c}{b}$, $g(z)$ is neither monotonically increasing or decreasing over its domain. Hence $g(z)$ is non-invertible function.

2. Let $g(z) = \sigma^{-1}(\sum_{i=1}^N w_i \sigma(a_i z + b_i))$, $0 < w_i < 1$, where $\sum_i w_i = 1$, $a_i > 0$, and $\sigma(x) = 1/(1 + \exp(-x))$ is the logistic sigmoid activation function and σ^{-1} is its inverse. Show that g is *strictly monotonically increasing* on its domain $(-\infty, \infty)$, which implies invertibility.

Let $y = h(z) = \sum_{i=1}^N w_i \sigma(a_i z + b_i)$, since $\sum_i w_i = 1$ and $\sigma(a_i z + b_i) \in [0, 1]$, for $z \in [-\infty, \infty]$,

$\sum_{i=1}^N w_i \sigma(a_i z + b_i) \in [0, 1]$. Calculating σ^{-1} ,

$$t = \sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$1 + \exp(-x) = \frac{1}{t}$$

$$\exp(-x) = \frac{1-t}{t}$$

$$-x = \log \frac{1-t}{t}$$

$$\boxed{\sigma^{-1}(t) = \log \frac{t}{1-t}}$$

$$\frac{d\sigma^{-1}(t)}{dt} = \frac{d \log \frac{t}{1-t}}{dt}$$

$$= \frac{d \log t}{dt} - \frac{d \log(1-t)}{dt}$$

$$= \frac{1}{t} + \frac{1}{1-t}$$

$$\frac{d\sigma^{-1}(t)}{dt} = \frac{1}{t(1-t)} \text{ (we will use this ahead.)}$$

$$\frac{dg(z)}{dz} = \frac{d\sigma^{-1}(h(z))}{dz}$$

$$= \frac{1}{h(z)(1-h(z))} \frac{dh(z)}{dz}$$

$$= \frac{1}{(\sum_{i=1}^N w_i \sigma(a_i z + b_i))(1 - (\sum_{i=1}^N w_i \sigma(a_i z + b_i)))} \frac{d}{dz} \left(\sum_{i=1}^N w_i \sigma(a_i z + b_i) \right)$$

$$= \frac{\sum_{i=1}^N w_i a_i \sigma(a_i z + b_i) (1 - \sigma(a_i z + b_i))}{\sum_{i=1}^N w_i \sigma(a_i z + b_i) (1 - (\sum_{i=1}^N w_i \sigma(a_i z + b_i)))}$$

$$\Rightarrow \boxed{\frac{dg(z)}{dz} = \frac{\sum_{i=1}^N w_i a_i \sigma(a_i z + b_i) (1 - \sigma(a_i z + b_i))}{\sum_{i=1}^N w_i \sigma(a_i z + b_i) (1 - (\sum_{i=1}^N w_i \sigma(a_i z + b_i)))}}$$

Since $a_i > 0$, we can say that $\frac{dg(z)}{dz} > 0$, which implies that $g(z)$ is monotonically increasing function. Implies that $g(z)$ is invertible!

3. Consider a residual function of the form $g(z) = z + f(z)$. Show that $df/dz > -1$ implies g is

invertible.

$$\begin{aligned}\frac{df(z)}{dz} &> -1 \\ \frac{df(z)}{dz} + 1 &> 0 \\ g(z) &= z + f(z) \\ \frac{dg(z)}{dz} &= \frac{df(z)}{dz} + 1 \\ \implies \boxed{\frac{dg(z)}{dz} > 0}\end{aligned}$$

Since $\frac{dg(z)}{dz} > 0$, we can say that $g(z)$ is monotonically increasing function, hence it is invertible.

4. Consider the following transformation:

$$g(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0) \quad (2)$$

where $\mathbf{z}_0 \in \mathbb{R}^D$, $\alpha \in \mathbb{R}_{>0}$, $\beta \in \mathbb{R}$, and $r = \|\mathbf{z} - \mathbf{z}_0\|_2$, $h(\alpha, r) = 1/(\alpha + r)$. Consider the following decomposition of $\mathbf{z} = \mathbf{z}_0 + r\tilde{\mathbf{z}}$. (i) Given $\mathbf{y} = g(\mathbf{z})$, show that $\beta \geq -\alpha$ is a sufficient condition to derive the unique r from equation (2). (ii) Given r and \mathbf{y} , show that equation (2) has a unique solution $\tilde{\mathbf{z}}$. Here, we have $(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0)$

$$\begin{aligned}g(\mathbf{z}) &= \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0) \\ &= \mathbf{z}_0 + r\tilde{\mathbf{z}} + \frac{\beta r\tilde{\mathbf{z}}}{\alpha + r} \\ g(\mathbf{z}) - \mathbf{z}_0 &= r\left(\tilde{\mathbf{z}} + \frac{\beta\tilde{\mathbf{z}}}{\alpha + r}\right) \\ \text{Taking norm on both sides,} \\ |g(\mathbf{z}) - \mathbf{z}_0| &= r\left(1 + \frac{\beta}{\alpha + r}\right) = f(r)\end{aligned}$$

For $f(r)$ to be invertible, it has to be non-decreasing in nature. Meaning $\frac{df(r)}{dr} \geq 0$.

$$\begin{aligned}\frac{df(r)}{dr} &\geq 0 \\ &= 1 + \frac{\beta}{\alpha + r} - \frac{\beta r}{(\alpha + r)^2} \geq 0 \\ &= \frac{(\alpha + r + \beta)(\alpha + r) - r\beta}{(\alpha + r)^2} \geq 0 \\ &= \frac{(\alpha + r)^2 + \alpha\beta + \cancel{r\beta} - \cancel{r\beta}}{(\alpha + r)^2} \geq 0 \implies \frac{\alpha\beta}{(\alpha + r)^2} \geq -1 \implies \boxed{\beta \geq \frac{-(\alpha + r)^2}{\alpha}}\end{aligned}$$

Since $r \geq 0$ it is sufficient to imply that $\boxed{\beta \geq -\alpha}$

Hence, it's a sufficient condition. To get unique solution of $\tilde{\mathbf{z}}$

$$\begin{aligned}
 g(\mathbf{z}) &= v\mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0) \\
 &= \mathbf{z}_0 + r\tilde{\mathbf{z}} + \frac{\beta r\tilde{\mathbf{z}}}{\alpha + r} \\
 \implies g(\mathbf{z}) - \mathbf{z}_0 &= r\tilde{\mathbf{z}}\left(1 + \frac{\beta}{\alpha + r}\right) \\
 \implies \tilde{\mathbf{z}} &= \frac{\mathbf{y} - \mathbf{z}_0}{r\left(1 + \frac{\beta}{\alpha + r}\right)}
 \end{aligned}$$

The above equation can be uniquely solved give r and \mathbf{y} .

Question 6 (4-3-6). In this question, we are concerned with analyzing the training dynamics of GANs. Consider the following value function

$$V(d, g) = dg \quad (3)$$

with $g \in \mathbb{R}$ and $d \in \mathbb{R}$. We will use this simple example to study the training dynamics of GANs.

1. Consider gradient descent/ascent with learning rate α as the optimization procedure to iteratively minimize $V(d, g)$ w.r.t. g and maximize $V(d, g)$ w.r.t. d . We will apply the gradient descent/ascent to update g and d simultaneously. What is the update rule of g and d ? Write your answer in the following form

$$[d_{k+1}, g_{k+1}]^\top = A[d_k, g_k]^\top$$

where A is a 2×2 matrix; i.e. specify the value of A . The update would look like below,

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

2. The optimization procedure you found in 6.1 characterizes a map which has a stationary point⁶, what are the coordinates of the stationary points? For stationary points, we have to set partial derivative of $V(d, g)$ with respect to d and g to 0. $\frac{\partial V(d, g)}{\partial d} = 0 \implies g = 0$ and $\frac{\partial V(d, g)}{\partial g} = 0 \implies d = 0$. Hence the stationary point is $(g, d) = (0, 0)$.
3. Analyze the eigenvalues of A and predict what will happen to d and g as you update them jointly. In other word, predict the behaviour of d_k and g_k as $k \rightarrow \infty$.
Eigenvalues of A are $1 + \alpha i$ and $1 - \alpha i$. The norm of eigenvalue is $\sqrt{1 + \alpha^2}$. The norm of the eigenvalue is greater than 1. Which denotes that the vector $[d_k, g_k]^\top$ is scaling up as training progresses. So for this kind of simultaneous training, the variables being optimized are gonna explode. Hence, such a training procedure will lead to divergence and not appropriate for training such minimax objectives.

6. 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). Source: https://en.wikipedia.org/wiki/Stationary_point