

Due Date: April 5th 23:59, 2019

Instructions

- For all questions, show your work!
- Starred questions are **hard** questions, not **bonus** questions.
- 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 Shawn Tan, Samuel Lavoie, and Chin-Wei Huang.**

This assignment covers mathematical and algorithmic techniques underlying the three most popular families of deep generative models, variational autoencoders (VAEs, Questions 1-3), autoregressive models (Question 4), and generative adversarial networks (GANs, Questions 5-7).

Question 1 (8-8). Reparameterization trick is a standard technique that makes the samples of a random variable differentiable. Consider a random vector $Z \in \mathbb{R}^K$ with a density function $q(\mathbf{z}; \phi)$. We want to find a deterministic function $\mathbf{g} : \mathbb{R}^K \rightarrow \mathbb{R}^K$ that depends on ϕ , to transform a random variable Z_0 having a ϕ -independent density function $q(\mathbf{z}_0)$, such that $\mathbf{g}(Z_0)$ has the same density as Z . Recall the change of density for a bijective, differentiable \mathbf{g} :

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

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 $\sigma \in \mathbb{R}_{>0}^K$. Show that $\mathbf{g}(\mathbf{z}_0)$ is distributed by $\mathcal{N}(\mu, \text{diag}(\sigma^2))$ using Equation (1).
2. Assume instead $\mathbf{g}(\mathbf{z}_0) = \mu + \mathbf{S}\mathbf{z}_0$, where \mathbf{S} is a non-singular $K \times K$ matrix. Derive the density of $\mathbf{g}(\mathbf{z}_0)$ using Equation (1).

Answer 1.

1. Define $\mathbf{z} = \mathbf{g}(\mathbf{z}_0)$. Let \mathbf{J} be the Jacobian matrix $\mathbf{J} = \frac{\partial \mathbf{g}(\mathbf{z}_0)}{\partial \mathbf{z}_0}$. Note that here \mathbf{J} is a diagonal matrix with $\mathbf{J}_{ii} = \sigma_i$. Since \mathbf{g} has inverse $\mathbf{g}^{-1}(\mathbf{z}) = (\mathbf{z} - \mu)/\sigma$ (where the division is elementwise), we can write $\mathbf{z}_0 = \mathbf{g}^{-1}(\mathbf{g}(\mathbf{z}_0)) = (\mathbf{z} - \mu)/\sigma$ and $\|\mathbf{z}_0\|^2 = \sum_{i=1}^K \frac{(z_i - \mu_i)^2}{\sigma_i^2}$ where $\mathbf{z} = \mathbf{g}(\mathbf{z}_0)$:

$$q(\mathbf{z}) = (2\pi)^{-\frac{K}{2}} e^{-\frac{\|\mathbf{z}_0\|^2}{2}} \prod_{i=1}^K \sigma_i^{-1} = \prod_{i=1}^K \left(\frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2} \frac{(z_i - \mu_i)^2}{\sigma_i^2}} \right)$$

The right hand side is the density of $\mathcal{N}(\mu, \text{diag}(\sigma^2))$.

2. Here $\mathbf{J} = \mathbf{S}$. Let $\mathbf{S} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ be the singular decomposition. Let $\Sigma = \mathbf{S}\mathbf{S}^\top = \mathbf{U}\mathbf{D}^2\mathbf{U}^\top$. We have $\mathbf{g}^{-1}(\mathbf{z}) = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top(\mathbf{z} - \mu)$ and

$$\|\mathbf{z}_0\|^2 = \mathbf{g}^{-1}(\mathbf{z})^\top \mathbf{g}^{-1}(\mathbf{z}) = (\mathbf{z} - \mu)^\top \mathbf{U}\mathbf{D}^{-2}\mathbf{U}^\top(\mathbf{z} - \mu) = (\mathbf{z} - \mu)^\top \Sigma^{-1}(\mathbf{z} - \mu)$$

where $\mathbf{z} = \mathbf{g}(\mathbf{z}_0)$. Again, by the change of variable formula,

$$\begin{aligned} q(\mathbf{z}) &= (2\pi)^{-\frac{K}{2}} e^{-\frac{\|\mathbf{z}_0\|^2}{2}} |\det \mathbf{S}|^{-1} \\ &= (\det 2\pi \mathbf{I}_K)^{-\frac{1}{2}} |\det \mathbf{S}\mathbf{S}^\top|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{z} - \mu)^\top \Sigma^{-1}(\mathbf{z} - \mu)} \\ &= (\det 2\pi \Sigma)^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{z} - \mu)^\top \Sigma^{-1}(\mathbf{z} - \mu)} \end{aligned}$$

which is the density of $\mathcal{N}(\mu, \Sigma)$ where $\Sigma = \mathbf{S}\mathbf{S}^\top$.

Question 2 (5-5-6). Consider a latent variable model $\mathbf{z} \sim p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ where $\mathbf{z} \in \mathbb{R}^K$, and $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$. 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} .¹ 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

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

for a fixed $q(\mathbf{z}|\mathbf{x})$, wrt the model parameter θ , gives the maximizer of the biased log marginal likelihood: $\arg \max_\theta \{\log p_\theta(\mathbf{x}) + B(\theta)\}$, where $B(\theta)$ is non-positive. Find $B(\theta)$.

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 of $\sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$ 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?
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)
 - (b) from the computational point of view (efficiency)
 - (c) in terms of memory (storage of parameters) <https://www.overleaf.com/13215018tppysmgzxp>

Answer 2.

1. Using the fact that $q(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ do not depend on θ and the identity $1 = p_\theta(\mathbf{z}|\mathbf{x})/p_\theta(\mathbf{z}|\mathbf{x})$, we have

$$\begin{aligned} \arg \max_\theta \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \right\} &= \arg \max_\theta \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})} \right] \right\} \\ &= \arg \max_\theta \left\{ \log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x})) \right\} \end{aligned}$$

That is, maximizer of the expected complete data log likelihood aims at maximizing the marginal likelihood and reducing the KL divergence between $q(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{z}|\mathbf{x})$; the latter gap causes the bias.

2. Denote by $\mathcal{L}[q(\mathbf{z})]$ the lower bound corresponding to the variational distribution $q(\mathbf{z})$. Due to the finiteness of the recognition model, the distributions q_ϕ can represent is a subset of \mathcal{Q} . Given q_i^* the maximizer of the ELBO (argmax over the family \mathcal{Q}), the inference gap of q_{ϕ^*} can be decomposed as the sum of two positive terms:

$$\log p_\theta(\mathbf{x}_i) - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)] = (\log p_\theta(\mathbf{x}_i) - \mathcal{L}[q_i^*(\mathbf{z})]) + (\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)])$$

1. 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.

Using the equality $\log p_\theta(\mathbf{x}) - \mathcal{L}[q(\mathbf{z}|\mathbf{x})] = D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$, we have

$$D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i)) = D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i)) + (\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)])$$

The last term is non-negative since q_i^* is the maximizer among \mathcal{Q} , so the KL on the LHS is no smaller than the KL on the RHS.

3. (a) Since the only term in the ELBO that depends on θ is the expected complete data log likelihood, there is a bias due to the KL divergence. As demonstrated by the last question, the recognition model cannot do better than the optimal q_i^* .
- (b) Computationally, the point of using a recognition model is to amortize the cost of inference, so that one does not need to update q_i for each data point \mathbf{x}_i until the approximation is good enough; instead one can simply use the output of the recognition model as a reasonable approximation, since the encoder, which is normally jointly trained with the decoder, is constantly updated.
- (c) In terms of number of parameters, amortization allows for $\mathcal{O}(1)$ storage (considering a fixed size recognition model) instead of learning a new $q_i(\mathbf{z})$ for each \mathbf{x}_i , which is $\mathcal{O}(n)$.

Question 3 (6-6). Since variational inference provides a lower-bound on the log marginal likelihood of the data, it gives us a biased estimate of the marginal likelihood. Therefore, methods of “tightening” the bound (i.e. finding a higher valid lower bound) may be desirable.

Consider a latent variable model with the joint $p(\mathbf{x}, \mathbf{h})$ where \mathbf{x} and \mathbf{h} are the observed and unobserved random variables, respectively. Now let $q(\mathbf{h})$ be a variational approximation to $p(\mathbf{h}|\mathbf{x})$. Define

$$\mathcal{L}_K = \mathbb{E}_{\mathbf{h}_j \sim q(\mathbf{h})} \left[\log \frac{1}{K} \sum_{j=1}^K \frac{p(\mathbf{x}, \mathbf{h}_j)}{q(\mathbf{h}_j)} \right]$$

Note that \mathcal{L}_1 is equivalent to the evidence lower bound (ELBO).

1. Show that \mathcal{L}_K is a lower bound of the log marginal likelihood $\log p(\mathbf{x})$.
2. Show that $\mathcal{L}_K \geq \mathcal{L}_1$; i.e. \mathcal{L}_K is a family of lower bounds tighter than the ELBO.

Answer 3.

1. Applying Jensen’s inequality and linearity of expectation gives

$$\mathbb{E} \left[\log \frac{1}{K} \sum_{j=1}^K \frac{p_j}{q_j} \right] \leq \log \frac{1}{K} \sum_{j=1}^K \mathbb{E}_{q_j} \left[\frac{p_j}{q_j} \right] = \log \frac{1}{K} \sum_{j=1}^K \int p(\mathbf{x}, \mathbf{h}_j) d\mathbf{h}_j = \log p(\mathbf{x})$$

2. Note that $a = \frac{1}{K} \sum_{j=1}^K a_j$ can be viewed as the expectation over a uniform measure; i.e. $j \sim \text{Unif}([1, \dots, K])$ and $a = \mathbb{E}_j[a_j]$. Now we can apply Jensen’s Inequality to obtain the lower bound:

$$\mathbb{E} \left[\log \frac{1}{K} \sum_{j=1}^K \frac{p_j}{q_j} \right] = \mathbb{E} \left[\log \mathbb{E}_j \left[\frac{p_j}{q_j} \right] \right] \geq \mathbb{E} \left[\mathbb{E}_j \left[\log \frac{p_j}{q_j} \right] \right] = \mathbb{E} \left[\log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h})} \right]$$

since \mathbf{h}_j ’s are identically distributed.

Question 4 (5-5-5-5). One way to enforce autoregressive conditioning is via masking the weight parameters.² Consider a two-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

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

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) 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

FIGURE 1 – 5×5 convolutional feature map.

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.

Answer 4. See Figure 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

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

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 – Receptive field under different masking schemes.

Question 5 (10). Let P_0 and P_1 be two probability distributions with densities f_0 and f_1 (respectively). This problem demonstrates that an optimal GAN Discriminator (i.e. one which is able to distinguish between examples from P_0 and P_1 with minimal NLL loss) can be used to express the probability density of a datapoint \mathbf{x} under f_1 , $f_1(\mathbf{x})$ in terms of $f_0(\mathbf{x})$.

Assume f_0 and f_1 have the same support. Show that $f_1(\mathbf{x})$ can be estimated by $f_0(\mathbf{x})D(\mathbf{x})/(1 - D(\mathbf{x}))$ by establishing the identity $f_1(\mathbf{x}) = f_0(\mathbf{x})D^*(\mathbf{x})/(1 - D^*(\mathbf{x}))$, where

$$D^* := \arg \max_D \mathbb{E}_{\mathbf{x} \sim P_1} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P_0} [\log(1 - D(\mathbf{x}))]$$

2. An example of this is the use of masking in the Transformer architecture (Problem 3 of TP2 practical part).

Answer 5. The given function to be maximized can be expressed as a functional of D , $G[D] := \int g(D(\mathbf{x}), \mathbf{x}) d\mathbf{x}$ where

$$g(D(\mathbf{x}), \mathbf{x}) := f_1(\mathbf{x}) \log D(\mathbf{x}) + f_0(\mathbf{x}) \log(1 - D(\mathbf{x}))$$

Setting the functional derivative to be zero yields

$$\frac{\delta G[D]}{\delta D} = \frac{\partial g(D(\mathbf{x}), \mathbf{x})}{\partial D} = \frac{f_1(\mathbf{x})}{D(\mathbf{x})} - \frac{f_0(\mathbf{x})}{1 - D(\mathbf{x})} = 0$$

solving which gives us $D^*(\mathbf{x}) = \frac{f_1(\mathbf{x})}{f_0(\mathbf{x}) + f_1(\mathbf{x})}$. Thus we can use D^* to estimate the density of f_1 by rearranging the terms, which yields $f_1 = f_0 D^* / (1 - D^*)$.

Question 6 (5-5-6). While generative adversarial networks were originally formulated as minimizing the Jensen-Shannon (JS)-divergence, the framework can be generalized to use other divergences, such as the Kullback–Leibler (KL)-divergence. In this exercise we see how KL can be approximated (bounded from below) via a function $T : \mathcal{X} \rightarrow \mathbb{R}$ (i.e. the discriminator). Let q and p be probability density functions and recall the definition of the KL divergence $D_{\text{KL}}(p||q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$.

*1. Let $R_1[T] := \mathbb{E}_p[T(x)] - E_q[e^{T(x)-1}]$.

- The convex conjugate of a function $f(u)$ is defined as $f^*(t) = \sup_{u \in \text{dom} f} ut - f(u)$. Show that the convex conjugate of $f(u) = u \log u$ is $f^*(t) = e^{t-1}$, and its biconjugate³, i.e. the convex conjugate of its convex conjugate, is $f^{**}(u) := (f^*)^*(u) = u \log u$.
- Use the fact found above to show that $D_{\text{KL}}(p||q) = \sup_T R_1[T]$, where the supremum is taken over the set of all (measurable) functions $\mathcal{X} \rightarrow \mathbb{R}$. Start from the following step

$$\sup_{T \in \mathbb{R}} \int p(x)T(x) - q(x)e^{T(x)-1} dx = \int \sup_{t \in \mathbb{R}} p(x)t - q(x)e^{t-1} dx$$

which you don't need to prove.

*2. Let $r(x) = e^{T(x)} / \mathbb{E}_q[e^{T(x)}]$ and $R_2[T] := \mathbb{E}_p[T(x)] - \log \mathbb{E}_q[e^{T(x)}]$.

- Verify that $r q$ is a proper density function, i.e. integrating to 1.
 - Show that $D_{\text{KL}}(p||q) \geq R_2[T]$, with equality if and only if $T(x) = \log(p(x)/q(x)) + c$ where c is a constant independent of x .
3. Compare the two representations of the KL divergence. For fixed $T(x)$, $p(x)$ and $q(x)$, which one of $R_1[T]$ and $R_2[T]$ is greater than or equal to the other?

Answer 6.

- This problem can be solved using the functional derivative. Alternatively, we use the technique of *convex conjugate*. Since the stationary point of $ut - u \log u$ is $u^* = e^{t-1}$, so the maximizer is $f^*(t) = e^{t-1}$. Likewise, the stationary point of $tu - e^{t-1}$ is $1 + \log u$, so the maximizer is $f^{**}(u) = u \log u$. Now, rewrite $R_1[T]$ as

$$\begin{aligned} \sup_{T(x)} \int p(x)T(x) - q(x)e^{T(x)-1} dx \\ &= \int \sup_t p(x)t - q(x)e^{t-1} dx \\ &= \int q(x) \sup_t \left(\frac{p(x)}{q(x)} t - e^{t-1} \right) dx \end{aligned}$$

3. More generally, the biconjugate of f is equal to itself if f is a lower semi-continuous convex function (this is known as the **Fenchel-Monreau Theorem**).

Thus, setting $u = p(x)/q(x)$ in the above integral yields

$$\int q(x) \frac{p(x)}{q(x)} \log \frac{p(x)}{q(x)} dx = D_{\text{KL}}(p||q)$$

This is known as the **f-divergence** representation of the KL divergence.

2. By definition, $q(x)r(x)$ is a normalized density function, since

$$\int q(x)r(x)dx = \frac{1}{\mathbb{E}_q[e^{T(x)}]} \int q(x)e^{T(x)}dx = \frac{\mathbb{E}_q[e^{T(x)}]}{\mathbb{E}_q[e^{T(x)}]} = 1$$

Now write $R_2[T]$ as

$$\mathbb{E}_p[T(x)] - \log \mathbb{E}_q[e^{T(x)}] = \mathbb{E}_p[\log e^{T(x)}] - \log \mathbb{E}_q[e^{T(x)}] = \mathbb{E}_p[\log e^{T(x)} / \mathbb{E}_q[e^{T(x)}]]$$

which is equal to $\mathbb{E}_p[\log r(x)]$. The gap between the KL and the lower bound is $D_{\text{KL}}(p||q) - \mathbb{E}[\log r(x)] = \mathbb{E}_p[\log p(x) - \log q(x)r(x)]$. Since qr is a normalized density function, the gap is equal to $D_{\text{KL}}(p||qr)$, which is non-negative. Furthermore, the KL divergence is zero if and only if $p = qr$, which means $r(x) = p(x)/q(x)$, and that $e^{T(x)} \propto p(x)/q(x)$. As a result, the optimal $T(x)$ is equal to the log likelihood ratio $\log(p(x)/q(x))$ up to some constant c .

3. We can express the difference as

$$R_1 - R_2 = \log \mathbb{E}_q[e^{T(x)}] - \mathbb{E}_q[e^{T(x)-1}] = \log \mathbb{E}_q[e^{T(x)-1}] + 1 - \mathbb{E}_q[e^{T(x)-1}] \leq 0$$

since the elementary inequality $t - 1 - \log t \geq 0$ holds for any $t > 0$. Thus $R_1 \leq R_2$.

Question 7 (10). Let $q, p : \mathcal{X} \rightarrow [0, \infty)$ be probability density functions with disjoint (i.e. non-overlapping) support; more formally, $\{x \in \mathcal{X} : p(x) > 0 \text{ and } q(x) > 0\} = \emptyset$. What is the Jensen Shannon Divergence (JSD) between p and q ? Recall that JSD is defined as $D_{\text{JS}}(p||q) = \frac{1}{2}D_{\text{KL}}(p||r) + \frac{1}{2}D_{\text{KL}}(q||r)$ where $r(x) = \frac{p(x) + q(x)}{2}$.

Answer 7. By definition $D_{\text{JS}}(p||q) = \frac{1}{2}\mathbb{E}_p[\log p - \log r] + \frac{1}{2}\mathbb{E}_q[\log q - \log r]$. The first KL is

$$\mathbb{E}_p[\log p] - \mathbb{E}_p[\log(p+q)/2] = \mathbb{E}_p[\log p] - \mathbb{E}_p[\log(p+q) - \log 2] = \log 2$$

since p and q has non-overlapping support. Likewise the second KL is also $\log 2$, so $D_{\text{JS}}(p||q) = \log 2$.