

Homework 2 Solution and Grading Policy

TA: Kai-Chun Chen

Homework Policy: (READ BEFORE YOU START TO WORK)

- Copying from other students' solution is not allowed. If caught, all involved students get 0 point on that particular homework. Caught twice, you will be asked to drop the course.
- Collaboration is welcome. You can work together with **at most one partner** on the homework problems which you find difficult. However, you should write down your own solution, not just copying from your partner's.
- Your partner should be the same for the entire homework.
- Put your collaborator's name beside the problems that you collaborate on.
- When citing known results from the assigned references, be as clear as possible.

1. (Mixture of random processes) [14]

In this problem we look at different ways to generate mixtures of random processes, and the entropy rate of the mixture of random processes. Consider two stationary random processes $\{X_0[i] \mid i \in \mathbb{N}\}$ and $\{X_1[i] \mid i \in \mathbb{N}\}$ taking values in disjoint alphabets \mathcal{X}_0 and \mathcal{X}_1 respectively. The two processes are independent from each other, that is, $\{X_0[i]\} \perp\!\!\!\perp \{X_1[i]\}$, and they have entropy rates \mathcal{H}_0 and \mathcal{H}_1 respectively. Let $\{\Theta_i \mid i \in \mathbb{N}\}$ be a **stationary** Bernoulli random process, independent of everything else.

- Let $\Theta_i = \Theta$ for all $i \in \mathbb{N}$, where $\Theta \sim \text{Ber}(q)$. Is the random process $\{X_{\Theta_i}[i]\}$ stationary? What is its entropy rate? [6]
- Let $\{\Theta_i\}$ be Markov with a probability transition matrix

$$P_{\Theta_2|\Theta_1} = \begin{bmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{bmatrix}, \text{ for } \alpha, \beta \in (0, 1).$$

Suppose that both $\{X_0[i]\}$ and $\{X_1[i]\}$ are i.i.d. processes in this problem. Is the random process $\{X_{\Theta_i}[i]\}$ stationary? What is its entropy rate? [8]

Solution:

a) Since $\{X_0[i]\}$ and $\{X_1[i]\}$ are stationary, we have

$$\begin{aligned} & P_{X_{\Theta_1}[1], \dots, X_{\Theta_n}[n]} \\ &= (1-q)P_{X_0[1], \dots, X_0[n]} + qP_{X_1[1], \dots, X_1[n]} \\ &= (1-q)P_{X_0[l+1], \dots, X_0[n]} + qP_{X_1[1], \dots, X_1[l+n]} \quad (\text{by stationariness}) \\ &= P_{X_{\Theta_{l+1}}[l+1], \dots, X_{\Theta_{l+n}}[l+n]} \end{aligned}$$

By definition, $\{X_{\Theta_i}[i]\}$ is stationary.

Let $Y_i = X_{\Theta_i}[i]$.

$$\begin{aligned} \mathcal{H}(X_{\Theta_i}[i]) &= \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1}) \\ &= \lim_{n \rightarrow \infty} H(Y_n, \Theta_n | Y^{n-1}, \Theta^{n-1}) \quad (\text{since } \mathcal{X}_0 \text{ and } \mathcal{X}_1 \text{ are disjoint}) \\ &= \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1}, \Theta^n) + H(\Theta_n | Y^{n-1}, \Theta^{n-1}) \quad (\text{by chain rule}) \\ &= \lim_{n \rightarrow \infty} H(Y_n | \Theta_n) + H(\Theta_n | \Theta_{n-1}) \\ &= \lim_{n \rightarrow \infty} H(Y_n | \Theta) + H(\Theta | \Theta) \\ &= q\mathcal{H}_1 + (1-q)\mathcal{H}_0 + 0 \end{aligned}$$

b) Since we know in advance that $\{\Theta_i\}$ is stationary, we can conclude that $\Pr\{\Theta_1 = 0\} = \frac{\beta}{\alpha+\beta}$ and $\Pr\{\Theta_1 = 1\} = \frac{\alpha}{\alpha+\beta}$.

As a result, $\{X_{\Theta_i}[i]\}$ can be shown to be stationary by decomposing $P_{X_{\Theta_1}[1], \dots, X_{\Theta_n}[n]}$ and $P_{X_{\Theta_{l+1}}[l+1], \dots, X_{\Theta_{l+n}}[l+n]}$ simply by showing that $\Pr\{\Theta_n = 0\} = \frac{\beta}{\alpha+\beta}$ and $\Pr\{\Theta_n = 1\} = \frac{\alpha}{\alpha+\beta}$ for all n .

Next, let $Y_i = X_{\Theta_i}[i]$. $\mathcal{H}(\{X_{\Theta_i}[i]\}) = \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1})$.

$$\begin{aligned} & H(Y_n | Y^{n-1}) \\ &= H(Y_n, \Theta_n | Y^{n-1}) \quad (\mathcal{X}_0 \cap \mathcal{X}_1 = \emptyset) \\ &= H(Y_n | \Theta_n, Y^{n-1}) + H(\Theta_n | Y^{n-1}) \quad (\text{by chain rule}) \\ &= \Pr\{\Theta_n = 1\}H(X_1[n] | X_1^{n-1}) + \Pr\{\Theta_n = 0\}H(X_0[n] | X_0^{n-1}) + H(\Theta_n | Y^{n-1}) \\ &= \Pr\{\Theta_n = 1\}H(X_1[n] | X_1^{n-1}) + \Pr\{\Theta_n = 0\}H(X_0[n] | X_0^{n-1}) + H(\Theta_n | \Theta^{n-1}) \quad (\mathcal{X}_0 \cap \mathcal{X}_1 = \emptyset) \\ &= \frac{\alpha}{\alpha+\beta}\mathcal{H}_1 + \frac{\beta}{\alpha+\beta}\mathcal{H}_0 + H(\Theta_2 | \Theta_1) \\ &= \frac{\alpha}{\alpha+\beta}(\mathcal{H}_1 + H_b(\beta)) + \frac{\beta}{\alpha+\beta}(\mathcal{H}_0 + H_b(\alpha)) \end{aligned}$$

Grading Policy:

- a) Stationaryness and reason [2] calculation of entropy rate [4]
- b) Stationaryness and reason [2] calculation of entropy rate [6]

2. (Binary hypothesis testing) [16]

Let X_1, X_2, \dots be a sequence of i.i.d. Bernoulli p random variables, that is,

$$\Pr\{X_i = 1\} = 1 - \Pr\{X_i = 0\} = p.$$

Based on the observations so far, the goal is of a decision maker to determine which of the following two hypotheses is true:

$$\mathcal{H}_0 : p = p_0$$

$$\mathcal{H}_1 : p = p_1$$

where $0 < p_0 < p_1 \leq 1/2$.

- a) (Warm-up) Consider the problem of making the decision based on X_1 .

Draw the optimal $(\pi_{1|0}, \pi_{0|1})$ trade-off curve. [4]

- b) Suppose the decision maker waits until an 1 appears and makes the decision based on the whole observed sequence. Sketch the optimal $(\pi_{1|0}, \pi_{0|1})$ trade-off curve. [4]

- c) Now suppose the decision maker waits until in total n 1's appear and makes the decision based on the whole observed sequence. Let $\varpi_{0|1}^*(n, \epsilon)$ denote the minimum type-II error probability subject to the constraint that the type-I error probability is not greater than ϵ , $0 < \epsilon < 1$. Does $\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{\varpi_{0|1}^*(n, \epsilon)}$ exist? If so, find it. Otherwise, show that the limit does not exist. [8]

Solution:

- a) By Neyman-Pearson theorem, the optimal test is randomized LRT. So, by leveraging the parameters of the randomized LRT, that is, τ and γ , we can derive the optimal trade-off curve. Note that the likelihood ratio can only take two values: $\frac{p_1}{p_0}, \frac{1-p_1}{1-p_0}$. Therefore, discuss the range of τ we get

$$\begin{cases} \pi_{1|0} = 1, \pi_{0|1} = 0, & 0 \leq \tau < \frac{1-p_1}{1-p_0} \\ \pi_{1|0} = p_0 + \gamma(1-p_0), \pi_{0|1} = (1-\gamma)(1-p_1) = \frac{1-p_1}{1-p_0}(1-\pi_{1|0}), & \tau = \frac{1-p_1}{1-p_0} \\ \pi_{1|0} = p_0, \pi_{0|1} = 1-p_1, & \frac{1-p_1}{1-p_0} < \tau < \frac{p_1}{p_0} \\ \pi_{1|0} = \gamma p_0, \pi_{0|1} = (1-\gamma)p_1 + (1-p_1) = 1 - \frac{p_1}{p_0}\pi_{1|0}, & \tau = \frac{p_1}{p_0} \\ \pi_{1|0} = 0, \pi_{0|1} = 1, & \tau > \frac{p_1}{p_0}. \end{cases}$$

We can then draw the trade-off curve using the equations derived above.

- b) Note that our observation can only be 1, 01, 001, 0001, \dots , let L be the length of the observation, we have

$$\mathcal{H}_0 : L \sim \text{Geo}(p_0)$$

$$\mathcal{H}_1 : L \sim \text{Geo}(p_1)$$

Similar to a), we can discuss the range of τ and get:

$$\begin{cases} \pi_{1|0} = 0, \pi_{0|1} = 1, \tau > \frac{p_1}{p_0} \\ \pi_{1|0} = \sum_{i=1}^{n-1} (1-p_0)^{i-1} p_0 + \gamma (1-p_0)^{n-1} p_0, \\ \pi_{0|1} = \sum_{i=n+1}^{\infty} (1-p_1)^{i-1} p_1 + (1-\gamma)(1-p_1)^{n-1} p_1, \tau = \frac{(1-p_1)^{n-1} p_1}{(1-p_0)^{n-1} p_0} \\ \pi_{1|0} = \sum_{i=1}^n (1-p_0)^{i-1} p_0, \pi_{0|1} = \sum_{i=n+1}^{\infty} (1-p_1)^{i-1} p_1, \frac{(1-p_1)^{n-1} p_1}{(1-p_0)^{n-1} p_0} > \tau > \frac{(1-p_1)^n p_1}{(1-p_0)^n p_0}. \end{cases}$$

And we can draw the trade-off curve using the equations derived above.

- c) The observation can be viewed as n i.i.d. geometric random variables. To see this, for any realization of observation, insert a “—” symbol in front of the sequence, also insert a “—” right after a “1”. For example, if $n = 4$ and the realization is 010001101, we write it as |01|0001|1|01|. Apparently, the length of the subsequence between two — is a geometric random variable. Hence, in this subproblem, we are testing $\text{Geo}(p_0)^{\otimes n}$ and $\text{Geo}(p_1)^{\otimes n}$. By Chernoff-Stein lemma,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{\varpi_{0|1}^*(n, \epsilon)} = D(\text{Geo}(p_0) \parallel \text{Geo}(p_1)) = \log \frac{p_0}{p_1} + \left(\frac{1-p_0}{p_0} \right) \log \frac{1-p_0}{1-p_1}.$$

Grading Policy:

- a) Specify the trade-off curve [2] Argue the optimality [2]
- b) Specify the trade-off curve [2] Argue the optimality [2]
- c) Formulate the problem as a hypothesis testing with n instances [3], Chernoff-Stein lemma [2], and calculation [3]

3. (Mixture of information divergences) [8]

For m discrete probability distributions P_1, P_2, \dots, P_m with the same support \mathcal{X} , consider the following minimization problem:

$$\min_{Q \in \mathcal{P}(\mathcal{X})} \sum_{i=1}^m \lambda_i D(P_i \parallel Q),$$

where $\mathcal{P}(\mathcal{X})$ denotes the collection of probability distributions over \mathcal{X} , $\sum_{i=1}^m \lambda_i = 1$, and $\lambda_i > 0$ for $i = 1, 2, \dots, m$. Show that $\sum_{i=1}^m \lambda_i P_i$ is a minimizer to the above problem.

Solution:

Let $\bar{P} = \sum_{i=1}^m \lambda_i P_i$, we have $\forall Q \in \mathcal{P}(\mathcal{X})$, that

$$\sum_{i=1}^m \lambda_i D(P_i \parallel Q) - \sum_{i=1}^m \lambda_i D(P_i \parallel \bar{P})$$

$$\begin{aligned}
&= \sum_{i=1}^m \sum_{x \in \mathcal{X}} \lambda_i P_i(x) \log \frac{\bar{P}(x)}{Q(x)} \\
&= \sum_{x \in \mathcal{X}} \left(\sum_{i=1}^m \lambda_i P_i(x) \right) \log \frac{\bar{P}(x)}{Q(x)} \\
&= \sum_{x \in \mathcal{X}} \bar{P}(x) \log \frac{\bar{P}(x)}{Q(x)} \\
&= D(\bar{P} \| Q) \geq 0.
\end{aligned}$$

Hence, \bar{P} is a minimizer.

Grading Policy

Reasonable Procedure [4] Correctness [4] Per mistake [-1]

4. (Rényi's divergence) [12]

Alfréd Rényi introduced the following generalization of information divergence called *Rényi's divergence of order α* (for simplicity, only deal with the discrete case):

$$D_\alpha(P \| Q) := \frac{1}{\alpha - 1} \log \left(\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha} \right), \quad \alpha \in (0, 1) \cup (1, \infty),$$

where P, Q are both probability distributions over a finite alphabet \mathcal{X} , and $\text{supp } P \subseteq \text{supp } Q$.

- (Non-negativity) Show that $D_\alpha(P \| Q) \geq 0$, with equality if and only if $P = Q$. [4]
- (Relation with KL divergence) Show that $D_\alpha(P \| Q) \geq D(P \| Q)$ for $\alpha > 1$ and $D_\alpha(P \| Q) \leq D(P \| Q)$ for $\alpha < 1$. Furthermore, $\lim_{\alpha \rightarrow 1} D_\alpha(P \| Q) = D(P \| Q)$. [4]
- (Data processing) Show that $D_\alpha(P \| Q)$ satisfies the data processing inequality. [4]

Solution:

- We divide the value of α into two cases:

If $\alpha \in (0, 1)$, we can lower bound $D_\alpha(P \| Q)$ by Hölder's inequality, which states that for any $p, q > 1$ and $\frac{1}{p} + \frac{1}{q} = 1$, we have

$$\sum_{i=1}^n |x_i y_i| \leq \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \left(\sum_{i=1}^n |y_i|^q \right)^{\frac{1}{q}},$$

the equality holds iff $(|x_1|^p, \dots, |x_n|^p)$ and $(|y_1|^q, \dots, |y_n|^q)$ are linearly dependent.

$$\begin{aligned}
D_\alpha(P \| Q) &= \frac{1}{\alpha - 1} \log \left(\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha} \right) \\
&\geq \frac{1}{\alpha - 1} \log \left(\sum_{a \in \mathcal{X}} (P(a)^\alpha)^{\frac{1}{\alpha}} \right)^\alpha \left(\sum_{a \in \mathcal{X}} (Q(a)^{1-\alpha})^{\frac{1}{1-\alpha}} \right)^{1-\alpha}
\end{aligned}$$

$$\begin{aligned}
& \text{(by Hölder's inequality with } p = \frac{1}{\alpha}, q = \frac{1}{1-\alpha} \text{)} \\
& = 0
\end{aligned}$$

If $\alpha \in (1, \infty)$, we can bound it by Jensen's inequality,

$$\begin{aligned}
D_\alpha(P\|Q) &= \frac{1}{\alpha-1} \log \left(\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha} \right) \\
&= \frac{1}{\alpha-1} \log \left(\sum_{a \in \mathcal{X}} P(a) \left(\frac{P(a)}{Q(a)} \right)^{\alpha-1} \right) \\
&= \frac{1}{\alpha-1} \log E_{X \sim P} \left[\left(\frac{P(X)}{Q(X)} \right)^{\alpha-1} \right] \\
&\geq \frac{1}{\alpha-1} E_{X \sim P} \left[\log \left(\frac{P(X)}{Q(X)} \right)^{\alpha-1} \right] \\
&= E_{X \sim P} \left[\log \frac{P(X)}{Q(X)} \right] = D(P\|Q) \geq 0
\end{aligned}$$

Furthermore, both equality holds iff $P(a) = Q(a)$ for all $a \in \mathcal{X}$, which is $P = Q$.

- b) In a), we have proved the case when $\alpha > 1$. For the case when $\alpha < 1$, we can directly obtain the result by substituting the inequality since $\frac{1}{\alpha-1} < 0$ here. Thus, it suffices to show that $\lim_{\alpha \rightarrow 1} D_\alpha(P\|Q) = D(P\|Q)$.

$$\begin{aligned}
\lim_{\alpha \rightarrow 1} D_\alpha(P\|Q) &= \lim_{\alpha \rightarrow 1} \frac{\log \left(\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha} \right)}{\alpha-1} \\
&\stackrel{H}{=} \lim_{\alpha \rightarrow 1} \frac{\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha} \log \left(\frac{P(a)}{Q(a)} \right)}{\sum_{a \in \mathcal{X}} P(a)^\alpha Q(a)^{1-\alpha}} \\
&= \sum_{a \in \mathcal{X}} P(a) \frac{P(a)}{Q(a)} = D(P\|Q)
\end{aligned}$$

- c) Follow the notation proving data processing inequality used in lecture, we have $\forall x, y$,

$$\frac{P_{X,Y}(x, y)}{Q_{X,Y}(x, y)} = \frac{P_X(x) W_{Y|X}(y|x)}{Q_X(x) W_{Y|X}(y|x)} = \frac{P_X(x)}{Q_X(x)}$$

As a result,

$$\begin{aligned}
D_\alpha(P_{X,Y}\|Q_{X,Y}) &= \frac{1}{\alpha-1} \sum_{a \in \mathcal{X}} \sum_{b \in \mathcal{Y}} P_{X,Y}(a, b) \left(\frac{P_{X,Y}(a, b)}{Q_{X,Y}(a, b)} \right)^{\alpha-1} \\
&= \frac{1}{\alpha-1} \sum_{a \in \mathcal{X}} P_X(a) \left(\frac{P_X(a)}{Q_X(a)} \right)^{\alpha-1} = D_\alpha(P_X\|Q_X)
\end{aligned}$$

Therefore,

$$\begin{aligned}
 D_\alpha(P_X \| Q_X) &= D_\alpha(P_{X,Y} \| Q_{X,Y}) \\
 &= \frac{1}{\alpha - 1} \mathbb{E}_{X,Y \sim Q_{X,Y}} \left[\left(\frac{P_{X,Y}(X,Y)}{Q_{X,Y}(X,Y)} \right)^\alpha \right] \\
 &= \frac{1}{\alpha - 1} \mathbb{E}_{Y \sim Q_Y} \left[\mathbb{E}_{X \sim Q_{X|Y=Y}} \left[\left(\frac{P_{X,Y}(X,Y)}{Q_{X,Y}(X,Y)} \right)^\alpha \middle| Y \right] \right] \\
 &\geq \frac{1}{\alpha - 1} \mathbb{E}_{Y \sim Q_Y} \left[\mathbb{E}_{X \sim Q_{X|Y=Y}} \left[\left(\frac{P_{X,Y}(X,Y)}{Q_{X,Y}(X,Y)} \right) \middle| Y \right]^\alpha \right] \\
 &= \frac{1}{\alpha - 1} \mathbb{E}_{Y \sim Q_Y} \left[\left(\frac{P_Y(Y)}{Q_Y(Y)} \right)^\alpha \right] = D_\alpha(P_Y \| Q_Y)
 \end{aligned}$$

Grading Policy

- a) Proof of inequality [3] Equivalent condition of equality [1]
- b) Proof of inequalities [2] Calculation of limit [2]
- c) Correct proof [4] Wrong for some α [-1]