

## ASSIGNMENT 6

### Exercise 1. (Convexity)

- a. For distributions  $p$  and  $q$  on a finite alphabet, show that  $D(p||q)$  is convex in the pair  $(p, q)$ . *i.e.*, if  $(p_1, q_1)$  and  $(p_2, q_2)$  are two pairs of pmfs, then,

$$D(\lambda p_1 + (1 - \lambda)p_2 || \lambda q_1 + (1 - \lambda)q_2) \leq \lambda D(p_1 || q_1) + (1 - \lambda)D(p_2 || q_2),$$

for all  $0 \leq \lambda \leq 1$ .

*Hint* – Suppose that the alphabet size is  $m$ . Then, the left-side is a sum of  $m$  terms. Apply log-sum inequality to a particular term and sum over all  $m$  terms.

- b. For  $(X, Y) \sim p(x)p(y|x)$ , show that  $I(X; Y)$  is a convex function of  $p(y|x)$  for fixed  $p(x)$ .

*Hint* –

- i. Consider  $p_1(y|x)$  and  $p_2(y|x)$  and their convex combination  $p_\lambda(y|x) = \lambda p_1(y|x) + (1 - \lambda)p_2(y|x)$ .
- ii. Write out the joint distribution  $p_\lambda(x, y)$  and the marginal  $p_\lambda(y)$ .
- iii. Consider the KL divergence between  $p_\lambda(x, y)$  and  $p(x)p_\lambda(y)$ .

**Exercise 2.** (Converse to the rate distortion theorem) Recall that the rate distortion function is given by

$$R(D) = \min_{p(\hat{x}|x): \sum_{x, \hat{x}} p(x)p(\hat{x}|x)d(x, \hat{x}) \leq D} I(X; \hat{X}).$$

- a. Prove that  $R(D)$  is a non-increasing convex function of  $D$ .

*Hint* – To prove that  $R(D)$  is convex, consider two rate distortion pairs,  $(R_1, D_1)$  and  $(R_2, D_2)$ , which lie on the rate distortion curve. Let the joint distributions that achieve these pairs be  $p_1(x, \hat{x}) = p(x)p_1(\hat{x}|x)$  and  $p_2(x, \hat{x}) = p(x)p_2(\hat{x}|x)$ . Consider the distribution

$$p_\lambda = \lambda p_1 + (1 - \lambda)p_2.$$

Since the distortion is a linear function of the distribution, we have

$$D(p_\lambda) = \lambda D_1 + (1 - \lambda)D_2.$$

Also,

$$I_{p_\lambda}(X; \hat{X}) \leq \lambda I_{p_1}(X; \hat{X}) + (1 - \lambda)I_{p_2}(X; \hat{X}).$$

- b. Consider any  $(2^{nR}, n)$  rate distortion code defined by functions  $f_n$  and  $g_n$ . Let  $\hat{X}^n = g_n(f_n(X^n))$  be the reproduced sequence corresponding to  $X^n$ . Justify the steps with labels

on the equality or the inequality signs.

$$\begin{aligned}
I(X^n; \hat{X}^n) &= H(X^n) - H(X^n | \hat{X}^n) \\
&\stackrel{(a)}{=} \sum_{i=1}^n H(X_i) - H(X^n | \hat{X}^n) \\
&\stackrel{(b)}{=} \sum_{i=1}^n H(X_i) - \sum_{i=1}^n H(X_i | \hat{X}^n, X_{i-1}, \dots, X_1) \\
&\stackrel{(c)}{\geq} \sum_{i=1}^n H(X_i) - \sum_{i=1}^n H(X_i | \hat{X}^i) \\
&= \sum_{i=1}^n I(X_i; \hat{X}_i).
\end{aligned}$$

- c. Assume that the expected distortion  $\mathbb{E}d(X^n, \hat{X}^n) \leq D$  for this code. Justify the steps with labels on the equality or the inequality signs.

$$\begin{aligned}
\sum_{i=1}^n R(\mathbb{E}d(X_i, \hat{X}_i)) &= n \sum_{i=1}^n \frac{1}{n} R(\mathbb{E}d(X_i, \hat{X}_i)) \\
&\stackrel{(a)}{\geq} nR \left( \frac{1}{n} \sum_{i=1}^n \mathbb{E}d(X_i, \hat{X}_i) \right) \\
&\stackrel{(b)}{=} nR(\mathbb{E}d(X^n, \hat{X}^n)).
\end{aligned}$$

- d. Using the results above, show that  $R \geq R(D)$ . *Hint* – Start with  $nR \geq H(\hat{X}^n)$ .

**Exercise 3.** (Uniquely decodable codes) Given an alphabet  $\mathcal{X} = \{1, \dots, m\}$  and a probability distribution  $P = (p_1, \dots, p_m)$  on  $\mathcal{X}$ , solve (using Lagrange multipliers) the following convex optimization problem:

$$\min_{\ell_1, \dots, \ell_m \in \mathbb{R}} \sum_{i=1}^m p_i \ell_i \quad \text{subject to} \quad \sum_{i=1}^m 2^{-\ell_i} \leq 1.$$

Conclude that for a uniquely decodable code, the minimum expected codeword length is greater than or equal to  $H(P)$ . Why is it *greater than or equal to* and not *equal to*?

**Exercise 4.** (Rényi entropy) For a distribution  $P$  on a finite alphabet  $\mathcal{X}$ , the Rényi entropy of order  $\alpha$ ,  $H_\alpha(P)$  is given by

$$H_\alpha(P) = \frac{1}{1-\alpha} \log \left( \sum_{i=1}^m p_i^\alpha \right).$$

for  $\alpha \geq 0$ ,  $\alpha \neq 1$ .

- a. Show that the Shannon entropy  $H(P)$  satisfies

$$H(P) = \lim_{\alpha \rightarrow 1} H_\alpha(P).$$

*Hint* – Use L'Hôpital's rule.

b. For i.i.d. random variables  $X$  and  $Y$  on  $\mathcal{X}$ , what is  $\mathbb{P}[X = Y]$  in terms of  $H_2(P)$ ?

c. In the limit as  $\alpha \rightarrow \infty$ ,  $H_\alpha$  converges to  $H_\infty$  defined by

$$H_\infty(P) = -\log \max_i P_i.$$

Show that  $H_2 \leq 2H_\infty$ .

d. Show that for a fixed  $P$ ,

$$H_0 \geq H_1 \geq H_2 \geq H_\infty.$$

**Exercise 5.** (List decoding Fano's inequality) Consider discrete random variables  $X$  and  $Y$  with  $X$  taking values in the set  $\{0, 1\}^k$ . Upon observing  $Y$  we produce a list  $L(Y)$  of size  $2^\ell$  such that

$$\mathbb{P}(X \in L(Y)) \geq 1 - \varepsilon.$$

Show that

$$H(X|Y) \leq \varepsilon k + (1 - \varepsilon)\ell + 1.$$

*Hint* – Define the random variable  $T = \mathbf{1}_{\{X \in L(Y)\}}$ . Expand  $H(X, Y, T)$  using chain rule.

**Exercise 6.** (Shearer's lemma) Shearer's lemma is a generalization of the basic inequality

$$H(X_1, \dots, X_n) \leq \sum_{i=1}^n H(X_i).$$

For  $S \subseteq [n] = \{1, 2, \dots\}$ , we write  $X_S = (X_i : i \in S)$ .

a. Prove the lemma: Let  $X_1, \dots, X_n$  be random variables. Let  $S_1, \dots, S_m \subseteq [n]$  be subsets such that each  $i \in [n]$  belongs to at least  $k$  sets. Then,

$$kH(X_1, \dots, X_n) \leq \sum_{j=1}^m H(X_{S_j}).$$

*Hint* – Let  $S_j = \{i_1, \dots, i_{s_j}\}$  with  $i_1 < \dots < i_{s_j}$ . Then,

$$\begin{aligned} H(X_{S_j}) &= H(X_{i_1}) + H(X_{i_2}|X_{i_1}) + \dots + H(X_{i_{s_j}}|X_{i_1}, \dots, X_{i_{s_j-1}}) \\ &\geq H(X_{i_1}|X_1, \dots, X_{i_1-1}) + H(X_{i_2}|X_1, \dots, X_{i_2-1}) + \dots + H(X_{i_{s_j}}|X_1, \dots, X_{i_{s_j-1}}). \end{aligned}$$

Sum the left side over  $j = 1$  to  $m$ .

b. Suppose  $n$  distinct points in  $\mathbb{R}^3$  have  $n_1$  distinct projections on the  $XY$ -plane,  $n_2$  distinct projections on the  $XZ$ -plane, and  $n_3$  distinct projections on the  $YZ$ -plane. For two different points, since all three projections cannot be the same, we have  $n \leq n_1 n_2 n_3$ . Using Shearer's lemma, show that

$$n \leq \sqrt{n_1 n_2 n_3}.$$

*Hint* – Let  $P = (X_1, X_2, X_3)$  be one of the  $n$  points picked uniformly at random. Then,  $P_1 = (X_1, X_2)$ ,  $P_2 = (X_1, X_3)$ , and  $P_3 = (X_2, X_3)$  are its three projections.

**Exercise 7.** (Shotgun DNA sequencing)<sup>1</sup> DNA sequencing is the basic workhorse of modern day biology and medicine. Shotgun sequencing is the dominant technique used: many randomly located short fragments called reads are extracted from the DNA sequence, and these reads are assembled to reconstruct the original sequence. A basic question is: given a sequencing technology and the statistics of the DNA sequence, what is the minimum number of reads required for reliable reconstruction?

The DNA sequence  $s = s_1s_2 \cdots s_G$  is modeled as an i.i.d. random process of length  $G$  with each symbol taking values according to a probability distribution  $p = (p_1, p_2, p_3, p_4)$  on the nucleotide alphabet  $\{A, C, G, T\}$ . A read is a substring of length  $L$  from the DNA sequence. The objective of DNA sequencing is to reconstruct the whole sequence  $s$  based on  $N$  reads from the sequence. The starting location of each read is uniformly distributed on the DNA sequence and are independent from one read to another. We seek to understand the fundamental limits on the two quantities  $N$  and  $L$ .

- a. *Covering*: Argue that for the perfect reconstruction of  $s$ , for a fixed  $L$ , the collection of reads should *cover* the entire sequence and hence a necessary condition is that  $N \geq G/L$ .
- b. *An improvement via the coupon collector problem*: The well-known “coupon collector problem” is the following. Suppose we repeatedly and independently sample a random variable that is uniformly distributed over  $\{1, 2, \dots, n\}$ . How many samples do we need to ensure the sampling of all  $n$  numbers? The answer to this question is roughly  $n \log n$  ([https://en.wikipedia.org/wiki/Coupon\\_collector%27s\\_problem](https://en.wikipedia.org/wiki/Coupon_collector%27s_problem)).

Now, consider a modified DNA read technique where in each read, you get to observe  $L$  independent locations (instead of contiguous locations). Can you use the coupon collector result to get an estimate on the necessary number of reads  $N$  for this modified problem? What does it say about the required number of reads for the original problem?

- c. Suppose we have two DNA sequences, the first sequence generated by a uniform distribution on  $\{A, C, T, G\}$  and the second by a distribution  $(0.5, 0.4, 0.05, 0.05)$ . The DNA sequences and the corresponding reads are as follows:

- 1. Sequence:  $ACTGCATAGT$ , Reads:  $TGC, CAT, ACT, TAG, AGT$ .
- 2. Sequence:  $ACACATACGC$ , Reads:  $ACA, CAC, TAC, ACG, CGC$

*Impossible to reconstruct?*

- i. Which among the two sequences can you reconstruct (uniquely) from the reads? Why?
- ii. Calculate the Rényi entropy of order 2 (see Ex.4) for both the distributions.



Fig. 4. Two pairs of interleaved repeats of length  $L - 1$  create ambiguity: from the reads, it is impossible to know whether the sequences  $x$  and  $y$  are as shown, or swapped.

<sup>1</sup>A. Motahari, G. Bresler, and D. Tse, “Information theory of DNA shotgun sequencing.” IEEE Transactions on Information Theory 59.10 (2013): 6273-6289.

- d. We observe that even if we have access to all length- $L$  reads of the sequence, *repeats* make reconstruction impossible (see figure). Denoting by  $S_i^L$  the length- $L$  subsequence starting at position  $i$ , and  $R_L$  the number of length- $L$  repeats, we have

$$\mathbb{E}[R_L] = \sum_{1 \leq i < j \leq G} \mathbb{P}[S_i^L = S_j^L].$$

Justify the following:

$$\mathbb{E}[R_L] > \left( \frac{G^2}{2} - GL \right) e^{-LH_2(P)}.$$

*Hint*– For a given sequence generated by  $(p_1, p_2, p_3, p_4)$ , what is the probability that two specific physically disjoint length- $\ell$  subsequences are identical? In the sum, drop the terms in which  $S_i^L$  and  $S_j^L$  overlap.

- e. *Phase transition*: For  $G \gg L$ , the above bound may be approximated as

$$\mathbb{E}[R_L] \approx \frac{G^2}{2} e^{-LH_2(P)}.$$

Let  $G, L \rightarrow \infty$  with  $L/\ln G = \bar{L}$ , a constant. Conclude that the expected number of repeats approaches zero if

$$\bar{L} > 2/H_2(P)$$

and approaches infinity if

$$\bar{L} < 2/H_2(P).$$

Interpret this result as a prescription for *how large  $L$  should be* in order for reconstruction to be successful. Observe that  $N$  does not play any role here.

- f. Assuming that  $\bar{L} > 2/H_2(P)$  and  $N$  equals the estimate obtained in part b, conclude that the number of reads (of length  $L$ ) per nucleotide, given by  $N/G$  is roughly  $H_2(P)$ .