

1 Distribution

$$\text{Bin}(n, p) : P_X(x) = \binom{n}{x} p^x q^{n-x} \text{ for } x \in [n]_0.$$

$$\mu = np, \sigma^2 = npq, H(X) = \frac{1}{2} \log(2\pi npq) + O(\frac{1}{n}).$$

$$\text{Pois}(\lambda) : P_X(x) = \frac{\lambda^x e^{-\lambda}}{x!} \text{ for } x \in \mathbb{N}_0.$$

$$\mu = \sigma^2 = \lambda.$$

$$\text{Geo}(p) : P_X(x) = q^{x-1} p \text{ for } x \in \mathbb{N}.$$

$$\mu = \frac{1}{p}, \sigma^2 = \frac{q}{p^2}, H(X) = -q \log q - p \log p.$$

$$\text{Exp}(\lambda) : f_X(x) = \lambda e^{-\lambda x} \text{ for } x \in \mathbb{R}_+^1.$$

$$\mu = \frac{1}{\lambda}, \sigma^2 = \frac{1}{\lambda^2}$$

$$\mathcal{N}(\mu, \sigma^2) : f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}.$$

$$h(X) = \frac{1}{2} \log(2\pi e \sigma^2).$$

$$\text{Lap}(\mu, b) : f_X(x) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}.$$

$$\sigma^2 = 2b^2, h(X) = \log(2be).$$

2 Markov Chain

$$X_1 \rightarrow X_2 \rightarrow \cdots \rightarrow X_n := \mathbf{X}_n, x^n, P_{X_{n+1}|X^n}(x_{n+1}|x^n) =$$

$$P_{X_{n+1}|X}(x_{n+1}|x_n).$$

$$\text{Stationary: } P_{X_1, \dots, X_n} = P_{X_{l+1}, \dots, X_{l+n}}, \forall n, l \in \mathbb{N}.$$

3 Central Limit Theorem

Khinchin WLLN: X_1, X_2, \dots , are i.i.d. with $\mathbb{E}[|X_i|] < \infty$, then $\forall \epsilon > 0$, $\lim_{n \rightarrow \infty} \Pr\{|X_n - \mu| \geq \epsilon\} = 0$.

Central limit theorem: X_1, X_2, \dots , are i.i.d. with $\mathbb{E}[|X_i|] < \infty$, then $\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z \sim \mathcal{N}(0, 1)$.

Berry-Esseen: X_1, X_2, \dots , are i.i.d. with $\mathbb{E}[|X_i - \mu|^3] = \mu_3 < \infty$. Let $Z_n := \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$, $Z \sim \mathcal{N}(0, 1)$. Then $|F_{Z_n}(z) - F_Z(z)| \leq c \frac{\mu_3}{\sigma^3 n^{1/2}}$, $\forall z \in \mathbb{R}$, $n \in \mathbb{N}$ for constant $c \in (0, 4, 0.5)$.

$$\text{Author: 斗得卿} \quad 1$$

$$\max_{S \in \mathcal{H}(X)} - \log \sum_{x \in S} P_X(x)^{1-\epsilon} P_X(x)^\epsilon.$$

$$\text{Theorem 11: } \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \log P_X^n(\pi, n) \right\} = \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \log P_X^n(n) \right\} = CI(P_X, P_X).$$

11 Delivering Information Reliably

$$BSC(p): \text{ flip the bit bit i.i.d. with probability } p \in (0, \frac{1}{2}).$$

12 Mutual Information

$$\begin{aligned} I(X; Y) &:= D(P_{XY} \| P_X \times P_Y). \\ \text{Exercise 1: } I(X; Y) &= \min_{Q_Y: D(P_Y|Q_Y) < \infty} D(P_{XY} \| Q_Y \times P_X). \\ I(X; Y|Z) &:= H(X|X; Y) - H(X|Y; Z). \\ \text{Chain rule: } I(X; Y^n) &= \sum_{i=1}^n I(X; Y_i|Y^{i-1}). \\ X - Y - Z, \text{ then } I(X; Y) &\geq I(X; Z). \\ X - Y - Z, \text{ then } I(X; Y) &\geq I(X; Y|Z). \end{aligned}$$

13 Noisy Channel Coding Theorem

$$\begin{aligned} \text{An } (n, k) \text{ code with } P^{(n)} &:= \Pr\{W \neq \hat{W}\} \leq \epsilon \text{ is called an } (n, k, \epsilon) \text{ code.} \\ k^*(n, k) \text{ is the largest } k \text{ s.t. } &\exists (n, k, \epsilon) \text{ code.} \\ C(n) &:= \lim_{n \rightarrow \infty} \frac{1}{n} k^*(n, \epsilon). \\ \text{Channel coding theorem for DMC without feedback:} \\ C(n) &= C^* := \max_{P_X} I(X; Y), \forall \epsilon \in (0, 1). \\ x^n \text{ is robust typical sequence: } &|P_{X^n}(a) - P_X(a)| \leq \epsilon P_X(a), \text{ where } P_{X^n}(a) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = a\}. \end{aligned}$$

$$\text{Author: 斗得卿} \quad 4$$

4 Representing An i.i.d. Sequence Almost Losslessly

$$\begin{aligned} \text{DMS: discrete memoryless source. } B(n, \epsilon) \text{ is an } \epsilon\text{-high-probability set: } &\Pr\{S^n \in B(n, \epsilon)\} \geq 1 - \epsilon. \\ \delta^{\text{-typical}} \text{ is } \delta\text{-typical: } &\frac{1}{n} \sum_{i=1}^n \log P_S(s_i) + H(S) \leq \delta. \\ \delta\text{-typical set } \mathcal{A}_n^{(\delta)}(S) := \{s^n | s^n \text{ is } \delta\text{-typical}\}. \\ \text{Properties of typical sequences and typical sets:} \\ \bullet \forall s^n \in \mathcal{A}_n^{(\delta)}(S), 2^{-n(H(S)+\delta)} &\leq \Pr\{S^n = s^n\} \leq 2^{-n(H(S)-\delta)}. \\ \bullet \Pr\{S^n \in \mathcal{A}_n^{(\delta)}(S)\} &\geq 1 - \epsilon \text{ for } n \text{ large enough.} \\ \bullet |\mathcal{A}_n^{(\delta)}(S)| &\leq 2^{n(H(S)+\delta)}. \\ \bullet |\mathcal{A}_n^{(\delta)}(S)| &\geq (1-\epsilon) 2^{n(H(S)-\delta)} \text{ for } n \text{ large enough.} \end{aligned}$$

$$\begin{aligned} s^n \rightarrow \hat{s}^n \rightarrow s^{*n} \text{ (} n, k \text{) code.} \\ (n, k, \epsilon, c) \text{ code: (} n, k \text{) code with } P^{(n)} &:= \Pr\{S^n \neq \hat{S}^n\} \leq \epsilon. \\ k^*(n, c): \text{ the smallest } k \text{ s.t. } &\exists (n, k, \epsilon) \text{ code.} \\ R^*(c) &:= \lim_{n \rightarrow \infty} \frac{k^*(n, c)}{n}. \end{aligned}$$

$$\begin{aligned} \text{A lossless source coding theorem for DMS: } R^*(c) &= H(S), \forall c \in (0, 1). \\ \text{AEP (Asymptotic Equipartition Property): Entropy} &\text{ determines the asymptotic size of a typical set, and determines the probability of a typical sequence asymptotically.} \end{aligned}$$

$$\begin{aligned} \text{H and if may be different: consider } X_1, X_2, \dots \text{ are i.i.d. and } X_{2k} = X_{2k+1}. \\ \text{If } \{X_i\} \text{ is stationary, then } H(X_n|X^{n-1}) \text{ is decreasing in } n. \\ \text{If } \{X_i\} \text{ is stationary, then } H(\{X_i\}) = \hat{H}(\{X_i\}). \\ \text{Stationary ergodic processes:} \\ \frac{1}{n} \sum_{i=1}^n f(X_{n+1}, \dots, X_{n+i}) \xrightarrow{\text{a.s.}} \mathbb{E}[f(X_{n+1}, \dots, X_{n+i})] \end{aligned}$$

5 Entropy

$$\begin{aligned} H(X|Y) &= \sum_y P_Y(y) H(X|Y = y) = \sum_{x,y} P_{XY}(x, y) \log \frac{1}{P_X(x|y)}. \\ 0 \leq H(X) \leq \log |X|, \text{ where } H(X) &= \log |X| \iff \end{aligned}$$

$$\text{Author: 斗得卿} \quad 1$$

$$\begin{aligned} X \text{ is uniform distributed over } \mathcal{X}. \\ H(X; Y) &= H(Y) + H(X|Y) = H(X) + H(Y|X). \\ H(X|Y) &\leq H(X), \text{ but } H(X|Y) = y \text{ may } > H(X). \\ H(X_1, \dots, X_n) &= \sum_{i=1}^n H(X_i|X_1, \dots, X_{i-1}). \\ H(X|Y; Z) &\leq H(X|Y). \\ \text{The above still holds for } h. \\ \text{Exercise 4: } H(X; Y; Z) &\leq H(X; Y) + H(X; Z) - H(X). \end{aligned}$$

$$\begin{aligned} \text{Concavity of Entropy: } H(\mathbf{p}) &:= -\sum_{i=1}^d p_i \log p_i \text{ is concave in } \mathbf{p}. \\ \text{That is, } H(\mathbf{p}_1 + (1-\lambda)\mathbf{p}_2) &\geq \lambda H(\mathbf{p}_1) + (1-\lambda)H(\mathbf{p}_2). \\ \text{Fano's inequality: } H(U|V) &\leq H_b(P_e) + P_e \log |\mathcal{U}|, \text{ where } P_e := \Pr\{U \neq V\}. \\ \Rightarrow \Pr\{U \neq V\} &\geq \frac{H(U|V) - 1}{\log |\mathcal{U}|}. \\ \text{Exercise 5: if } U, V \text{ both take values in } \mathcal{U}, \text{ then } \\ H(U|V) &\leq H_b(P_e) + P_e \log(|\mathcal{U}| - 1). \end{aligned}$$

6 Representing A Sequence with Memory Almost Losslessly

$$\begin{aligned} \text{Entropy rate:} \\ \bullet \hat{H}(\{X_i\}) &:= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) \text{ if exists.} \\ \bullet \hat{H}(\{X_i\}) &:= \lim_{n \rightarrow \infty} H(X_n|X^{n-1}) \text{ if exists.} \end{aligned}$$

$$\begin{aligned} \text{If } \{X_i\} \text{ is stationary, then } H(X_n|X^{n-1}) \text{ is decreasing in } n. \\ \text{If } \{X_i\} \text{ is stationary, then } H(\{X_i\}) = \hat{H}(\{X_i\}). \\ \lim_{n \rightarrow \infty} \Pr\{X_{n+1} = y|X_n = x\} = \pi(y), \forall x, y \in \mathcal{X}, \text{ where } \pi(\cdot) \text{ is the unique steady-state distribution.} \\ \text{If } P_{X_n} = \pi, \text{ the Markov process becomes a stationary process.} \end{aligned}$$

$$\text{Author: 斗得卿} \quad 2$$

$$\begin{aligned} D_{\max} &:= \min_{S \in \mathcal{H}(X)} \mathbb{E}[d(S, \hat{S})]. \\ D_{\max} &= \min_{S \in \mathcal{H}(X)} \mathbb{E}[D(S, \hat{S})]. \\ R(D) &= R^*(D) := \min_{P_{X,Y}: R_0(X,Y) \leq D} I(X; Y). \\ R^*(D_{\max}) &\leq H(S), R^*(D) = 0 \text{ if } D \geq D_{\max}. \\ \text{Ber}(p) \text{ source: } R(D) &= \begin{cases} H(p) - H_0(D), & \text{if } 0 \leq D \leq \min\{p, 1-p\} \\ 0, & \text{if } D > \min\{p, 1-p\} \end{cases} \\ \text{Gaussian source: } R(D) &= \begin{cases} \frac{1}{2} \log \frac{\sigma^2}{D}, & \text{if } 0 \leq D \leq \sigma^2 \\ 0, & \text{if } D > \sigma^2. \end{cases} \\ R(D) &\leq R^*(D). \end{aligned}$$

$$\begin{aligned} D_{\max} &:= \min_{S \in \mathcal{H}(X)} \mathbb{E}[d(S, \hat{S})]. \\ D_{\max} &= \min_{S \in \mathcal{H}(X)} \mathbb{E}[D(S, \hat{S})]. \\ R(D) &= R^*(D) := \min_{P_{X,Y}: R_0(X,Y) \leq D} I(X; Y). \\ R^*(D_{\max}) &\leq H(S), R^*(D) = 0 \text{ if } D \geq D_{\max}. \\ \text{Ber}(p) \text{ source: } R(D) &= \begin{cases} H(p) - H_0(D), & \text{if } 0 \leq D \leq \min\{p, 1-p\} \\ 0, & \text{if } D > \min\{p, 1-p\} \end{cases} \\ \text{Gaussian source: } R(D) &= \begin{cases} \frac{1}{2} \log \frac{\sigma^2}{D}, & \text{if } 0 \leq D \leq \sigma^2 \\ 0, & \text{if } D > \sigma^2. \end{cases} \\ R(D) &\leq R^*(D). \end{aligned}$$

$$\text{Author: 斗得卿} \quad 5$$

$$\begin{aligned} \text{For a homogeneous, irreducible, and aperiodic Markov process } \{X_i\}, H(\{X_i\}) = \hat{H}(\{X_i\}) = H(X_1|X_0), P_{X_1} = \sum_{x \in \mathcal{X}} \pi(x) H(X_1|X_0 = x), \text{ where } \pi \text{ is the unique steady-state distribution.} \end{aligned}$$

$$\begin{aligned} \text{Chernoff-Sid lemma: } \lim_{n \rightarrow \infty} \frac{1}{n} \log \omega_n(n, \epsilon) = D(P_0 \| P_1). \\ \text{Typical set:} \end{aligned}$$

7 Information for Continuous Distributions

$$\begin{aligned} \text{The covariance of } n\text{-dimensional } X \text{ is } k_n, \text{ then } h(X) = h(X^{(c)}) = \frac{1}{2} \log((2\pi e)^n \det(k)). \end{aligned}$$

8 Learning a Bit of Information

$$\begin{aligned} \pi_{\theta|0}(\phi): \text{ false alarm, false positive, false rejection, type I error.} \\ \pi_{\theta|1}(\phi): \text{ miss detection, false negative, false acceptance, type II error.} \\ A_\phi(\phi): \text{ acceptance region of } H_0. \\ \text{Likelihood ratio } LR(x) &:= \frac{P_1(x)}{P_0(x)}, \text{ Log likelihood ratio } LLR(x) := \log LR(x). \\ \text{Likelihood ratio test (LRT) with parameter } \tau \in \mathbb{R}_+^d &\text{ is } \phi_\tau^{LRT}(x) = \mathbb{1}\{LR(x) > \tau\}. \end{aligned}$$

$$\begin{aligned} \text{(Randomized) LRT } \phi_\tau, x \mapsto \begin{cases} 1, & \text{if } LR(x) > \tau \\ \gamma, & \text{if } LR(x) = \tau \\ 0, & \text{if } LR(x) < \tau \end{cases} \\ \text{Neyman-Pearson problem: minimize } \pi_{0|0}(\phi) \text{ (subject to } \pi_{1|0}(\phi) \leq \epsilon. \\ \text{Neyman-Pearson: LRT is optimal.} \\ \text{Generalized to } n \text{ i.i.d.: } \phi_{\tau, n}^{LRT}(x^n) = \text{Chernoff Information: } CI(P_0, P_1) := \end{aligned}$$

$$\begin{aligned} P_0, P_1 \text{ are given.} \\ P_1(a) = \sum_{b \in \mathcal{A}} P_1(b|a) P_1(b). \\ \text{Exercise 6: } D(P_1 \| P_0) \text{ is a continuous and strictly increasing function of } \lambda \text{ for } \lambda \in [0, 1]. \\ P_1^*(\pi) (= (n_0, \pi_1)), (n) := \min\{n_0 \pi_0^*(n) + \pi_1 n_1^*(n)\}. \\ P_1^*(n) = \min\{\max\{\pi_0^*, \pi_1^*\} \}. \\ \text{Chernoff Information: } CI(P_0, P_1) := \end{aligned}$$

$$\text{Author: 斗得卿} \quad 3$$

Homework 1 Solution

TA: Heng-Chien Lion *

r1942067@ntu.edu.tw

1. (Another kind of typical sequences) [18]

$$\begin{aligned} \text{In this problem, let us consider another kind of typical sequences defined as follows.} \\ \text{Definition. For } \gamma \in (0, 1), \text{ a sequence } s^n \text{ is called } \gamma\text{-typical with respect to a DMS } S \sim P_S \text{ if} \\ |\pi(a|s^n) - P_S(a)| \leq \gamma P_S(a), \forall a \in \mathcal{S}, \\ \text{where } \pi(a|s^n) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{s_i = a\}. \text{ The } \gamma\text{-typical set} \\ \mathcal{T}_\gamma^{(n)}(S) := \{s^n \in \mathcal{S}^n | s^n \text{ is } \gamma\text{-typical with respect to } S\}. \end{aligned}$$

In the following, please assume $|S| < \infty$.

$$\begin{aligned} \text{a) Show that the typical sequence and typical set defined above also satisfy AEP (Proposition 1 of Unit 1) with } \mathcal{A}_n^{(c)}(S) \text{ replace by } \mathcal{T}_\gamma^{(n)}(S), \text{ and the } \delta \text{ in properties 1, 3, and 4 replaced by something, denoted by } \zeta(\gamma), \text{ depending on } \gamma. \text{ Specify this } \zeta(\gamma). \\ \text{b) Show that } \mathcal{T}_\gamma^{(n)}(S) \subseteq \mathcal{A}_n^{(c)}(S) \text{ where } \delta = \zeta(\gamma) \text{ found in a).} \\ \text{c) Find an alphabet } \mathcal{S}, \text{ a reference probability mass function } P_S, \text{ and } \gamma \text{ such that } \forall \delta' > 0, n \in \mathbb{N}, \mathcal{A}_n^{(\delta')}(S) \not\subseteq \mathcal{T}_\gamma^{(n)}(S). \end{aligned}$$

Remark. From b) and c) we see that the typicality defined in this problem is *stronger* than that defined in the lecture. Hence, they are called strong typicality and weak typicality respectively in the literature.

$$\begin{aligned} \text{Solution:} \\ \bullet \text{Property 1 } \forall s^n \in \mathcal{T}_\gamma^{(n)}, \text{ we have } \forall a \in \mathcal{S}, |\pi(a|s^n) - P_S(a)| \leq \gamma P_S(a) \\ \frac{1}{n} \log \Pr\{S^n = s^n\} = \sum_{a \in \mathcal{S}} \pi(a|s^n) \log P_S(a) \end{aligned}$$

*With contributions by Wen-Shao Ho and Chen-Hao Hsiao.

$$\begin{aligned} \mathcal{T}_\gamma^{(n)} \subseteq \mathcal{A}_\gamma^n \text{ for } \delta = \zeta(\gamma) = \gamma H(S) \\ \text{c) For some } \gamma < 1, \text{ consider } \mathcal{S} = \{0, 1\} \text{ and } P_S(0) = P_S(1) = \frac{1}{2}. \text{ Consider } s^n = 0^n. \\ H(S) = 1 \Rightarrow 0^n \in \mathcal{A}_\gamma^n, \forall \gamma > 0. \text{ But } 0^n \notin \mathcal{T}_\gamma^{(n)}. \\ \text{Grading Policy:} \\ \text{a) Property 1 words [2] points, Property 2 words [3] points, and Property 3, 4 word [1] point in total.} \\ \text{b) Correct arguments about set relations [4] and correct arithmetic details [2].} \\ \text{c) A correct example [3], and justification [3].} \end{aligned}$$

2. (Finer asymptote for lossless source coding achievability) [20]

$$\begin{aligned} \text{Consider a discrete memoryless source } S_i \stackrel{\text{i.i.d.}}{\sim} P_S, i = 1, 2, \dots, \text{ where } P_S \text{ is the PMF of the source } S. \text{ Let } R > H(S), \text{ that is,} \\ R = H(S) + \delta, \end{aligned}$$

$$\begin{aligned} \text{where } \delta > 0 \text{ denotes a constant. Then, it can be shown as a corollary of the lossless source coding theorem in our lecture that there exists a sequence of } \{n_i, \lfloor n_i R \rfloor\} \text{ codes such that } \forall \epsilon > 0, \\ P_{n_i}^{(n_i)} \leq \epsilon \text{ for } n \text{ sufficiently large.} \end{aligned}$$

$$\begin{aligned} \text{Notably, the gap to the fundamental limit, } \delta, \text{ is a constant not depending on } n. \\ \text{Suppose we would like to achieve (1) for a given } \epsilon \in (0, 1/2). \text{ Is it possible to derive a finer asymptote for } R = H(S), \text{ the gap to the fundamental limit? In this problem, we are going to show that } R = H(S) = \Theta(n^{-1/2}) \text{ suffices.} \end{aligned}$$

$$\begin{aligned} \text{a) (Warm-up) Let } \zeta(S) > 0 \text{ denote the standard deviation of } \log \frac{1}{P_S(s)} \text{ when } S \sim P_S \text{ and } \Phi(\cdot) \text{ denote the CDF of a standard normal RV. Use the central limit theorem to prove the following} \\ \lim_{n \rightarrow \infty} \Pr \left\{ \prod_{i=1}^n P_S(S_i) \geq 2^{-n(H(S) + \zeta(S) \delta)} \right\} = \Phi(\delta) \end{aligned}$$

With the above, if we define a set of length- n sequences $\mathcal{B}_\delta^{(n)}(S)$ as follows:

$$\mathcal{B}_\delta^{(n)}(S) = \left\{ s^n \mid \prod_{i=1}^n P_S(s_i) \geq 2^{-n(H(S) + \zeta(S) \delta)} \right\},$$

$$\begin{aligned} \text{one can control the probability of each sequence in } \mathcal{B}_\delta^{(n)}(S) \text{ from below and hence can control the cardinality of this set from above. Also, we know that } \Pr\{S^n \in \mathcal{B}_\delta^{(n)}(S)\} \rightarrow \Phi(\delta) \text{ as } n \rightarrow \infty \text{ from Part a). It is hence tempting to us label all the sequences in } \mathcal{B}_\delta^{(n)}(S) \text{ and give up the rest as a source encoding scheme. However, to upper bound the error probability, knowing its limit as } n \rightarrow \infty \text{ is not enough. Berry-Esseen theorem is a standard refinement of the CLT.} \end{aligned}$$

$$\begin{aligned} \text{b) Show that (1) can be attained using the aforementioned scheme if the rate approaches } H(S) \text{ from above as } n \rightarrow \infty \text{ in the following manner:} \\ R_n = H(S) + n^{-1/2} \zeta(S) \Phi^{-1}(\epsilon) + \zeta_n \end{aligned}$$

$$\text{where } \zeta_n = O(n^{-1}) \text{ denotes a positive sequence tends to zero not slower than } n^{-1}. \quad [12]$$

$$\begin{aligned} \text{Remark. The above is not optimal – the optimal asymptote of the rate (when } \zeta(S) > 0) \text{ is} \\ R_n = H(S) + n^{-1/2} \zeta(S) \Phi^{-1}(\epsilon) - \frac{\log \zeta(S)}{2n} + O(1/n). \end{aligned}$$

$$\begin{aligned} \text{Solution:} \\ \text{a) Central limit theorem states that for i.i.d random variables } X_1, \dots, X_n \text{ with finite mean } \mu, \text{ variance } \sigma^2, \text{ As } n \rightarrow \infty, \text{ the random variable } \sqrt{n}(X_n - \mu) \text{ converges in distribution to } \mathcal{N}(0, \sigma^2), \text{ where } X_n := \frac{1}{n} \sum_{i=1}^n X_i. \text{ Namely, for any real number } z, \\ \lim_{n \rightarrow \infty} \Pr \left\{ \sqrt{n}(X_n - \mu) \leq z \right\} = \Phi \left(\frac{z}{\sigma} \right). \\ \text{Choose } X_i = \log \frac{1}{P_S(S_i)} \text{ and } Z = \delta \zeta(S) \text{ in the Central limit theorem:} \\ \lim_{n \rightarrow \infty} \Pr \left\{ \sqrt{n} \left(-\frac{1}{n} \log \prod_{i=1}^n P_S(S_i) - \mathbb{E} \left[\log \frac{1}{P_S(S)} \right] \right) \leq \delta \zeta(S) \right\} = \Phi(\delta) \\ \Leftrightarrow \lim_{n \rightarrow \infty} \Pr \left\{ \left(-\frac{1}{n} \log \prod_{i=1}^n P_S(S_i) \right) \leq H(S) + n^{-1/2} \delta \zeta(S) \right\} = \Phi(\delta) \\ \Leftrightarrow \lim_{n \rightarrow \infty} \Pr \left\{ \prod_{i=1}^n P_S(S_i) \geq 2^{-n(H(S) + n^{-1/2} \delta \zeta(S))} \right\} = \Phi(\delta). \end{aligned}$$

$$\begin{aligned} \text{b) Note that we can now label all the sequences in } \mathcal{B}_\delta^{(n)}(S) \text{ and give up the rest, we want to find } \delta(\epsilon) \text{ such that when } n \text{ is sufficiently large,} \\ \Pr \left\{ S^n \in \mathcal{B}_{\delta(n)}^{(n)}(S) \right\} \geq 1 - \epsilon. \end{aligned}$$

Combine Berry-Esseen theorem and subproblem a), we know that

$$\Pr \left\{ S^n \in \mathcal{B}_{\delta(n)}^{(n)}(S) \right\} \geq \Phi(\delta(\epsilon)) - \epsilon \frac{\rho^3}{\zeta(S)^3} n^{-1/2}.$$

Therefore, it suffices to find $\delta(\epsilon)$ such that

$$\Phi(\delta(\epsilon)) - \epsilon \frac{\rho^3}{\zeta(S)^3} n^{-1/2} \geq 1 - \epsilon \Leftrightarrow \delta(\epsilon) \geq \Phi^{-1} \left(1 - \epsilon + \epsilon \frac{\rho^3}{\zeta(S)^3} n^{-1/2} \right).$$

The order of $\Phi^{-1}(1 - \epsilon + \epsilon \frac{\rho^3}{\zeta(S)^3} n^{-1/2})$ can be further characterized using Taylor expansion at $1 - \epsilon$ on the function $\Phi^{-1}(x)$:

$$\Phi^{-1}(x) = \Phi^{-1}(1 - \epsilon) + \frac{\Phi^{-1}(1 - \epsilon)}{1} (x - (1 - \epsilon)) + O(x^2).$$

Plug in $x = 1 - \epsilon + e^{\frac{\epsilon^2}{(nS)^2}} n^{-1/2}$ we have

$$\Phi^{-1}\left(1 - \epsilon + \frac{\epsilon^2}{(nS)^2}\right) = \Phi^{-1}(1 - \epsilon) + O(n^{-1/2}) = -\Phi^{-1}(\epsilon) + O(n^{-1/2}).$$

The result above indicates that if we choose $\delta(\epsilon) = -\Phi^{-1}(\epsilon) + O(n^{-1/2})$, we can guarantee

$$\Pr\left\{S^n \in \mathcal{B}_{\delta(\epsilon)}^{(n)}(S)\right\} \geq 1 - \epsilon.$$

Similar to the proof of the properties of the typical set, for n sufficiently large:

$$\begin{aligned} 1 &\geq \Pr\left\{S^n \in \mathcal{B}_{\delta(\epsilon)}^{(n)}(S)\right\} \geq |\mathcal{B}_{\delta(\epsilon)}^{(n)}(S)| 2^{-n(H(S) + n^{-1/2}\delta_0(S))} \\ \Rightarrow |\mathcal{B}_{\delta(\epsilon)}^{(n)}(S)| &\leq 2^{n(H(S) + n^{-1/2}\delta_0(S))} = 2^{n(H(S) - n^{-1/2}\delta_0(S) + o(n^{-1}))}. \end{aligned}$$

Therefore, the choice

$$R_n = H(S) - n^{-1/2}\delta_0(S)\Phi^{-1}(\epsilon) + O(n^{-1}),$$

makes it possible to encode all the sequences in $\mathcal{B}_{\delta(\epsilon)}^{(n)}(S)$ and also satisfies the requirement that $\Pr\left\{S^n \in \mathcal{B}_{\delta(\epsilon)}^{(n)}(S)\right\} \geq 1 - \epsilon$.

Grading Policy:

- Justification for your application of central limit theorem [4], other details [4].
- Overall logic in scheme and analysis [4], application of Berry-Esseen theorem [4], other details such as mean value theorem or Taylor's theorem [4].

3. (An alternative lossless source coding theorem) [12]

For a discrete memoryless source $\{S_i | i \in \mathbb{N}\}$, consider a sequence of $(n, \lfloor nR \rfloor)$ source codes indexed by $n = 1, 2, \dots$ with compression rate $R > 0$.

Prove the following statements.

- If $R > H(S)$, there exist a sequence of $(n, \lfloor nR \rfloor)$ codes with

$$\lim_{n \rightarrow \infty} P_{\epsilon}^{(n)} = 0.$$

In other words, the probability of error can be driven to zero as $n \rightarrow \infty$.

- If $R < H(S)$, for any sequence of $(n, \lfloor nR \rfloor)$ codes, the sequence of error probabilities must converge to 1, that is,

$$\lim_{n \rightarrow \infty} P_{\epsilon}^{(n)} = 1.$$

$$H(X) = \sum_{i=1}^{\infty} p_i \log q_i = \sum_{i=1}^{\infty} p_i \left[\log \frac{p_i}{p_i} + \log(\mu_X - 1) \right] = \mu_X \mathbf{h}_X\left(\frac{1}{\mu_X}\right)$$

Equality holds iff $p_i = q_i \forall i$.

Grading Policy:

- Correct arithmetic [4] and application of Jensen's inequality on the right variable [4].
- Specify α and β [8].
- Arithmetic [3] and selection of $\{p_i\}_{i=1}^{\infty}$ [3].

2. (Entropy of a random variable with a infinitely countable support) [14]

In the lecture, we define the entropy for a random variable with a finite alphabet \mathcal{X} (in fact a finite support $\text{supp}_x = \text{infset}_x$). For a random variable X that has an infinitely countable support, sometimes $H(X)$ is finite and sometimes $H(X)$ becomes infinite. In this problem we look at an example.

- Consider an infinite series

$$\sum_{n=1}^{\infty} \frac{1}{n(\log n)^{\alpha}}$$

where $\alpha \geq 0$. Use the *integral test for convergence* to show that the series converges if and only if $\alpha > 1$.

- Let s_n denote the above series if the series converges. Let us define a random variable $X_n \in \{2, 3, \dots\}$ with PMF

$$P_{X_n}(n) = \frac{1}{s_n n (\log n)^{\alpha}}.$$

Show that $H(X_n)$ exists if $\alpha > 2$ and it diverges to ∞ if $1 < \alpha \leq 2$.

Solution:

$$a) \int_2^{\infty} \frac{dx}{x^2 (\log x)^{\alpha}} = \begin{cases} \frac{1}{(n-1)} & (\alpha > 1) \\ \infty & (0 \leq \alpha \leq 1) \end{cases}$$

By the integral test for convergence, the series converges if and only if $\alpha > 1$. Many students derive wrong integral when $\alpha = 1$.

- Let $s_n := \sum_{i=2}^{\infty} \frac{1}{i(\log i)^{\alpha}}$ for $\alpha > 1$.

$$\begin{aligned} H(X_n) &= \sum_{n=2}^{\infty} P_{X_n}(n) \log \frac{1}{P_{X_n}(n)} \\ &= \log s_n + \frac{1}{s_n} \sum_{n=2}^{\infty} \frac{1}{n(\log n)^{\alpha-1}} + \frac{1}{s_n} \sum_{n=2}^{\infty} \frac{\alpha \log(\log n)}{n(\log n)^{\alpha}} \end{aligned}$$

Homework 3 Simple Solution

TA: Heng-Chien Liou*

1. (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:

$$\begin{aligned} \mathcal{H}_0: & p = p_0 \\ \mathcal{H}_1: & p = p_1 \end{aligned}$$

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

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

Draw the optimal (π_{10}, π_{01}) trade-off curve.

- Suppose the decision maker waits until an 1 appears and makes the decision based on the whole observed sequence. Sketch the optimal (π_{10}, π_{01}) trade-off curve.

- Now suppose the decision maker waits until in total n 1's appear and makes the decision based on the whole observed sequence. Let $\pi_{\epsilon}^{(n)}(n, \epsilon)$ 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 $\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi_{\epsilon}^{(n)}(n, \epsilon)$ exist? If so, find it. Otherwise, show that the limit does not exist.

Solution:

- By Neyman-Pearson theorem, the optimal test is randomized LRT. Note that the likelihood ratio can only take two values: $\frac{p_1}{p_0}$ and $\frac{1-p_1}{1-p_0}$. Therefore, discuss the range of τ we get

$$\begin{cases} \pi_{10} = 0, \pi_{01} = 0, & 0 \leq \tau < \frac{1-p_0}{1-p_1} \\ \pi_{10} = p_0 + \gamma(1-p_0), \pi_{01} = (1-\gamma)(1-p_1) = \frac{1-p_0}{1-p_1}(1-\pi_{10}), & \tau = \frac{1-p_0}{1-p_1} \\ \pi_{10} = p_0, \pi_{01} = 1-p_1, & \frac{1-p_0}{1-p_1} < \tau < \frac{p_0}{p_1} \\ \pi_{10} = \gamma p_0, \pi_{01} = (1-\gamma)p_1 + (1-p_1) = 1 - \frac{p_0}{p_1} \pi_{10}, & \tau = \frac{p_0}{p_1} \\ \pi_{10} = 0, \pi_{01} = 1, & \tau > \frac{p_0}{p_1} \end{cases}$$

*with contribution by Chen-Hao Hsiao

Solution:

- Let $\delta = R - H(S)$. A simple scheme is to use $\mathcal{A}_{2\delta}^{(n)}$ and standard asymptotic equipartition property for coding. Here we try to prove it using the lossless source coding theorem in the lecture.

Consider a sequence $\{\epsilon_n\}_{n=1}^{\infty} > 0$ such that $\lim_{n \rightarrow \infty} \epsilon_n = 0$. For any $n, \epsilon_n > 0$, consider $2k^*$ such that for any $k \geq k^*$ (n, ϵ_n), there exists a (n, k) code with $P_{\epsilon}^{(n)} < \epsilon_n$. By the lossless source coding theorem in the lecture and the Archimedean property, $\exists n_0 > \max_{i=1}^k \{n_i\}$ such that $\forall n \geq n_0, k^* (n, \epsilon_n) < n(H(S) + \delta/2) < \lfloor nR \rfloor$. $\{n_m\}_{m=1}^{\infty}$ then is a strictly increasing sequence of integers.

Therefore, we can construct a sequence of $(n, \lfloor nR \rfloor)$ codes as follows: for n such that $n_m \leq n < n_{m+1}$ select $k = k^* (n, \epsilon_m)$ and one of the corresponding (n, k) code $(\text{enc}_n, \text{dec}_n)$ such that $P_{\epsilon}^{(n)} < \epsilon_m$. By the above justification, $k < \lfloor nR \rfloor$. $\{(\text{enc}_n, \text{dec}_n)\}_n$ is then a desirable code with the property.

- $\mathcal{B}(S) = R$. Again, a simple strategy is to use $\mathcal{A}_{2\delta}^{(n)}$ and standard asymptotic equipartition property to characterize the error probability for any sequence of codes. Here we try to prove it using the lossless source coding theorem in the lecture.

We attempt to prove by contradiction. Suppose there exist a sequence of $(n, \lfloor nR \rfloor)$ codes with error probabilities not converging to 1, namely, for some $\epsilon < 1$, $\forall m \in \mathbb{N}$ $\exists m' (m) \geq m$ such that $P_{\epsilon}^{(m'(m))} < 1 - \epsilon$.

$\forall m \in \mathbb{N}$, the $(n'(m), \lfloor n'(m)R \rfloor)$ code from the above sequences, by definition, is then a $(n'(m), \lfloor n'(m)R \rfloor, 1 - \epsilon)$ code.

$$P_{\epsilon}^{(n'(m))} < 1 - \epsilon, \forall m \in \mathbb{N} \quad (1)$$

On the other hand, for such ϵ , by the lossless source coding theorem, we know that $\exists n_{\delta, \epsilon} \in \mathbb{N}$ such that $\forall n \geq n_{\delta, \epsilon}$,

$$k^* (n, 1 - \epsilon) > n(H(S) - \delta/2) > nR \quad (2)$$

From (1), for $m \geq n_{\delta, \epsilon}$, $\exists m' (m) \geq m \geq n_{\delta, \epsilon}$ and a corresponding $(n'(m), \lfloor n'(m)R \rfloor)$ code with error smaller than $1 - \epsilon$. But from (2), for such $n'(m) \geq n_{\delta, \epsilon}$, the smallest length required is $k^* (n'(m), 1 - \epsilon) > n'(m)R$. This leads to contradiction.

Grading Policy:

- A correct scheme [3] and the justification [3]; alternatively, a correct proof of existence [6].
- A proof that applies for all codes [4] and other details [2].

$$\begin{cases} \leq \infty & (\alpha > 2) \\ = \infty & (1 < \alpha \leq 2) \end{cases}$$

Grading Policy:

- Correct statement of the integral test [2] and correct evaluation of the integral [5].
- Correct evaluation of the relevant integrals [3] and correct argument of the series's convergence and divergence [4].

3. (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 $\{X_0[i] | i \in \mathbb{N}\}$ and $\{X_1[i] | i \in \mathbb{N}\}$ taking values in disjoint alphabets \mathcal{X}_0 and \mathcal{X}_1 respectively. The two processes are independent from each other, that is, $\{X_0[i]\} \perp \{X_1[i]\}$, and they have entropy rates \mathcal{H}_0 and \mathcal{H}_1 respectively. Let $\{\theta_i | i \in \mathbb{N}\}$ 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?
- Let $\{\Theta_i\}$ be Markov with a probability transition matrix

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

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?

Solution:

- Since $\{X_0[i]\} \in \mathbb{N}$ and $\{X_1[i]\} \in \mathbb{N}$ are stationary processes, $\forall n, l \in \mathbb{N}$,

$$\begin{aligned} P_{X_{\theta_1}[1], \dots, X_{\theta_n}[n]} &= (1 - q)P_{X_0[1], \dots, X_0[n]} + qP_{X_1[1], \dots, X_1[n]} \\ &= (1 - q)P_{X_0[1], \dots, X_0[n]} + qP_{X_1[1], \dots, X_1[n]} \\ &= P_{X_{\theta_1}[1], \dots, X_{\theta_n}[n]} \end{aligned}$$

So $\{X_{\theta_i}[i]\}$ is stationary.

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

$$\begin{aligned} \mathcal{H}(X_{\theta_i}[i]) &= \lim_{n \rightarrow \infty} H\{Y^{n-1}\} \\ &= \lim_{n \rightarrow \infty} H(Y_n, \Theta_n | Y^{n-1}, \Theta^{n-1}) \quad (\mathcal{X}_0 \cap \mathcal{X}_1 = \emptyset) \end{aligned}$$

Homework 2 Simple Solution

TA: Heng-Chien Liou*

1. (Maximum entropy) [22]

In the lecture we show that for a random variable taking values in a *finite* alphabet \mathcal{X} , the maximum entropy is $\log |\mathcal{X}|$, and the maximizing distribution is $\text{Unif}(\mathcal{X})$. In this problem, we consider a random variable X taking values in the countable alphabet \mathbb{N} (the set of positive integers) with mean $\mu < \infty$. Please use the steps below to prove that its entropy is upper bounded as follows:

$$H(X) \leq \mu \mathbf{h}_X(\mu_X^{-1}).$$

Notation: The PMF of X is denoted by $P_X(i) \equiv p_i$, $i \in \mathbb{N}$, and $\mathbb{E}[X] \equiv \mu_X = \sum_{i=1}^{\infty} ip_i$.

- Use the concavity of $\log(\cdot)$ and Jensen's inequality to show that for any non-negative sequence $\{q_i\}_{i=1}^{\infty}$ with $\sum_{i=1}^{\infty} q_i = 1$,

$$H(X) \leq \sum_{i=1}^{\infty} p_i \log q_i. \quad [8]$$

- Find a non-negative sequence $\{q_i\}_{i=1}^{\infty}$ and α, β such that $\sum_{i=1}^{\infty} q_i = 1$, $\sum_{i=1}^{\infty} iq_i = \mu_X$, and $-\log q_i = \alpha i + \beta \quad \forall i \in \mathbb{N}$.

- Complete the proof by plugging in the $\{q_i\}_{i=1}^{\infty}$ found in Part b) into the upper bound of $H(X)$ found in Part a). Show that this upper bound on $H(X)$ is attainable with an appropriate choice of $\{p_i\}_{i=1}^{\infty}$ and hence it is the maximum entropy of random variables taking values in positive integers with expected value being μ_X .

Solution:

- $H(X) + \sum_{i=1}^{\infty} p_i \log q_i = \sum_{i=1}^{\infty} p_i \log \frac{q_i}{p_i}$. By Jensen, $\sum_{i=1}^{\infty} p_i \log \frac{q_i}{p_i} \leq \log \sum_{i=1}^{\infty} p_i \frac{q_i}{p_i} = 0$.

$$b) \begin{cases} 1 = \sum_{i=1}^{\infty} q_i = \frac{e^{-(\alpha+\beta)}}{1-e^{-(\alpha+\beta)}} \\ \mu_X = \sum_{i=1}^{\infty} iq_i = \frac{e^{-(\alpha+\beta)}}{(1-e^{-(\alpha+\beta)})^2} \end{cases} \implies \begin{cases} \alpha = \log \frac{\mu_X}{\mu_X - 1} \\ \beta = \log(\mu_X - 1) \end{cases}$$

$$c) \mathbf{h}_X\left(\frac{1}{\mu_X}\right) = \frac{1}{\mu_X} \log \mu_X + \frac{\mu_X - 1}{\mu_X} \log \frac{\mu_X}{\mu_X - 1}$$

*with contributions by Wen-Shao Hsiao

$$\begin{aligned} &= \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1}, \Theta^n) + H(\Theta_n | Y^{n-1}, \Theta^{n-1}) \quad (\text{chain rule}) \\ &= \lim_{n \rightarrow \infty} H(Y_n | \Theta_n) + H(\Theta_n | \Theta^{n-1}) \\ &= \lim_{n \rightarrow \infty} H(Y_n | \Theta) + H(\Theta | \Theta) \\ &= q\mathcal{H}_1 + (1 - q)\mathcal{H}_0 + 0 \end{aligned}$$

- $\{X_{\theta_i}[i]\}$ is stationary. Let $Y_i = X_{\theta_i}[i]$. $\mathbb{E}[\{H(X_{\theta_i}[i])\}] = \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1})$.

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

Grading Policy:

- Stationary argument [2], chain rule and entropy calculation [3], asymptotic [1].
- Stationary argument [2], chain rule and entropy calculation [6].

Let the *posterior probabilities* be the conditional distribution of θ given $X^n = x^n$:

$$\pi_0^{(n)}(x^n) = \Pr\{\theta = 0 | X^n = x^n\}, \quad \pi_1^{(n)}(x^n) = \Pr\{\theta = 1 | X^n = x^n\}.$$

- Derive the expressions of $\pi_0^{(n)}(x^n)$ and $\pi_1^{(n)}(x^n)$ in terms of $\pi_0^{(0)}, \pi_1^{(0)}, P_0, P_1$.
- Consider $\pi_0^{(n)}(X^n)$ and $\pi_1^{(n)}(X^n)$ as random variables, because they are functions of the random sequence X^n . Use the Strong Law of Large Numbers to show that if \mathcal{H}_0 is true, then with probability 1,

$$\pi_0^{(n)}(X^n) \rightarrow 1, \quad -\frac{1}{n} \log \pi_1^{(n)}(X^n) \rightarrow D(P_0 | P_1) \quad \text{as } n \rightarrow \infty. \quad [8]$$

Solution:

- Denote $x^n = (x_1, x_2, \dots, x_n, \dots, x_N)$

$$\begin{aligned} \pi_0^{(n)}(X^n) &= \Pr\{\theta = 0 | X^n = x^n\} = \frac{\Pr\{\theta = 0, X^n = x^n\}}{\Pr\{X^n = x^n\}} \\ &= \frac{\pi_0^{(0)} \Pr\{X^n = x^n | \theta = 0\}}{\pi_0^{(0)} \Pr\{X^n = x^n | \theta = 0\} + \pi_1^{(0)} \Pr\{X^n = x^n | \theta = 1\}} \\ &= \frac{\pi_0^{(0)} \prod_{i=1}^n P_0(x_i)}{\pi_0^{(0)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(0)} \prod_{i=1}^n P_1(x_i)} \\ \pi_1^{(n)}(X^n) &= \frac{\pi_1^{(0)} \prod_{i=1}^n P_1(x_i)}{\pi_0^{(0)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(0)} \prod_{i=1}^n P_1(x_i)} \end{aligned}$$

- Recall the basic Bayes' theorem.

$$\begin{aligned} \pi_0^{(n)}(X^n) &= \frac{\pi_0^{(0)} \prod_{i=1}^n P_0(x_i)}{\pi_0^{(0)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(0)} \prod_{i=1}^n P_1(x_i)} \\ &= \frac{\pi_0^{(0)} 2^{-n(\frac{1}{2} \sum_{i=1}^n \log P_0(x_i))}}{\pi_0^{(0)} 2^{-n(\frac{1}{2} \sum_{i=$$

$$\Pr \left\{ \frac{1}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}} < \pi_1^{(n)}(X^*) < \frac{1}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}} \right\} = 1$$

Hence, with probability 1, $\pi_1^{(n)}(X^*) \rightarrow 1$ as $n \rightarrow \infty$.

Following this derivation, for every $n \geq N_\epsilon$, with probability 1,

$$-\frac{1}{n} \log \left(\frac{e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}} \right) < -\frac{1}{n} \log \pi_1^{(n)}(X^*) < -\frac{1}{n} \log \left(\frac{e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}} \right)$$

Furthermore,

$$\begin{aligned} -\frac{1}{n} \log \left(\frac{e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}} \right) &= \mathbb{D}(P_0|P_1) + \epsilon - \frac{\log e}{n} + \frac{1}{n} \log(1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) + \epsilon)}) \\ &\leq \mathbb{D}(P_0|P_1) + \epsilon - \frac{\log e}{n} + \frac{1}{n} \log(1 + e^{2^{-N_\epsilon}(\mathbb{D}(P_0|P_1) + \epsilon)}) \\ &= \mathbb{D}(P_0|P_1) + \epsilon + O\left(\frac{1}{n}\right) \\ -\frac{1}{n} \log \left(\frac{e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}}{1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}} \right) &= \mathbb{D}(P_0|P_1) - \epsilon - \frac{\log e}{n} + \frac{1}{n} \log(1 + e^{2^{-n}(\mathbb{D}(P_0|P_1) - \epsilon)}) \\ &\geq \mathbb{D}(P_0|P_1) - \epsilon - \frac{\log e}{n} + \frac{1}{n} \log(1 + e^{2^{-N_\epsilon}(\mathbb{D}(P_0|P_1) - \epsilon)}) \\ &= \mathbb{D}(P_0|P_1) - \epsilon + O\left(\frac{1}{n}\right) \end{aligned}$$

Hence, with probability 1,

$$\mathbb{D}(P_0|P_1) - \epsilon \leq \lim_{n \rightarrow \infty} -\frac{1}{n} \log \pi_1^{(n)}(X^*) \leq \mathbb{D}(P_0|P_1) + \epsilon,$$

where $\epsilon \in (0, \mathbb{D}(P_0|P_1))$ arbitrarily. Hence,

$$\Pr \left\{ \lim_{n \rightarrow \infty} -\frac{1}{n} \log \pi_1^{(n)}(X^*) = \mathbb{D}(P_0|P_1) \right\} = 1$$

Grading Policy:

- $\pi_1^{(n)}(x^*)$ [2] and $\pi_1^{(n)}(x^*)$ [2].
- Application of SLN and argument for almost sure convergence [4], convergence for $\pi_1^{(n)}(X^*)$ [2], convergence for $-(1/n) \log \pi_1^{(n)}(X^*)$ [2].

3. Minimizing information divergence) [22]

- Let $\mathcal{P}(\mathcal{N})$ denote the collection of all probability distributions over \mathcal{N} and $G(p) \in \mathcal{P}(\mathcal{N})$ be a geometric distribution with parameter $p \in (0, 1)$:

$$X \sim G(p) \iff \Pr\{X = n\} = (1-p)p^{n-1}, \quad n \in \mathcal{N} = \{1, 2, \dots\}.$$

4

Homework 4 Simple Solution

TA: Heng-Chien Liou*

1. (Mutual information) [8]

- How much information does the length of a sequence give about the content of a sequence? In this problem, let us consider a sequence of i.i.d. $\text{Ber}(1/2)$ random variables X_1, X_2, \dots . Let N_0 be a random variable taking value at 6 with probability $1/3$ and at 12 with probability $2/3$. Furthermore, N_0 is independent of the sequence $\{X_i \mid i = 1, 2, \dots\}$. Also let N_1 denote the length of the sequence when the first "1" appears. Obviously N_1 is also a random variable.

$$\text{Compute } \mathbb{I}(N_0; X_1, X_2, \dots, X_{N_0}) \text{ and } \mathbb{I}(N_1; X_1, X_2, \dots, X_{N_1}). \quad [4]$$

- Consider a sequence of n binary random variables (X_1, X_2, \dots, X_n) . Each sequence with an even number of 1's has probability $2^{-(n-1)}$, and each sequence with an odd number of 1's has probability 0. Compute the following:

$$\mathbb{I}(X_1; X_2), \quad \mathbb{I}(X_1; X_3|X_1), \quad \mathbb{I}(X_3; X_4|X_1, X_2), \quad \dots, \quad \mathbb{I}(X_{n-1}; X_n|X_1, X_2, \dots, X_{n-2}). \quad [4]$$

Solution:

- By definition:

$$\begin{aligned} \mathbb{I}(N_0; X_1, X_2, \dots, X_{N_0}) &= \mathbb{H}(N_0) - \mathbb{H}(N_0|X_1, X_2, \dots, X_{N_0}) = \mathbb{H}(N_0) \\ &= \frac{1}{3} \log 3 + \frac{2}{3} \log \left(\frac{3}{2} \right) = \log 3 - \frac{2}{3}. \end{aligned}$$

Similarly, $\mathbb{I}(N_1; X_1, X_2, \dots, X_{N_1}) = \mathbb{H}(N_1)$. We only need to calculate $\mathbb{H}(N_1)$. By problem 4 b) of HW1, $\mathbb{H}(N_1) = 2$.

- We will first show that X_1, X_2, \dots, X_{n-1} are mutually independent. Note that for all $x_1, \dots, x_{n-1} \in \{0, 1\}^{n-1}$,

$$\begin{aligned} \Pr\{X_1 = x_1, \dots, X_{n-1} = x_{n-1}\} &= \Pr\{X_1 = x_1, \dots, X_{n-1} = x_{n-1}, X_n = 0\} \\ &\quad + \Pr\{X_1 = x_1, \dots, X_{n-1} = x_{n-1}, X_n = 1\} \\ &= 2^{-(n-1)} \end{aligned}$$

since either $x_1 \oplus \dots \oplus x_{n-1} \oplus 0 = 0$ or $x_1 \oplus \dots \oplus x_{n-1} \oplus 1 = 0$.

*with contribution by Chen-Hao Hsiao and Wen-Shao Ho

1

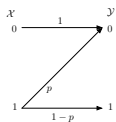
$$\begin{aligned} &= \max_{P_X} \{ \mathbb{H}(I) + \mathbb{I}(X^{(q)}; Y^{(q)}|I) \} & I \text{ is deterministic of } X^{(q)} \\ &= \max_{P_X} \sum_{i=1}^t P_X(i) \left[\log \frac{1}{P_X(i)} + \mathbb{I}(X^{(q)}; Y^{(q)}|I = i) \right] \\ &= \max_{P_X} \sum_{i=1}^t P_X(i) \log \frac{2^{C(i)}}{P_X(i)} & \text{Jensen's inequality} \\ &= \log \sum_{i=1}^t 2^{C(i)} \end{aligned}$$

The equality of the Jensen's inequality holds iff $\frac{2^{C(i)}}{2^{C(i)}} = \text{constant}$, $\forall i = 1, \dots, t$.

Let $P_X^{(i)}(x)$ be the optimal input probability distribution for the i -th channel. The optimal input probability distribution for the sum channel is

$$P_{X^{(q)}}(x) = \sum_{i=1}^t \Pr\{X^{(q)} = x | I = i\} P_X^{(i)}(x) = \frac{\sum_{i=1}^t 2^{C(i)} P_X^{(i)}(x) \mathbb{I}(I(x) = i)}{\sum_{i=1}^t 2^{C(i)}} \quad \forall x \in X^{(q)}.$$

4. (Z channel) [14]



The Z channel (depicted above) is one of the simplest asymmetric channel with its channel law described as follows.

$$P_{Y|X} = \begin{bmatrix} 1 & 0 \\ p & 1-p \end{bmatrix}.$$

In the following, let us assume $p = 1/2$.

- Find the capacity of the Z channel and a capacity achieving input distribution P_X^* . Also find P_Y^* , the output distribution induced by the input distribution P_X^* . [8]
- Is the capacity achieving input distribution of the Z channel unique? [2]
- Recall that $C = \mathbb{D}(P_{Y|X} \| P_Y^* | P_X^*)$ and can be viewed as a weighted average of $\{\mathbb{D}(P_{Y|X}(\cdot|i) \| P_Y^*(\cdot)) \mid i \in \mathcal{X}\}$. For the Z channel, derive $\mathbb{D}(P_{Y|X}(\cdot|0) \| P_Y^*(\cdot))$ and $\mathbb{D}(P_{Y|X}(\cdot|1) \| P_Y^*(\cdot))$. [4]

4

Under the constraint that $P \in \mathcal{P}(\mathcal{N})$ and $\mathbb{E}_{X \sim P}[X] = \sum_{i=1}^n x P(x) = \mu > 1$, find the minimum value of $\mathbb{D}(P|G(p))$ and a minimizing distribution. [12]

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

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

where $\mathcal{P}(\mathcal{X})$ denotes the collection of probability distributions over \mathcal{X} . Find a minimizer to the above problem. [10]

Solution:

- Expanding the Kullback-Leibler divergence over countable alphabet gives:

$$\begin{aligned} \mathbb{D}(P|G(p)) &= \sum_{x=1}^{\infty} P(x) \log \frac{P(x)}{(1-p)p^{x-1}} \\ &= \sum_{x=1}^{\infty} P(x) \log P(x) - \log(1-p) - \log p \sum_{x=1}^{\infty} P(x)(x-1) \\ &= -\mathbb{H}(P) - \log(1-p) - (\mu-1) \log p \\ &\geq -\mu \log \left(\frac{1}{\mu} \right) - \log(1-p) - (\mu-1) \log p, \end{aligned}$$

where the inequality comes from Problem 1 of Homework 2, and $P^* = G(1-\mu^{-1})$ is the minimizing distribution.

- Let $\bar{P} = \frac{1}{m} \sum_{i=1}^m P_i$, $\bar{P} \in \mathcal{P}(\mathcal{X})$.

$$\forall Q \in \mathcal{P}(\mathcal{X}),$$

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

Hence, \bar{P} is a minimizer.

5

Therefore, for $\mathcal{I} \subseteq \{1, 2, \dots, n-1\}$, by calculating the marginal distribution, we have

$$\Pr \left\{ \bigcap_{i \in \mathcal{I}} X_i = x_i \right\} = \prod_{i \in \mathcal{I}} \Pr\{X_i = x_i\} = 2^{-|\mathcal{I}|}.$$

By the mutual independence derived above, for $k = 1, 2, \dots, n-2$,

$$\begin{aligned} \mathbb{I}(X_k; X_{k+1}|X_1, \dots, X_{k-1}) &= \mathbb{H}(X_k|X_1, \dots, X_{k-1}) - \mathbb{H}(X_k|X_1, \dots, X_{k-1}, X_{k+1}) \\ &= \mathbb{H}(X_k) - \mathbb{H}(X_k) = 0. \end{aligned}$$

Also,

$$\begin{aligned} \mathbb{I}(X_{n-1}; X_n|X_1, \dots, X_{n-2}) &= \mathbb{H}(X_{n-1}|X_1, \dots, X_{n-2}) - \mathbb{H}(X_{n-1}|X_1, \dots, X_{n-2}, X_n) \\ &= \mathbb{H}(X_{n-1}) - \mathbb{H}(X_1 \oplus \dots \oplus X_{n-2} \oplus X_n|X_1, \dots, X_{n-2}, X_n) \\ &= 1 - 0 = 1. \end{aligned}$$

We use $\mathbb{H}(f(X)|X) = 0$ in the second equation.

2. (Data processing) [12]

- Let $X_1 - X_2 - X_3 - X_4$ form a Markov chain. Prove that

$$\mathbb{I}(X_1; X_3) + \mathbb{I}(X_2; X_4) \leq \mathbb{I}(X_1; X_4) + \mathbb{I}(X_2; X_3). \quad [6]$$

- Let $X_1 - X_2 - (X_3, X_4)$ form a Markov chain. Prove that

$$\mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_4) \leq \mathbb{I}(X_1; X_2) + \mathbb{I}(X_3; X_4). \quad [6]$$

Solution:

-

$$\begin{aligned} \mathbb{I}(X_1; X_3) + \mathbb{I}(X_2; X_4) &= \mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_3; X_4) - \mathbb{I}(X_1; X_4|X_3) & \text{chain rule} \\ &= \mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_3; X_4) & \text{markov} \\ &= \mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_4) + \mathbb{I}(X_3; X_4|X_1) & \text{chain rule} \\ &\leq \mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_4) + \mathbb{I}(X_3; X_4|X_1) & \text{data processing} \\ &= \mathbb{I}(X_1; X_4) + \mathbb{I}(X_3; X_4) & \text{chain rule} \end{aligned}$$

$$\begin{aligned} \mathbb{I}(X_1; X_3) + \mathbb{I}(X_1; X_4) &= \mathbb{I}(X_1; X_3; X_4) - \mathbb{I}(X_1; X_4|X_3) + \mathbb{I}(X_1; X_4) & \text{chain rule} \\ &\leq \mathbb{I}(X_1; X_2) - \mathbb{I}(X_1; X_4|X_3) + \mathbb{I}(X_1; X_4) & \text{data processing} \\ &= \mathbb{I}(X_1; X_2) + \mathbb{I}(X_3; X_4) & \text{chain rule} \end{aligned}$$

2

Solution:

- Note that \mathcal{X} is an alphabet of size 2, hence we can model any input distribution as $P_X = \text{Ber}(q)$, $0 \leq q \leq 1$. For given $X \sim \text{Ber}(q)$, we then have $Y \sim \text{Ber}((1-p)q)$.

$$\begin{aligned} C(q) &:= \mathbb{I}(X; Y) = \mathbb{H}(Y) - \mathbb{H}(Y|X) \\ &= \mathbb{h}_b((1-p)q) - P_X(0)\mathbb{H}(Y|X=0) - P_X(1)\mathbb{H}(Y|X=1) \\ &= \mathbb{h}_b((1-p)q) - q\mathbb{h}_b(p) \end{aligned}$$

$$\begin{aligned} \frac{d}{dq} C(q) &= \frac{d}{dq} \mathbb{h}_b((1-p)q) - q\mathbb{h}_b(p) \\ &= (1-p) \log \frac{1 - (1-p)q}{(1-p)q} - \mathbb{h}_b(p) \end{aligned}$$

$$C'(0) = 0$$

$$C'(1) = \mathbb{h}_b(1-p) - \mathbb{h}_b(p)$$

by the extreme value theorem, if we restrict $p = 1/2$, it becomes clear that $q = 2/5$ is the unique maximizer. Hence

$$C = \max_X \mathbb{I}(X; Y) = \max_{P_X} C(q) = C(2/5) = \mathbb{h}_b(1/5) - 2/5 = \log 5 - 2.$$

And $P_X^* = \text{Ber}(2/5)$ and $P_Y^* = \text{Ber}(1/5)$.

- Yes, it is unique.
-

$$\begin{aligned} \mathbb{D}(P_{Y|X}(\cdot|0) \| P_Y^*(\cdot)) &= 1 \log \frac{1}{1-1/5} + 0 \log \frac{0}{1/5} = \log 5 - 2 \\ \mathbb{D}(P_{Y|X}(\cdot|1) \| P_Y^*(\cdot)) &= \frac{1}{2} \log \frac{1/2}{1-1/5} + \frac{1}{2} \log \frac{1/2}{1/5} = \log 5 - 2 \\ \mathbb{D}(P_{Y|X} \| P_Y^* | P_X^*) &= \log 5 - 2 \end{aligned}$$

5

Grading Policy

- Find the tight lower bound [8] and find minimizer to achieve minimum [4].
- Specify the minimizer [3] and justify it [7].

6

$$\begin{aligned} &-(\mathbb{I}(X_1; X_4) + \mathbb{I}(X_3; X_4|X_1) - \mathbb{I}(X_3; X_4)) & \text{chain rule} \\ &\leq \mathbb{I}(X_1; X_4) + \mathbb{I}(X_3; X_4) & \text{nonnegative} \end{aligned}$$

3. (Sum Channel) [16]

Consider i DMCS

$$\left\{ (X^{(i)}, P_{Y|X}^{(i)}, Y^{(i)}) \mid i = 1, 2, \dots, I \right\},$$

where DMCS $(X^{(i)}, P_{Y|X}^{(i)}, Y^{(i)})$ has channel capacity $C^{(i)}$, for $1 \leq i \leq I$. The channel input alphabets are disjoint, and so are the channel output alphabets, that is,

$$X^{(i)} \cap X^{(j)} = Y^{(i)} \cap Y^{(j)} = \emptyset, \quad \forall i \neq j.$$

The **sum channel** $(X^{(q)}, P_{Y|X}^{(q)}, Y^{(q)})$ associated to these channels is defined as follows:

- Input alphabet is the union $X^{(q)} := \bigcup_{i=1}^I X^{(i)}$ of the individual input alphabets.
- Output alphabet is the union $Y^{(q)} := \bigcup_{i=1}^I Y^{(i)}$ of the respective output alphabets.
- At each time slot the transmitter chooses to use *one and only one* of the I channels to transmit a symbol, that is,

$$P^{(q)}(y|x) := \begin{cases} P_{Y|X}^{(i)}(y|x), & \text{if } x \in X^{(i)} \text{ and } y \in Y^{(i)} \\ 0, & \text{otherwise} \end{cases}$$

- Introduce a random variable I indicating which DMCS is used in the sum channel, that is,

$$I = i \quad \text{if } X \in X^{(i)}, i = 1, 2, \dots, I.$$

Show that for the sum channel $P_{Y|X}^{(q)}(X; Y) = \mathbb{I}(X; Y|I) + \mathbb{H}(I)$. [4]

- Find the capacity of the sum channel in terms of $\{C^{(i)} \mid i = 1, 2, \dots, I\}$. [6]
- Find the optimal input probability distribution for the sum channel in terms of the optimal input probability distributions for the individual channels. [6]

Solution:

- Let $X^{(q)}$ be the r.v. of input sum channel
- Let $Y^{(q)}$ be the r.v. of output sum channel

$$\begin{aligned} C^{(q)} &= \max_{P_X^{(q)}} \mathbb{I}(X^{(q)}; Y^{(q)}) \\ &= \max_{P_X^{(q)}} \mathbb{I}(X^{(q)}; Y^{(q)}, I) & X^{(q)} \text{ disjoint} \\ &= \max_{P_X^{(q)}} \{ \mathbb{I}(X^{(q)}; I) + \mathbb{I}(X^{(q)}; Y^{(q)}|I) \} & \text{chain rule} \end{aligned}$$

3

Homework 5

TA: Heng-Chien Liou*

1. (Information divergence) [10]

Consider two probability density functions $f(\cdot)$ and $g(\cdot)$. Let μ_1 and μ_2 denote the mean of f and g respectively. Let σ_f^2 and σ_g^2 denote the variance of f and g respectively.

- Compute $\mathbb{D}(f|g)$ in the following cases: (1) both f and g are Gaussian; (2) both f and g are Laplace. [6]
- If $\mu_1 = \mu_2$, which of the above cases gives the largest/smallest KL divergence? Your answer may depend on σ_1, σ_2 . [2]
- If $\sigma_1 = \sigma_2$, which of the above cases gives the largest/smallest KL divergence? Your answer may depend on μ_1, μ_2 . [2]

Solution:

- For f and g both being Gaussian, we have

$$\begin{aligned} \mathbb{D}(f|g) &= \int_{-\infty}^{\infty} f(x) \log \frac{f(x)}{g(x)} dx \\ &= \int_{-\infty}^{\infty} f(x) \left(\frac{(x-\mu_1)^2}{2\sigma_1^2} \log e - \frac{(x-\mu_2)^2}{2\sigma_2^2} \log e + \log \frac{\sigma_2}{\sigma_1} \right) dx \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{1}{2} \log e + \int_{-\infty}^{\infty} f(x) \frac{(x-\mu_2)^2}{2\sigma_2^2} \log e dx \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{1}{2} \log e + \frac{1}{2\sigma_2^2} (\sigma_f^2 + \mu_1^2 - 2\mu_1\mu_2 + \mu_2^2) \log e \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_f^2 - \sigma_g^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} \log e \end{aligned}$$

For f and g both being Laplace, we have

$$\$$

$$\begin{aligned}
 &= \log \frac{\sigma_2}{\sigma_1} - \log e + \log e \int_{\mu_1}^{\mu_2} \frac{\sqrt{2} |x - \mu_2|}{\sigma_2} dx \\
 \text{For } \mu_1 \geq \mu_2 & \\
 &= \int_{\mu_1}^{\mu_2} \frac{\sqrt{2} |x - \mu_2|}{\sigma_2} dx \\
 &= \int_{\mu_1}^{\mu_2} \frac{|x - \mu_2|}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx \\
 &= \int_{-\infty}^{\mu_2} \frac{-(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx + \int_{\mu_2}^{\mu_1} \frac{(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx \\
 &\quad + \int_{\mu_1}^{\mu_2} \frac{(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx \\
 &= -\frac{\sqrt{2}(x - \mu_2) - \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{-\infty}^{\mu_2} + \frac{\sqrt{2}(x - \mu_2) - \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{\mu_2}^{\mu_1} \\
 &\quad - \frac{\sqrt{2}(x - \mu_2) + \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{\mu_1}^{\mu_2} \\
 &= \frac{\sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} + \frac{\sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} + \frac{\sqrt{2}(\mu_1 - \mu_2) - \sigma_1}{2\sigma_2} + \frac{\sqrt{2}(\mu_1 - \mu_2) + \sigma_1}{2\sigma_2} \\
 &= \frac{\sigma_1}{\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} + \frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2} = \frac{\sigma_1}{\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} + \frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}
 \end{aligned}$$

Similarly, for $\mu_1 < \mu_2$

$$\begin{aligned}
 &= \int_{\mu_1}^{\mu_2} \frac{\sqrt{2} |x - \mu_2|}{\sigma_2} dx \\
 &= \int_{-\infty}^{\mu_1} \frac{-(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx + \int_{\mu_1}^{\mu_2} \frac{-(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx \\
 &\quad + \int_{\mu_2}^{\mu_1} \frac{(x - \mu_2)}{\sigma_1 \sigma_2} e^{\frac{\sqrt{2} |x - \mu_2|}{\sigma_2}} dx \\
 &= -\frac{\sqrt{2}(x - \mu_2) - \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{-\infty}^{\mu_1} + \frac{\sqrt{2}(x - \mu_2) - \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{\mu_1}^{\mu_2} \\
 &\quad - \frac{\sqrt{2}(x - \mu_2) + \sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(x - \mu_2)}{\sigma_2}} \Big|_{\mu_2}^{\mu_1} \\
 &= \frac{\sqrt{2}(\mu_1 - \mu_2) - \sigma_1}{2\sigma_2} - \frac{\sqrt{2}(\mu_1 - \mu_2) + \sigma_1}{2\sigma_2} + \frac{\sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} + \frac{\sigma_1}{2\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} \\
 &= -\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2} + \frac{\sigma_1}{\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} = \frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2} + \frac{\sigma_1}{\sigma_2} e^{\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}}
 \end{aligned}$$

AWGN channel with average output power constraint

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[Y_i^2] \leq B.$$

where $Y = X + Z$, $Z \perp X$, and $Z \sim \mathcal{N}(0, \sigma^2)$.

Evaluate the channel capacity $C(B)$.

[4]

Solution:

a) Follow the hint, let $\hat{h}(x) := \mathbb{E}[h(x, Y)]$, we have

$$\frac{1}{n} \sum_{i=1}^n \hat{h}(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y_i} [\hat{h}(x_i, Y_i)] \leq B.$$

Therefore, it can be viewed as an input-only cost constraint.

Furthermore, the capacity is

$$\begin{aligned}
 C(B) &= \max_{P_{X, \mathbb{E}_{P_{X,Y}}[h(X,Y)] \leq B}} I(X; Y) = \max_{P_{X, \mathbb{E}_{P_{X,Y}}[h(X,Y)] \leq B}} \max_{P_{Y|X} \in \mathcal{P}_{Y|X}[\mathbb{E}_{P_{X,Y}}[h(X,Y)] \leq B]} I(X; Y) \\
 &= \max_{P_{X, \mathbb{E}_{P_{X,Y}}[h(X,Y)] \leq B}} I(X; Y).
 \end{aligned}$$

b) The capacity can be directed calculated via a), note that for random variable Y satisfying $\mathbb{E}_{P_{Y,Y}}[Y^2] \leq B$:

$$h(Y) - h(Y|X) \leq \frac{1}{2} \log(2\pi e B) - h(Y) = \frac{1}{2} \log\left(\frac{B}{\sigma^2}\right).$$

The equality can be achieved by choosing $X \sim \mathcal{N}(0, B - \sigma^2)$. Hence,

$$C(B) = \frac{1}{2} \log\left(\frac{B}{\sigma^2}\right)$$

4. (Compression with guarantee on the cross-entropy loss) [20]

Consider a discrete memoryless source $S \sim \pi$ with a finite alphabet $\mathcal{S} = \{1, 2, \dots, k\}$, $|\mathcal{S}| = k < \infty$. The encoder aims to compress the source so that the decoder can give good estimates of the source sequence. In many applications, however, the decoder may not want to give a deterministic estimate. Instead, for each symbol s_k in a length- n sequence s^n , its goal is to produce a *probability vector* \mathbf{q}_k in the k -dimensional probability simplex \mathcal{P}_k , where the i -th coordinate, $q_i(\ell)$, stands for the probability of $s_i = \ell$ that the decoder believes in based on what it receives from the encoder. A standard way to quantify the loss is the empirical cross entropy loss

$$\ell_{\text{CE}}(s^n, \mathbf{q}^n) = \sum_{i=1}^n \frac{1}{n} \log \frac{1}{q_i(s_i)}$$

$$\begin{aligned}
 &\leq \min_{(S, \hat{S})} \left\{ I(S; \hat{S}) \left| \log \frac{1}{P_{\hat{S}|S}(S|\hat{S})} \right| \leq D \text{ and } S \sim \pi \right\} \\
 &= \min_{(S, \hat{S})} \left\{ I(S; \hat{S}) \left| H(\hat{S}|\hat{S}) \leq D \text{ and } S \sim \pi \right. \right\} \\
 \text{d) } R(D) &= H(\pi) - D \text{ since} \\
 I(S; \hat{S}) &= H(S) - H(S|\hat{S}) \\
 &= H(\pi) - H(S|\hat{S}) \\
 &\geq H(\pi) - D.
 \end{aligned}$$

The equality holds since $0 \leq D \leq H(\pi)$, $\exists P_{\hat{S}|S}$ s.t. $H(\hat{S}|\hat{S}) = D$.

Hence,

$$D(f|g) = \log \frac{\sigma_2}{\sigma_1} - \log e + \log e \left(\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2} + \frac{\sigma_1}{\sigma_2} e^{-\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma_2}} \right)$$

b) For $\mu_1 = \mu_2$

$$D(f|g) = \begin{cases} \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 - \sigma_2^2}{2\sigma_2^2} \log e & \text{Gaussian} \\ \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1}{\sigma_2} \log e - \log e & \text{Laplace} \end{cases}$$

Simple calculations based on σ_2/σ_1 and the property of quadratic forms shows that pair of Gaussian always admit a larger divergence if $\mu_1 = \mu_2$.

c) For $\sigma_1 = \sigma_2 = \sigma$,

$$D(f|g) = \begin{cases} \frac{(\mu_1 - \mu_2)^2}{2\sigma^2} \log e & \text{Gaussian} \\ -\log e + \log e \left(\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma} + e^{-\frac{\sqrt{2}(\mu_1 - \mu_2)}{\sigma}} \right) & \text{Laplace} \end{cases}$$

Let $y = |\mu_1 - \mu_2|/\sigma$, then we can compare

$$D(f|g) = \begin{cases} \frac{y^2}{2} \log e & \text{Gaussian} \\ -\log e + \log e \left(\sqrt{2}y + e^{-\sqrt{2}y} \right) & \text{Laplace} \end{cases}$$

The equation $\frac{y^2}{2} - (-1 + \sqrt{2}y + e^{-\sqrt{2}y}) = 0$ has two real solutions, one of them being zero. We denote another one as $y_0 > 0$. For $|\mu_1 - \mu_2|/\sigma \leq y_0$, a pair of Laplace distribution admit a divergence greater than or equal to that of Gaussian distribution. For $|\mu_1 - \mu_2|/\sigma > y_0$, a pair of Gaussian has a larger divergence.

2. (Differential entropy) [10]

- Consider a Laplace random variable $X \sim \text{Lap}(\mu, b)$, that is, the probability density function of X is $f_X(x) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}$, $x \in \mathbb{R}$. Compute its differential entropy $h(X)$. [4]
- Consider a problem of maximizing differential entropy $h(X)$ subject to the constraint that $\mathbb{E}[|X|] \leq B$. Find the maximum differential entropy and show that a zero-mean Laplace distributed X attains the maximum value. [6]

Solution:

a)

$$\begin{aligned}
 h(X) &= \mathbb{E}_X \left[\log \frac{1}{f_X(x)} \right] \\
 &= \log 2b + \log e \left(\int_{-\infty}^{\mu} \frac{\mu - x}{2b^2} e^{-\frac{\mu-x}{b}} dx + \int_{\mu}^{\infty} \frac{x - \mu}{2b^2} e^{-\frac{x-\mu}{b}} dx \right)
 \end{aligned}$$

Note that it can be viewed as the average distortion per symbol when the distortion function is set to be

$$d: \mathcal{S} \times \mathcal{P}_d \rightarrow [0, \infty), (s, q) \mapsto d(s, q) = \log \frac{1}{q(s)}.$$

Hence, one can study a lossy source coding problem to understand how to represent a memoryless source with the smallest rate so that the decoder can declare an estimation probability vector with the empirical cross entropy loss not greater than a prescribed level D . By the lossy source coding theorem, the rate is given by the following rate distortion function:

$$R(D) = \inf_{\{I(S; Q) \mid \mathbb{E} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \}}$$

a) Show that for the lossy source coding problem, $D_{\min} = 0$ and $D_{\max} = H(\pi)$.

b) Show that for any jointly distributed $(S, Q) \sim P$,

$$H(S|Q) \leq \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right].$$

Then, argue that $R(D) \geq H(\pi) - D$, for $0 \leq D \leq H(\pi)$.

c) Show that for $0 \leq D \leq H(\pi)$,

$$R(D) \leq \min_{(S, \hat{S}), S \sim \pi} \left\{ I(S; \hat{S}) \mid H(\hat{S}|\hat{S}) \leq D \text{ and } S \sim \pi \right\}.$$

d) Show that for $0 \leq D \leq H(\pi)$,

$$R(D) = \min_{(S, \hat{S})} \left\{ I(S; \hat{S}) \mid H(\hat{S}|\hat{S}) \leq D \text{ and } S \sim \pi \right\} = H(\pi) - D.$$

Hence, $R(D) = \max(0, H(\pi) - D)$.

Solution:

- For D_{\min} , one could choose $q(s) = \begin{cases} 1 & \text{if } s \\ 0 & \text{otherwise} \end{cases}$
 $\mathbb{E}_q[d(S, q(S))] = \mathbb{E}_q[\log 1] = 0$, so $D_{\min} = 0$.
 For D_{\max} , since

$$\begin{aligned}
 \mathbb{E}_S \left[\log \frac{1}{\pi(S)} \right] - \mathbb{E}_S \left[\log \frac{1}{q(S)} \right] &= \mathbb{E}_S \left[\log \frac{q(S)}{\pi(S)} \right] \\
 &\leq \log \mathbb{E}_S \left[\frac{q(S)}{\pi(S)} \right] = 0 \quad \forall q.
 \end{aligned}$$

$$D_{\max} = \min_q \mathbb{E}_q[d(S, q)] = \mathbb{E}_S \left[\log \frac{1}{\pi(S)} \right] = H(\pi).$$

Information Theory HW1

許博翔

September 21, 2023

Problem 1. I'll prove (b) first, and then use (b) to prove (a) for convenience.

- Suppose that $s^n = (s_1, s_2, \dots, s_n) \in \mathcal{T}^{(n)}(S)$.
 By the definition of $\mathcal{T}^{(n)}(S)$, $\forall a \in \mathcal{S}$, $\left| \frac{1}{n} \sum_{i=1}^n \mathbb{I}(s_i = a) - P_S(a) \right| \leq \gamma P_S(a)$.
 $\Rightarrow \forall a \in \mathcal{S}$, $\left| \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) \mathbb{I}(s_i = a) - P_S(a) \log P_S(a) \right| \leq \gamma P_S(a) \log P_S(a)$.
 $\Rightarrow \sum_{a \in \mathcal{S}} \left| \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) \mathbb{I}(s_i = a) - P_S(a) \log P_S(a) \right| \leq \sum_{a \in \mathcal{S}} \gamma P_S(a) \log P_S(a)$.
 By triangular inequality,
 $\left| \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) + H(S) \right|$
 $= \left| \sum_{a \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) \mathbb{I}(s_i = a) - \sum_{a \in \mathcal{S}} P_S(a) \log P_S(a) \right|$
 $\leq \sum_{a \in \mathcal{S}} \left| \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) \mathbb{I}(s_i = a) - P_S(a) \log P_S(a) \right|$
 $\leq \sum_{a \in \mathcal{S}} \gamma P_S(a) \log P_S(a) = -\gamma H(S)$.

Taking $d = \xi(\gamma) = -\gamma H(S)$, and we get $\left| \frac{1}{n} \sum_{i=1}^n \log P_S(s_i) + H(S) \right| \leq d$, which means $s^n \in \mathcal{A}_d^{(n)}(S)$.
 $\therefore \mathcal{T}_d^{(n)}(S) \subseteq \mathcal{A}_d^{(n)}(S)$.

(a) Recall from (b), we take $d = \xi(\gamma) = -\gamma H(S)$.

The 4 properties in the proposition are:

- The original property is: $\forall s^n \in \mathcal{A}_d^{(n)}(S)$, $2^{-n(H(S)+d)} \leq \Pr\{S^n = s^n\} \leq 2^{-n(H(S)-d)}$.

$$\begin{aligned}
 &= \log 2b + \log e \left(\int_0^{\infty} \frac{t}{2} e^{-t} dt + \int_0^{\infty} \frac{t}{2} e^{-t} dt \right) \\
 &= \log 2b + \log e \left(\int_0^{\infty} t e^{-t} dt \right) \\
 &= \log 2b + \log e \left(-te^{-t} \Big|_0^{\infty} + \int_0^{\infty} e^{-t} dt \right) \\
 &= \log 2be
 \end{aligned}$$

We can also observe that translation doesn't change differential entropy.

b) Let $X \sim \text{Lap}(0, b)$. Its differential entropy $h(X) = \log(2be)$ and $\mathbb{E}[|X|] = b$.

Let Y be a random variable with $\mathbb{E}[|Y|] = b$.

Then we have $h(X) = \mathbb{E}_X[-\log f_X(X)] = \mathbb{E}_X\left[-\frac{1}{b} + \left|\frac{X}{b}\right|\right] = \mathbb{E}_Y[-\log f_Y(Y)]$.

Then

$$\begin{aligned}
 h(Y) - h(X) &= \mathbb{E}_Y \left[\log \frac{f_X(Y)}{f_Y(Y)} \right] \\
 &\leq \log \mathbb{E}_Y \left[\frac{f_X(Y)}{f_Y(Y)} \right] \quad \text{Jensen's inequality} \\
 &= 0
 \end{aligned}$$

$\text{Lap}(0, b)$ maximizes $h(Y)$ for given $\mathbb{E}[|Y|] = b$, and the maximum value is $\log(2be)$.

Since $\log(2be)$ is increasing in $0 \leq b \leq B$, the Laplace distribution $\text{Lap}(0, B)$ also maximizes $h(S)$ for given $\mathbb{E}[|S|] \leq B$, and the maximum differential entropy is $\log(2Be)$.

3. (Channel Coding with Input-Output Cost Constraint) [10]

In this problem we explore channel coding with input and output cost constraint.

- Consider a DMC $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$. Let $b: \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$ be an input-output cost function. Suppose the channel coding has to satisfy the following average cost constraint: for each codeword x^n ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y_i} [b(x_i, Y_i)] \leq B.$$

Note that Y_i follows distribution $P_{Y|X}(\cdot|x_i)$.

Argue that the problem is equivalent to another channel coding problem with a properly defined input-only cost function. Show that the capacity-cost function is

$$C(B) = \max_{P_{X, \mathbb{E}_{P_{X,Y}}[b(X,Y)] \leq B} I(X; Y). \quad [6]$$

Hint: Consider the input-only cost function $\hat{b}(x) := \mathbb{E}[b(x, Y)]$, and check that the steps in the proof of DMC with input cost in the lecture are still valid.

b) Using discretization techniques, the above DMC result can be extended to continuous memoryless channels. With the extension (no need to prove it here), let us consider an

b) Denote the law of P as $P_{S,Q}$ and its marginal as Q as P_Q

$$\begin{aligned}
 H(S|Q) - \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right] &= h(S; Q) - h(Q) - \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right] \\
 &= \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{Q(S)P_Q(Q)}{P_{S,Q}(S, Q)} \right] \\
 &\leq \log \mathbb{E}_{(S,Q) \sim P} \left[\frac{Q(S)P_Q(Q)}{P_{S,Q}(S, Q)} \right] \\
 &= \log \left(\sum_{s \in \mathcal{S}} \int_{q \in \mathcal{P}_Q} q(s) P_Q(q) \right) = 0
 \end{aligned}$$

$$\begin{aligned}
 R(D) &= \inf_{(S, Q)} \left\{ I(S; Q) \mid \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
 &= H(\pi) - \sup \left\{ H(S|Q) \mid \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
 &\geq H(\pi) - D, \quad \forall D_{\min} = 0 \leq D \leq D_{\max} = H(\pi)
 \end{aligned}$$

c) The key to this problem is to make the observation that, for all pair of random variable (S, \hat{S}) over $S \times \mathcal{S}$, we can associate a random vector $\mathbf{Q}_{S, \hat{S}}$ over \mathcal{P}_S such that

$$\begin{aligned}
 \mathbf{Q}_{S, \hat{S}} &= \mathbf{q} = P_{\hat{S}|S}(\cdot|s) \text{ if } S = s \\
 \text{equivalently } \mathbf{Q}_{S, \hat{S}} &= \sum_{i=1}^k P_{\hat{S}|S}(\cdot|i) \mathbb{I} \{ \hat{S} = i \}
 \end{aligned}$$

a mixture of k vectors \mathbf{q}_i determined by the outcome of S . Note that this specify a Markov chain $S - S - \mathbf{Q}_{S, \hat{S}}$. For this Markov chain, by the data processing inequality, $I(S; \mathbf{Q}_{S, \hat{S}}) \leq I(S; \hat{S})$. And

$$\mathbb{E}_{(S, \hat{S})} \left[\log \frac{1}{Q_{S, \hat{S}}(S)} \right] = \mathbb{E}_{(S, \hat{S})} \left[\log \frac{1}{P_{\hat{S}|S}(S|\hat{S})} \right]$$

Hence,

$$\begin{aligned}
 R(D) &= \inf_{(S, Q)} \left\{ I(S; Q) \mid \mathbb{E}_{(S,Q) \sim P} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
 &= \inf_{(S, \hat{S})} \left\{ I(S; \hat{S}) \mid \mathbb{E}_{(S, \hat{S})} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
 &\leq \min_{(S, \hat{S})} \left\{ I(S; \hat{S}) \mid \mathbb{E}_{(S, \hat{S})} \left[\log \frac{1}{Q(S)} \right] \leq D \text{ and } S \sim \pi \text{ and } \mathbf{Q} = \mathbf{Q}_{S, \hat{S}} \right\}
 \end{aligned}$$

\therefore from (b) we know that $\forall s^n \in \mathcal{T}^{(n)}(S)$, $s^n \in \mathcal{A}^{(n)}(S)$.

$\therefore 2^{-n(H(S)+d)} \leq \Pr\{S^n = s^n\} \leq 2^{-n(H(S)-d)}$.

(2) Let $\mathcal{A}_d(n) := \{s^n \in \mathcal{S}^n : |s|(\mathbf{a}^n) - P_S(a)| \leq \gamma P_S(a)\}$.

Since $S \sim P_S$ is a DMS, the random variables $\{X_i\}_{i=1}^{\infty}$ where $X_i := \mathbb{I}\{S_i = a\}$ are i.i.d.

The average of X_i , denote as $\mu = \Pr\{S_i = a\} = P_S(a)$.

$$X_n = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{S_i = a\} = \pi(\mathbf{a}^{(n)}).$$

Take $\epsilon > \gamma P_S(a)$.

By the weak law of large numbers, $\lim_{n \rightarrow \infty} \Pr$

Problem 2.

(a) Define $X_i = \log \frac{1}{P_S(S_i)}$. Since S_i are i.i.d, X_i are also i.i.d.

$$\text{Since } P_S(S_i) \leq 1, \text{ we get that } \log \frac{1}{P_S(S_i)} \geq 0.$$

$$\Rightarrow \mathbb{E}[|X_i|] = \mathbb{E}[X_i] = \mathbb{E}\left[\log \frac{1}{P_S(S_i)}\right] = H(S) < \infty.$$

$$\begin{aligned} \prod_{i=1}^n P_S(S_i) &\geq 2^{-n(H(S) + \epsilon^{-1/2}k_S(S))} \\ \Leftrightarrow \prod_{i=1}^n \frac{1}{P_S(S_i)} &\leq 2^{n(H(S) + \epsilon^{-1/2}k_S(S))} \\ \Leftrightarrow \sum_{i=1}^n \log \frac{1}{P_S(S_i)} &\leq n(H(S) + \epsilon^{-1/2}k_S(S)) \\ \Leftrightarrow \left(\frac{1}{n} \sum_{i=1}^n X_i\right) - H(S) &\leq \epsilon^{-1/2}k_S(S) \\ \Leftrightarrow \frac{\sqrt{n}(\bar{X}_n - \mathbb{E}[X_i])}{\sqrt{n}} &\leq \delta. \end{aligned}$$

By central limit theorem, $\frac{\sqrt{n}(\bar{X}_n - \mathbb{E}[X_i])}{\sqrt{n}} \xrightarrow{d} Z \sim \mathcal{N}(0, 1)$ as $n \rightarrow \infty$.

$$\Rightarrow \Pr\left\{\prod_{i=1}^n P_S(S_i) \geq 2^{-n(H(S) + \epsilon^{-1/2}k_S(S))}\right\} = \Pr\left\{\frac{\sqrt{n}(\bar{X}_n - \mathbb{E}[X_i])}{\sqrt{n}} \leq \delta\right\}$$

$$\rightarrow \Pr(Z \leq \delta) = \Phi(\delta) \text{ as } n \rightarrow \infty.$$

(b) Let $Z \sim \mathcal{N}(0, 1)$, by Berry-Esseen theorem, $|\Pr\{S^n \in \mathcal{B}_\delta^{(n)}(S)\} - \Pr\{Z \leq \delta\}| = \left|\Pr\left\{\frac{\sqrt{n}(\bar{X}_n - \mathbb{E}[X_i])}{\sqrt{n}} \leq \delta\right\} - \Pr\{Z \leq \delta\}\right| \leq cn^{-1/2}$ for some constant $c > 0$.

$$\Rightarrow \Pr\{S^n \in \mathcal{B}_\delta^{(n)}(S)\} \geq \Pr\{Z \leq \delta\} - cn^{-1/2} = \Phi(\delta) - cn^{-1/2}.$$

Take $\delta = \Phi^{-1}(1 - \epsilon + cn^{-1/2}) = -\Phi^{-1}(\epsilon - cn^{-1/2})$, we get that $\Pr\{S^n \in \mathcal{B}_\delta^{(n)}(S)\} \geq 1 - \epsilon$.

$$\Rightarrow \Pr\{S^n \notin \mathcal{B}_\delta^{(n)}(S)\} \leq \epsilon.$$

Since $\frac{d\Phi^{-1}(x)}{dx} = \frac{1}{\frac{d\Phi}{d\Phi^{-1}(x)}} = \frac{1}{\frac{1}{\sqrt{2\pi}}e^{-(\Phi^{-1}(x))^2/2}} = \sqrt{2\pi}e^{(\Phi^{-1}(x))^2/2}$, there is $\Phi^{-1}(\epsilon - cn^{-1/2}) \approx \Phi^{-1}(\epsilon) - \sqrt{2\pi}e^{(\Phi^{-1}(\epsilon))^2/2}cn^{-1/2} = \Phi^{-1}(\epsilon) - O(n^{-1/2})$ for n sufficiently large.

$$\Rightarrow \delta = -\Phi^{-1}(\epsilon) + \zeta'_n, \text{ where } \zeta'_n = O(n^{-1/2}).$$

Lemma 2.1. $\exists \zeta_n = O(n^{-1})$ s.t. $nk \leq [n(\delta + \zeta_n)]$.

Proof. Consider $\zeta_n = \frac{1}{n}$, we get that $[n(\delta + \zeta_n)] = [nk + \frac{1}{n}] = [nk] + 1 \geq$

Author: 许博翔

3

Information Theory HW2

许博翔

October 5, 2023

Problem 1.

(a) Define $Q_X(x) = q_x$.

$$H(X) + \sum_{i=1}^r p_i \log q_i = \sum_{i=1}^r \mathbb{E}\left[\log \frac{Q_X}{P_X}\right] \leq \sum_{i=1}^r \log \mathbb{E}\left[\frac{Q_X}{P_X}\right] = \log \left(\sum_{i=1}^r p_i \frac{q_i}{p_i}\right) = \log \left(\sum_{i=1}^r q_i\right) = \log 1 = 0.$$

$$\therefore H(X) \leq -\sum_{i=1}^r p_i \log q_i.$$

(b) $-\log q_i$ is an arithmetic sequence $\Rightarrow q_i$ is a geometric sequence.

Suppose that $q_i = q_0 r^i$, where $1 < r < 1$ and $q_0 > 0$.

$$\therefore 1 = \sum_{i=1}^r q_i = \frac{q_0 r^i}{1-r}$$

$$\Rightarrow q_0 = \frac{1-r}{r}$$

$$\therefore \mu_X = \sum_{i=1}^r i q_i = \sum_{i=1}^r i q_0 r^i = \sum_{i=1}^r i q_0 r^i = \sum_{j=1}^r \sum_{i=j}^r q_0 r^i = \sum_{j=1}^r q_0 r^j = \frac{q_0 r^j}{(1-r)^2}$$

$$\Rightarrow \frac{1}{1-r} = \mu_X$$

$$\therefore r = 1 - \frac{1}{\mu_X} = \frac{\mu_X - 1}{\mu_X}, \quad q_0 = \frac{\frac{1}{\mu_X}}{1 - \frac{1}{\mu_X}} = \frac{1}{\mu_X - 1}.$$

$$-\log q_i = -\log q_0 r^i = -\log q_0 - i \log r.$$

Take $\alpha = -\log r = \log(\mu_X) - \log(\mu_X - 1)$, $\beta = -\log q_0 = \log(\mu_X - 1)$ satisfies the conditions.

$$\therefore \text{the answer is } q_i = \frac{(\mu_X - 1)^{i-1}}{(\mu_X)^i}, \quad \alpha = \log(\mu_X) - \log(\mu_X - 1), \quad \beta = \log(\mu_X - 1).$$

(c) $-\sum_{i=1}^{\infty} p_i \log q_i = \sum_{i=1}^{\infty} p_i (\alpha i + \beta) = \alpha \mu_X + \beta = \log(\mu_X) \mu_X - \log(\mu_X - 1) \mu_X + \log(\mu_X - 1) = \mu_X (\log(\mu_X) - (1 - \frac{1}{\mu_X}) \log(\mu_X - 1)) = \mu_X (\log(\frac{\mu_X}{\mu_X - 1}) - \frac{1}{\mu_X} \log \frac{\mu_X}{\mu_X - 1}) +$

Information Theory HW2

$$\begin{aligned} & - \sum_{x_1, \theta_1, x_2, \theta_2} P_{X_{01}|1}(\theta_1, x_1) P_{\theta_1|\theta_1}(\theta_2|\theta_1) P_{X_{02}|2}(x_2) \log(P_{X_{02}|2}(x_2)) \\ & = - \sum_{\theta_1, \theta_2, x_1} P_{X_{01}|1}(\theta_1, x_1) P_{\theta_1|\theta_1}(\theta_2|\theta_1) \log(P_{\theta_1|\theta_1}(\theta_2|\theta_1)) \\ & + \sum_{x_1, \theta_1, x_2} P_{X_{01}|1}(\theta_1, x_1) P_{\theta_1|\theta_1}(\theta_2|\theta_1) H(X_{02}|2) \\ & = - \sum_{\theta_1, \theta_2} P_{\theta_1}(\theta_1) P_{\theta_2|\theta_1}(\theta_2|\theta_1) \log(P_{\theta_2|\theta_1}(\theta_2|\theta_1)) \\ & + \sum_{\theta_1, \theta_2} P_{\theta_1}(\theta_1) P_{\theta_2|\theta_1}(\theta_2|\theta_1) H(X_{02}|2) \\ & = -(1-q)(1-\alpha) \log(1-\alpha) - (1-q)\alpha \log(\alpha) - q\beta \log(\beta) - q(1-\beta) \log(1-\beta) \\ & - H(X_{02}|2) - (1-q)(1-\alpha) + q\beta + H(X_{02}|2) - (1-q)\alpha + (1-\beta) \\ & (X_{02}|2) \text{ sum i.i.d. over } \theta_1, \theta_2 \text{ in } H(X_{02}|2) = H(X_{02}|2) \\ & (1-q)H_{02}(\alpha) + qH_{02}(\beta) + H_{02}((1-q)(1-\alpha) + q\beta) + H_{02}((1-q)\alpha + q(1-\beta)) \\ & = \frac{\beta}{\alpha + \beta} H_{02}(\alpha) + \frac{\alpha}{\alpha + \beta} H_{02}(\beta) + H_{02}\left(\frac{\beta}{\alpha + \beta}(1-\alpha) + \frac{\alpha}{\alpha + \beta}\beta\right) + H_{02}\left(\frac{\beta}{\alpha + \beta}\alpha + \frac{\alpha}{\alpha + \beta}(1-\beta)\right) \\ & = \frac{\beta}{\alpha + \beta} (H_{02}(\alpha) + H_{02}(\beta)) + H_{02}\left(\frac{\alpha}{\alpha + \beta}(H_{02}(\beta) + H_{02}(\alpha))\right). \end{aligned}$$

Author: 许博翔

4

nk.

$$\text{Since } \sum_{s^* \in \mathcal{B}_\delta^{(n)}(S)} P_{S^n}(s^*) \leq \sum_{s^*} P_{S^n}(s^*) = 1,$$

$$\text{and if } s^* \in B, \text{ then } P_{S^n}(s^*) = \prod_{i=1}^n P_{S_i}(s_i) \geq 2^{-n(H(S) + \epsilon^{-1/2}k_S(S))}.$$

$$\therefore |\mathcal{B}_\delta^{(n)}(S)| 2^{-n(H(S) + \epsilon^{-1/2}k_S(S))} = \sum_{s^* \in \mathcal{B}_\delta^{(n)}(S)} 2^{-n(H(S) + \epsilon^{-1/2}k_S(S))} \leq \sum_{s^* \in \mathcal{B}_\delta^{(n)}(S)} P_{S^n}(s^*) \leq 1.$$

$$\Rightarrow |\mathcal{B}_\delta^{(n)}(S)| \leq 2^{n(H(S) + \epsilon^{-1/2}k_S(S))}.$$

By **Lemma (2.1)**, there exists $\zeta'_n \in O(n^{-1})$ s.t. $n(H(S) + \epsilon^{-1/2}k_S(S)) \leq [n(H(S) + \epsilon^{-1/2}k_S(S) + \zeta'_n)]$.

Take $R = H(S) + \epsilon^{-1/2}k_S(S)\delta + \zeta'_n = H(S) - n^{-1/2}k_S(S)\Phi^{-1}(\epsilon) + n^{-1/2}k_S(S)\zeta'_n + \zeta'_n$.

Since $n^{-1/2}k_S(S)\zeta'_n = O(n^{-1})$, we get that $R = H(S) - n^{-1/2}k_S(S)\Phi^{-1}(\epsilon) + \zeta'_n$ for some $\zeta'_n = O(n^{-1})$.

Therefore, $\mathcal{B}_\delta^{(n)}(S)$ is an $(n, [nR])$ code with $P_{S^n}^{(n)} \leq \epsilon$.

Problem 3.

(a) Let $\delta \in (0, R - H(S))$, and $\mathcal{A}_\delta^{(n)}(S)$ be the δ -typical set defined in Definition 1.

By the third property of Proposition 1, we know that $|\mathcal{A}_\delta^{(n)}(S)| \leq 2^{n(H(S) + \delta)}$ $H(S) + \delta \leq R$ and $H(S) + \delta \leq R - 1$ for a large enough $2^{n\delta}$ for n large enough.

$$\Rightarrow \mathcal{A}_\delta^{(n)}(S) \text{ is an } (n, [nR]) \text{ code.}$$

By the second property of Proposition 1, we know that $\forall \epsilon > 0, \exists N$ s.t. $\forall n \geq N, P_{S^n}^{(n)} = \Pr\{S^n \notin \mathcal{A}_\delta^{(n)}(S)\} \leq \epsilon$.

Since $P_{S^n}^{(n)} \geq 0$, therefore by the definition of limits, $\lim_{n \rightarrow \infty} P_{S^n}^{(n)} = 0$.

\therefore such sequence exists, and it is $\mathcal{A}_\delta^{(n)}(S)$.

(b) For a given $(n, [nR])$ code, let $\mathcal{B}^{(n)}$ denote the range of the decoding function.

Let $\delta \in (0, H(S) - R)$, and $\mathcal{A}_\delta^{(n)}(S)$ be the δ -typical set defined in Definition 1.

By the first property of Proposition 1, we know that $\forall s^* \in \mathcal{A}_\delta^{(n)}(S), \Pr\{S^n = s^*\} \leq 2^{-n(H(S) - \delta)}$.

$$\Rightarrow \Pr\{S^n \in \mathcal{A}_\delta^{(n)}(S) \cap \mathcal{B}^{(n)}\} = \sum_{s^* \in \mathcal{A}_\delta^{(n)}(S) \cap \mathcal{B}^{(n)}} \Pr\{S^n = s^*\}$$

Author: 许博翔

4

Information Theory HW2

$$\begin{aligned} \frac{1}{\mu_X} \log(\mu_X) &= \mu_X \left(-\left(1 - \frac{1}{\mu_X}\right) \log\left(1 - \frac{1}{\mu_X}\right) - \frac{1}{\mu_X} \log\left(\frac{1}{\mu_X}\right)\right) = \mu_X h\left(\mu_X^{-1}\right). \\ \therefore H(X) &\leq \mu_X h\left(\mu_X^{-1}\right), \text{ and the equation holds when } p_i = q_i \text{ for all } i, \text{ that is,} \\ X &\sim \text{Geo}\left(\frac{1}{\mu_X}\right) \text{ is the geometric distribution.} \end{aligned}$$

Problem 2.

$$\begin{aligned} \text{(a)} \int_2^\infty \frac{1}{x(\log x)^\alpha} dx &= \int_{x=2}^\infty (\log x)^{-\alpha} d(\log x) \\ &= \begin{cases} \frac{1}{1-\alpha} (\log x)^{1-\alpha} \Big|_2^\infty, & \text{if } \alpha \neq 1, \text{ which converges } \Leftrightarrow 1-\alpha < 0 \Leftrightarrow \alpha > 1, \\ \log \log x \Big|_2^\infty, & \text{if } \alpha = 1, \text{ which does not converge} \end{cases} \\ \therefore \sum_{n=2}^\infty \frac{1}{n(\log n)^\alpha} &\text{ converges } \Leftrightarrow \alpha > 1. \end{aligned}$$

(b) First, we know that the series converges $\Leftrightarrow \alpha > 1$, so we only consider $\alpha > 1$.

$$H(X_n) = -\mathbb{E}[\log P_{X_n}] = \sum_{n=2}^\infty \frac{1}{s_n n (\log n)^\alpha} \log(s_n n (\log n)^\alpha) = \sum_{n=2}^\infty \frac{\log s_n}{s_n n (\log n)^\alpha} + \sum_{n=2}^\infty \frac{\alpha \log \log n}{s_n n (\log n)^{\alpha-1}}$$

$$\text{For } \alpha \leq 2, \text{ since } H(X_n) > \sum_{n=2}^\infty \frac{1}{s_n n (\log n)^{\alpha-1}} \rightarrow \infty \text{ from (a); therefore } H(X_n) \text{ diverges to } \infty.$$

$$\text{For } \alpha > 2, \text{ since } H(X_n) < \sum_{n=2}^\infty \frac{\log s_n}{s_n n (\log n)^\alpha} + \sum_{n=2}^\infty \frac{1}{s_n n (\log n)^{\alpha-1}} + \sum_{n=2}^\infty \frac{\alpha \log \log n}{s_n n (\log n)^{\alpha-1}}$$

$$\log \log n \leq \log n \text{ for } n \geq 2 \Rightarrow \sum_{n=2}^\infty \frac{\alpha \log \log n}{s_n n (\log n)^{\alpha-1}} \leq \sum_{n=2}^\infty \frac{1}{s_n n (\log n)^{\alpha-1}} + \sum_{n=2}^\infty \frac{\alpha}{s_n n (\log n)^{\alpha-1}} = \log \log s_n + \frac{(1+\alpha)s_n - 1}{s_n} < \infty,$$

$$\text{and } \sum_{n=2}^m \frac{1}{s_n n (\log n)^\alpha} \log(s_n n (\log n)^\alpha) \text{ is increasing as } m \text{ increases.}$$

$$\Rightarrow H(X_n) = \sum_{n=2}^\infty \frac{1}{s_n n (\log n)^\alpha} \log(s_n n (\log n)^\alpha) \text{ converges.}$$

$$\therefore H(X_n) \text{ exists if } \alpha > 2, \text{ and diverges to } \infty \text{ if } 1 < \alpha \leq 2.$$

Problem 3. Note that $P_{X_{01}|1}(\theta_1, x_1)$ is defined as $P_{\theta_1}(\theta_1 \wedge X_{01}[1] = x_1)$, while $P_{X_{02}|2}(x_2)$ is defined as $\Pr\{X_{02}[1] = x_2\}$.

Since $X_{01}[1]$ and θ_1 are independent, there is $P_{X_{01}|1}(\theta_1, x_1) = P_{\theta_1}(\theta_1) P_{X_{01}|1}(x_1)$.

Author: 许博翔

2

Information Theory HW3

许博翔

October 19, 2023

Note that in this homework, I'll use the following definition:

Problem 1, 2: If $P(p)$, then $P(x) = p(1-p)^{x-1}$.

Problem 3: If $P = G(p)$, then $P(x) = (1-p)^{p^{x-1}}$, which is the definition given in the homework.

$$\exp_2(x) := 2^x.$$

Problem 1.

$$\text{(a) Consider } \phi_{\gamma, \tau}(x) := \begin{cases} 1, & \text{if } LR(x) > \tau \\ \gamma, & \text{if } LR(x) = \tau \\ 0, & \text{if } LR(x) < \tau \end{cases}$$

$$LR(0) = \frac{P_1(0)}{P_1(0)} = \frac{1-p_1}{1-p_1}$$

$$LR(1) = \frac{P_1(1)}{P_1(1)} = \frac{p_1}{p_1}$$

$$\therefore p_0 < p_1.$$

$$\therefore LR(1) = \frac{p_1}{p_1} > 1 > \frac{1-p_1}{1-p_0} = LR(0).$$

By Neyman-Pearson theorem, $\phi_{\gamma, \tau}$ is optimal.

$$\pi_{10}(\phi_{\gamma, \tau}) = P_1(LR(X) > \tau) + \gamma P_1(LR(X) = \tau).$$

$$\pi_{01}(\phi_{\gamma, \tau}) = P_1(LR(X) < \tau) + (1-\gamma) P_1(LR(X) = \tau).$$

We only need to consider the cases $\tau = LR(x)$ for some x , since other cases can be reduced to these cases by setting γ properly.

$$\text{For } \tau = LR(0), \pi_{10} = P_1(1) + \gamma P_1(0) = p_0 + \gamma(1-p_1); \pi_{01} = 0 + (1-\gamma) P_1(0) = (1-\gamma)(1-p_1).$$

$$\text{For } \tau = LR(1), \pi_{10} = 0 + \gamma P_1(1) = \gamma p_0; \pi_{01} = P_1(0) + (1-\gamma) P_1(1) =$$

$$\begin{aligned} & \leq \sum_{s^* \in \mathcal{B}_\delta^{(n)}(S)} 2^{-n(H(S) - \delta)} \leq \sum_{s^* \in \mathcal{B}_\delta^{(n)}(S)} 2^{-n(H(S) - R - \delta)} \\ & = |\mathcal{B}_\delta^{(n)}(S)| 2^{-n(H(S) - R - \delta)} \leq 2^{n\delta} 2^{-n(H(S) - R - \delta)} \\ & \text{Since } H(S) - R - \delta > 0 \text{ by definition of } \delta, \text{ we get that} \\ & \lim_{n \rightarrow \infty} P_{S^n}^{(n)} = \lim_{n \rightarrow \infty} \Pr\{S^n \in \mathcal{A}_\delta^{(n)}(S) \cap \mathcal{B}^{(n)}\} \geq \lim_{n \rightarrow \infty} (1 - 2^{-n(H(S) - R - \delta)}) = 1. \end{aligned}$$

On the other hand, $P_{S^n}^{(n)} \leq 1$, so there is $\lim_{n \rightarrow \infty} P_{S^n}^{(n)} = 1$.

Information Theory HW2

$$\begin{aligned} \text{(a)} \because \forall i, n \in \mathbb{N}, P_{X_{01}|1}(X_{01}[i], X_{02}[i], \dots, X_{0n}[i]) &= P_{X_{01}|1}(X_{01}[i], X_{02}[i], \dots, X_{0n}[i]) \\ X_{01} \text{ is stationary w.r.t. } \theta &\Leftrightarrow \theta \text{ is } i \text{ i.i.d. } P_{X_{01}|1}(X_{01}[i], X_{02}[i], \dots, X_{0n}[i]) = P_{X_{01}|1}(1+i, X_{02}[i], \dots, X_{0n}[i]) = \dots \\ \therefore \{X_{01}[i]\} &\text{ is stationary.} \end{aligned}$$

$$\begin{aligned} & \text{By the definition of entropy rates,} \\ \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[\log P_{X_{01}|1}(X_{01}[1], X_{02}[1], \dots, X_{0n}[1])] &= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_{01}[1], X_{02}[1], \dots, X_{0n}[1]) = \mathcal{H}_X. \\ \Rightarrow \mathcal{H}(\{X_{01}[i]\}) &= \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[\log (X_{01}[1], X_{02}[1], \dots, X_{0n}[1])] \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[\log P_{X_{01}|1}(X_{01}[1], X_{02}[1], \dots, X_{0n}[1])] \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} (\Pr(\theta = 0) \mathbb{E}[\log \Pr(\theta = 0) P_{X_{01}|1}(X_{02}[1], \dots, X_{0n}[1]) \\ &+ \Pr(\theta = 1) \mathbb{E}[\log \Pr(\theta = 1) P_{X_{01}|1}(X_{02}[1], X_{03}[1], \dots, X_{0n}[1])]) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} ((1-q) \mathbb{E}[\log(1-q) + \log P_{X_{01}|1}(X_{02}[1], \dots, X_{0n}[1]) + q \mathbb{E}[\log q + \log P_{X_{01}|1}(X_{02}[1], X_{03}[1], \dots, X_{0n}[1])]) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} ((1-q) \log(1-q) + q \log q) + (1-q) \mathcal{H}_0 + q \mathcal{H}_1 = (1-q) \mathcal{H}_0 + q \mathcal{H}_1. \end{aligned}$$

(b) Suppose $\Theta_1 \sim \text{Ber}(q)$.

$$\text{Since } \{\Theta_i\} \text{ is stationary, } (1-q \quad q) = \begin{pmatrix} 1-q & \alpha \\ \beta & 1-\beta \end{pmatrix}.$$

$$\Rightarrow 1-q = (1-q)(1-\alpha) + q\beta$$

$$\Rightarrow \alpha(1-q) = q\beta$$

$$\Rightarrow q = \frac{\alpha}{\alpha + \beta}$$

$$\therefore P_{X_{01}|1}(X_{01}[1], X_{02}[1]) = \sum_{i=1}^\infty P_{X_{01}|1}(X_{01}[i], X_{02}[i], \dots, X_{0n}[i]) = P_{X_{01}|1}(1, X_{02}[1], X_{03}[1], \dots, X_{0n}[1])$$

$$= P_{\theta_1} P_{X_{01}|1} \prod_{i=2}^\infty P_{\theta_i, \theta_{i-1}} P_{X_{01}|1}(X_{02}[i], X_{03}[i], \dots, X_{0n}[i]) = P_{X_{01}|1} \prod_{i=2}^\infty P_{X_{01}|1}(X_{02}[i], X_{03}[i], \dots, X_{0n}[i])$$

$$X_{01}[i] \text{ is independent of } \frac{X_{01}[j]}{X_{01}[i]} \text{ for any } (j, j') \neq (i, i) \Rightarrow P_{\theta_1, \theta_{i-1}} P_{X_{01}|1}(X_{02}[i], X_{03}[i], \dots, X_{0n}[i]) = P_{\theta_1, \theta_{i-1}} P_{X_{01}|1}(X_{02}[i], X_{03}[i], \dots, X_{0n}[i])$$

$$\Rightarrow \forall i, n \in \mathbb{N}, P_{X_{01}|1}(X_{01}[1], X_{02}[1], \dots, X_{0n}[1]) = P_{\theta_1, \theta_{i-1}} P_{X_{01}|1}(X_{02}[1], X_{03}[1], \dots, X_{0n}[1])$$

$$P_{\theta_1, \theta_{i-1}} P_{X_{01}|1}(X_{02}[1], X_{03}[1], \dots, X_{0n}[1]) = P_{X_{01}|1}(1+i, X_{02}[1], \dots, X_{0n}[1]) = P_{X_{01}|1}(1+i, X_{02}[1], \dots, X_{0n}[1])$$

$$\therefore \{X_{01}[i]\} \text{ is stationary.}$$

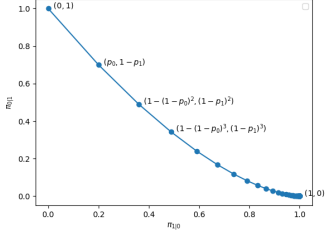
By theorem 11, $\mathcal{H}(\{X_{01}[i]\}) = H(X_{01}[2]|X_{01}[1])$

By Neyman-Pearson theorem, $\phi_{\gamma,\tau}$ is optimal.

We only need to consider the cases $\tau = LR(y)$ for some y , since other cases can be reduced to these cases by setting γ properly.

Since $LR(y)$ is decreasing, for $\tau = LR(y)$, $\pi_{i0}(\phi_{\gamma,\tau}) = P_0\{Y < y\} + \gamma P_0\{Y = y\} = 1 - (1 - p_0)^{i-1} + \gamma p_0(1 - p_0)^{i-1} = 1 - (1 - p_0)^{i-1}(1 - \gamma p_0)$.
 $\pi_{i0}(\phi_{\gamma,\tau}) = P_1\{Y > y\} + (1 - \gamma)P_1\{Y = y\} = (1 - p_1)^i + (1 - \gamma)p_1(1 - p_1)^{i-1} = (1 - \gamma p_1)(1 - p_1)^{i-1}$.

For each y , it forms a segment, where the intersection of the segments formed by y and $y + 1$ is $(1 - (1 - p_0)^i, (1 - p_1)^i)$, which can be calculated by setting γ in the segment formed by y to 1 or in the other segment to 0.



(c) Let Y_i be the random variable denoting the length of the sequence between the $(i-1)$ -th 1 and the i -th 1 (including the i -th 1 and excluding the $(i-1)$ -th 1). One can see that Y_i are i.i.d. and $Y_i \sim G(p_0)$.

Clearly, $Z = Y_1 + Y_2 + \dots + Y_n$ is the random variable of the length of the observed sequence.

Let $Q_0 = G(p_0)$, $Q_1 = G(p_1)$.

Author: 许博翔

3

$$\begin{aligned} \text{From Chernoff-Stein lemma, } \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{0,n}(n, \epsilon) &= \mathbb{E}_{Y \sim G(p_0)} \left[\log \frac{Q_0(Y)}{Q_1(Y)} \right] = \\ &= \sum_{i=1}^{\infty} p_0(1-p_0)^{i-1} \log \frac{p_0(1-p_0)^{i-1}}{p_1(1-p_1)^{i-1}} = \sum_{i=1}^{\infty} p_0(1-p_0)^{i-1} \log \frac{p_0}{p_1} + \sum_{i=1}^{\infty} (i-1)p_0(1-p_0)^{i-1} \log \frac{1-p_0}{1-p_1} = \\ &= p_0 \log \frac{p_0}{p_1} + p_0 \log \frac{1-p_0}{1-p_1} + \sum_{i=1}^{\infty} (i-1)p_0(1-p_0)^{i-1} \log \frac{1-p_0}{1-p_1} = \\ &= \log \frac{p_0}{p_1} + p_0 \log \frac{1-p_0}{1-p_1} + \sum_{i=1}^{\infty} \sum_{j=i+1}^{\infty} (1-p_0)^{i-1} \log \frac{p_0}{p_1} + p_0 \log \frac{1-p_0}{1-p_1} \sum_{j=1}^{\infty} (1-p_0)^j = \\ &= \log \frac{p_0}{p_1} + p_0 \log \left(\frac{1-p_0}{1-p_1} \right) \frac{1-p_0}{p_0} = \log \frac{p_0}{p_1} + \left(\frac{1}{p_0} - 1 \right) \log \frac{1-p_0}{1-p_1} \end{aligned}$$

Problem 2.

$$\begin{aligned} \text{(a)} \quad \pi_0^{(n)}(x^n) &= \Pr\{\Theta = 0 | X^n = x^n\} = \Pr\{X_i \stackrel{i.i.d.}{\sim} P_0 | X^n = x^n\} = \frac{\Pr\{X_i \stackrel{i.i.d.}{\sim} P_0 \wedge X^n = x^n\}}{\Pr\{X_i \stackrel{i.i.d.}{\sim} P_0 \wedge X^n = x^n\} \vee \{X_i \stackrel{i.i.d.}{\sim} P_1 \wedge X^n = x^n\}} = \frac{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(x_i)} \\ &= \frac{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(x_i)} \\ \text{Similarly, } \pi_1^{(n)}(x^n) &= \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(x_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(x_i)} \end{aligned}$$

$$\begin{aligned} \text{(b)} \quad & -\frac{1}{n} \log \pi_0^{(n)} \prod_{i=1}^n P_0(X_i) = -\frac{1}{n} \left(\log \pi_0^{(n)} + \sum_{i=1}^n \log P_0(X_i) \right) \xrightarrow{\text{SLLN}} -\frac{1}{n} \log \pi_0^{(n)} - \\ & \mathbb{E}[\log(P_0(X))] \xrightarrow{\text{law of large numbers}} -\mathbb{E}[\log(P_0(X))] = H(X) \text{ as } n \rightarrow \infty. \\ \text{From HW2 we know that } H(X) &\leq -\sum_{i=1}^n P_0(i) \log P_0(i), \text{ with equality} \iff \\ & P_1 \sim P_0. \\ & -\frac{1}{n} \log \pi_1^{(n)} \prod_{i=1}^n P_1(X_i) = -\frac{1}{n} \left(\log \pi_1^{(n)} + \sum_{i=1}^n \log P_1(X_i) \right) \xrightarrow{\text{SLLN}} -\frac{1}{n} \log \pi_1^{(n)} - \\ & \mathbb{E}[\log(P_1(X))] \xrightarrow{\text{law of large numbers}} -\mathbb{E}[\log(P_1(X))] > H(X) \text{ as } n \rightarrow \infty. \end{aligned}$$

$$\begin{aligned} & \frac{\pi_0^{(n)} \prod_{i=1}^n P_0(X_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(X_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(X_i)} \rightarrow \exp_2(n \mathbb{E}[\log(P_1(X))] + n H(X)) = \exp_2(\mathbb{E}[\log(P_1(X))]) + \\ & H(X) \text{ as } \mathbb{E}[\log(P_1(X))] > 0 \text{ as } n \rightarrow \infty. \\ \therefore \pi_0^{(n)}(X^n) &= \frac{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(x_i)} = \frac{1}{1 + \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(x_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(x_i)}} \rightarrow \frac{1}{1+0} = 1 \text{ as } n \rightarrow \infty. \end{aligned}$$

Author: 许博翔

4

Information Theory HW3

许博翔

November 2, 2023

Problem 1.

- (a) Since N_0 is deterministic from X_1, X_2, \dots, X_{N_0} , N_1 is deterministic from X_1, X_2, \dots, X_{N_0} , there is $I(N_0; X_1, \dots, X_{N_0}) = H(N_0) = \frac{1}{3} \log 3 + \frac{2}{3} \log 3 - 1 = \log 3 - \frac{2}{3}$, $I(N_1; X_1, \dots, X_{N_1}) = H(N_1) = \sum_{i=1}^n \frac{1}{2^i} = \sum_{i=1}^n \frac{1}{2^i} \sum_{j=i}^n \frac{1}{2^j} = \sum_{j=1}^n \sum_{i=j}^n \frac{1}{2^j} = \sum_{j=1}^n \frac{1}{2^{j-1}} = 2$.
- (b) Let's assume $n \geq 2$ (because for $n = 1$ there is nothing to be computed).
 Claim: X_1, X_2, \dots, X_{n-1} are mutually independent $\text{Ber}(\frac{1}{2})$.
- Proof.* $\forall x \in \{0, 1\}^{n-1}$, there is exactly one $x^* \in \{0, 1\}^n$ (which is $x_1, x_2, \dots, x_{n-1}, x_n \oplus \dots \oplus x_{n-1}$) s.t. $2 \cdot \sum_{i=1}^n x_i^* \equiv 1 \pmod{2}$ and $\forall 1 \leq i \leq n-1$, $x_i^* = x_i$.
 $\therefore \Pr\{(X_1, \dots, X_{n-1}) = x\} = \Pr\{(X_1, \dots, X_n) = x^*\} = 2^{-(n-1)}$.
 $\Rightarrow (X_1, \dots, X_{n-1})$ is a uniform distribution on $\{0, 1\}^{n-1}$, which means X_1, X_2, \dots, X_{n-1} are mutually independent $\text{Ber}(\frac{1}{2})$. ■

Similarly, for any distinct i_1, i_2, \dots, i_{n-1} , $X_{i_1}, \dots, X_{i_{n-1}}$ are mutually independent.

$$\begin{aligned} \text{Let } 1 \leq i \leq n-1, \\ I(X_i; X_{i+1}, \dots, X_{n-1}) &= H(X_i | X_1, \dots, X_{i-1}) - H(X_i | X_1, \dots, X_{i-1}, X_{i+1}) \\ X_1, \dots, X_i &\text{ are mutually independent} \\ H(X_i) &= H(X_i | X_1, \dots, X_{i-1}, X_{i+1}) \\ X_1, \dots, X_{i-1} &\text{ are mutually independent if } i < n-1 \\ \begin{cases} H(X_i) - H(X_i) = 0, & \text{if } i < n-1 \\ H(X_i) - H(X_i | X_1 \oplus \dots \oplus X_{i-1} \oplus X_{i+1} | X_1, \dots, X_{i-1}, X_{i+1}) = H(X_i) - 0 = 1, & \text{if } i = n-1 \end{cases} \end{aligned}$$

Author: 许博翔

6

Problem 3.

- (a) Let $X_i \in \mathcal{X}^{(i)}$.
 $I(X; Y) = H(X) - H(X|Y) \stackrel{I \text{ is deterministic from } Y}{=} H(X, I) - H(X|Y) = H(X|I) + H(I) - H(X|Y) \stackrel{I \text{ is deterministic from } Y}{=} H(X|I) + H(I) - H(X|Y, I) = I(X; Y|I) + H(I)$.
- (b) The capacity is $\max_P I(X; Y) = \max_P E_{(X,Y) \sim P_{X,Y}} (\log \frac{P_{1,Y}(Y|X)}{P_Y(Y)}) = \max_P \sum_{i=1}^I P_Y(i) (I(X_i; Y) - \log P_Y(i)) = \sum_{i=1}^I P_Y(i) (C^{(i)} + H(I))$.
- (c) Consider the distribution: $P_Y(i) = \frac{2^{C^{(i)}}}{\sum_{j=1}^I 2^{C^{(j)}}}$.
 $\sum_{i=1}^I P_Y(i) C^{(i)} + H(I) = \sum_{i=1}^I P_Y(i) \log \frac{2^{C^{(i)}}}{\sum_{j=1}^I 2^{C^{(j)}}} + \sum_{i=1}^I P_Y(i) \log \sum_{j=1}^I 2^{C^{(j)}} = \sum_{i=1}^I P_Y(i) \log \frac{P_Y(i)}{P_Y(i)} + \log \sum_{j=1}^I 2^{C^{(j)}} + \log \sum_{j=1}^I 2^{C^{(j)}} \geq \log \sum_{j=1}^I 2^{C^{(j)}}$, with equality $\iff D(P_Y || P_Y) = 0 \iff P_Y = P_Y$.
 \therefore the capacity is $\log \sum_{j=1}^I 2^{C^{(j)}}$, and the distribution P_Y is $P_Y(i) = \frac{2^{C^{(i)}}}{\sum_{j=1}^I 2^{C^{(j)}}}$.

Problem 4.

- (a) Suppose that $X \sim \text{Ber}(q)$.
 $\Rightarrow P_X(0) = 1 - q + pq = 1 - \frac{1}{2}q$, $P_X(1) = q(1 - \frac{1}{2}q)$
 $I(X; Y) = H(X) + H(Y) - H(X, Y) = -q \log q - (1 - q) \log(1 - q) - \frac{1}{2}q \log(\frac{1}{2}q) - (1 - \frac{1}{2}q) \log(1 - \frac{1}{2}q) + (1 - q) \log(1 - q) + 2 \cdot \frac{1}{2}q \log(\frac{1}{2}q) = -q \log q - \frac{1}{2}q \log(\frac{1}{2}q) - (1 - \frac{1}{2}q) \log(1 - \frac{1}{2}q) - (1 - \frac{1}{2}q) \log(1 - \frac{1}{2}q)$.
 Let $\frac{dI(X; Y)}{dq} = -1 - \frac{1}{2} \log(\frac{1}{2}q) - \frac{1}{2} \log e + \frac{1}{2} \log(1 - \frac{1}{2}q) + \frac{1}{2} \log e = -1 + \frac{1}{2} \log \frac{1 - \frac{1}{2}q}{\frac{1}{2}q} = 0$.
 $\Rightarrow \log \frac{1 - \frac{1}{2}q}{\frac{1}{2}q} = 2$.
 $\Rightarrow \frac{1 - \frac{1}{2}q}{\frac{1}{2}q} = 4$.
 $\Rightarrow q = \frac{2}{5}$.

Author: 许博翔

3

$n \rightarrow \infty$.

$$\begin{aligned} \text{As what we computed above, for any constant } c > 0, & -\frac{1}{n} \log \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(X_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(X_i)} \\ &= H(X) + \mathbb{E}[\log(P_1(X))] + \frac{1}{n} \log c + \text{a constant} = H(X) + \mathbb{E}[\log(P_1(X))] = D(P_1 || P_0). \\ \therefore \log \text{ is an increasing function, and } & \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(X_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(X_i)} < \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(X_i)}{\pi_0^{(n)} \prod_{i=1}^n P_0(X_i) + \pi_1^{(n)} \prod_{i=1}^n P_1(X_i)} \\ &= \pi_1^{(n)}(X^n) < \frac{\pi_1^{(n)} \prod_{i=1}^n P_1(X_i)}{2 \pi_0^{(n)} \prod_{i=1}^n P_0(X_i)} \text{ when } n \rightarrow \infty. \\ \therefore \text{ by squeeze theorem, } & -\frac{1}{n} \log \pi_1^{(n)}(X^n) \rightarrow D(P_1 || P_0) \text{ as } n \rightarrow \infty. \end{aligned}$$

Problem 3.

- (a) Let $X \sim P$.
 $D(P||G(p)) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)} = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{(1-p)^{p^{x-1}}} = H(X) - \mathbb{E}[\log((1-p)^{p^{x-1}})] = H(X) - \log(1-p) - \mathbb{E}[\log(p)] = H(X) - \log(1-p) - \log(p) \mathbb{E}[X-1] = H(X) - \log(1-p) + \log p - p \log p$.
 $\frac{d}{dp} D(P||G(p)) = \frac{1}{1-p} + \frac{1}{p} - \frac{1}{p^2} = \frac{1 - (1-p)p}{p^2(1-p)}$, which equals to 0 $\iff \frac{1}{1-p} = \mu \iff p = 1 - \frac{1}{\mu}$.
 One can also verify that if $p < 1 - \frac{1}{\mu}$, $\frac{d}{dp} D(P||G(p)) < 0$ and if $p > 1 - \frac{1}{\mu}$, $\frac{d}{dp} D(P||G(p)) > 0$.
 \therefore the minimum possible value of $D(P||G(p))$ occurs when $p = 1 - \frac{1}{\mu}$, that is, the distribution is $G(1 - \frac{1}{\mu})$, and $D(P||G(p)) = H(X) - \log \mu + (1 - \mu) \log(1 - \mu)$.

- (b) Let $X_i \sim P, Y \sim R$ where $R(y) := \frac{1}{\sum_{j=1}^m R(j)} P(y)$.
 From HW2 we know that $H(R) \leq -\sum_{j=1}^m R(j) \log Q(j)$, with equality $\iff Q \sim R$.
 $\Rightarrow \sum_{i=1}^m D(P_i || Q) = \sum_{i=1}^m \left(H(X_i) - \sum_{j=1}^m R(j) \log Q(j) \right)$

Author: 许博翔

5

Problem 2.

- (a) $I(X_1; X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1)$
 $I(X_1; X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1, X_2) - I(X_1; X_2 | X_1)$
 $\Rightarrow I(X_1; X_1) + I(X_1; X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1, X_2) - I(X_1; X_2 | X_1) + I(X_1; X_2) = I(X_1; X_2) + I(X_1; X_2) - I(X_1; X_2 | X_1, X_2) \leq I(X_1; X_2) + I(X_1; X_2)$
- (b) It's equivalent to two Markov's chains: $X_1 - X_2 - X_3, X_1 - X_2 - X_4$.
 $I(X_1; X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1)$
 $I(X_1; X_1) = I(X_1; X_2) - I(X_1; X_2 | X_1)$
 $I(X_1; X_2) + I(X_1; X_2) \geq I(X_1; X_2) + I(X_1; X_1) - I(X_1; X_1 | X_2) \geq I(X_1; X_2) + I(X_1; X_1) - I(X_1; X_2) = I(X_1; X_1)$

Author: 许博翔

2

Information Theory HW5

许博翔

November 23, 2023

Problem 1.

- (a) (1) From Gaussian integral, we know that $\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$.
 $\int_{-\infty}^{\infty} x e^{-x^2} dx = \int_{-\infty}^{\infty} \frac{1}{2} e^{-x^2} d(x^2) = -\frac{1}{2} e^{-x^2} + c$.
 $\lim_{x \rightarrow \infty} x e^{-x^2} = \lim_{x \rightarrow \infty} \frac{x}{e^{x^2}} = \lim_{x \rightarrow \infty} \frac{1}{2x e^{x^2}} = 0$.
 $\lim_{x \rightarrow -\infty} x e^{-x^2} = \lim_{x \rightarrow -\infty} \frac{x}{e^{x^2}} = \lim_{x \rightarrow -\infty} \frac{1}{2x e^{x^2}} = 0$.
 $\int_{-\infty}^{\infty} x^2 e^{-x^2} dx = \int_{-\infty}^{\infty} x e^{-x^2} \cdot x dx = \frac{1}{2} e^{-x^2} \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} -\frac{1}{2} e^{-x^2} \cdot 1 dx = 0 + \frac{1}{2} \sqrt{\pi} = \frac{1}{2} \sqrt{\pi}$.
 $f(x) = \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2}, g(x) = \frac{1}{\sigma_2 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_2}{\sigma_2} \right)^2}$.
 $D(f||g) = \int_{-\infty}^{\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2} \log \left(\frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2} \cdot \frac{1}{\sigma_2 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_2}{\sigma_2} \right)^2} \right) dx$
 $= \int_{-\infty}^{\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2} \left(\log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{1}{2} \log e \left(-\left(\frac{x - \mu_1}{\sigma_1} \right)^2 + \left(\frac{x - \mu_2}{\sigma_2} \right)^2 \right) \right) dx$
 $= \int_{-\infty}^{\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2} \left(\log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{1}{2} \log e \left(\frac{x - \mu_1}{\sigma_1} \right)^2 \right) dx$
 $+ \frac{1}{2} \log e \left(\frac{x - \mu_2}{\sigma_2} \right)^2 + \log e \left(\frac{\mu_1 - \mu_2}{\sigma_1} \right) \left(\frac{x - \mu_1}{\sigma_1} \right) + \frac{1}{2} \log e \left(\frac{\mu_1 - \mu_2}{\sigma_2} \right)^2 dx$
 $= \log \left(\frac{\sigma_2}{\sigma_1} \right) - \frac{1}{2} \log e + \frac{1}{2} \log e \left(\frac{\sigma_1^2}{\sigma_2^2} + (\mu_1 - \mu_2)^2 \right)$
 $= \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\log e}{2\sigma_1^2} (\sigma_1^2 - \sigma_2^2 + (\mu_1 - \mu_2)^2)$.
- (2) $f(x) = \frac{1}{\sqrt{2\pi} \sigma_1} e^{-\frac{1}{2} \left(\frac{x - \mu_1}{\sigma_1} \right)^2}, g(x) = \frac{1}{\sqrt{2\pi} \sigma_2} e^{-\frac{1}{2} \left(\frac{x - \mu_2}{\sigma_2} \right)^2}$.
 $\int_{-\infty}^{\infty} x e^{-x^2} dx = e^{-x^2} - \int_{-\infty}^{\infty} e^{-x^2} dx = -(x+1)e^{-x^2} + c$.
 $\int_{-\infty}^{\infty} x e^{-x^2} dx = -e^{-x^2} x - \int_{-\infty}^{\infty} -e^{-x^2} dx = -(x+1)e^{-x^2} + c$.
 $\lim_{x \rightarrow \infty} e^{-x^2} x = 0$.

Author: 许博翔

4

$$\begin{aligned}
\int_{-\infty}^{\infty} |x-a|e^{-|x-b|}dx &= \int_{-\infty}^{\infty} |x+b-a|e^{-|x|}dx, \\
\text{If } e := a-b \geq 0, \text{ then } \int_{-\infty}^{\infty} |x+b-a|e^{-|x|}dx &= \int_{-\infty}^e (e-x)e^x dx + \int_0^{\infty} (e+x)e^{-x} dx \\
&= e + 1 + (-e^{-e} + e) + ((e+1)e^{-e} - 1) + (e+1)e^{-e} - e^{-e} = 2e + 2e^{-e}, \\
\text{If } e < 0, \text{ then } \int_{-\infty}^{\infty} |x-c|e^{-|x|}dx &= \int_{-\infty}^{\infty} |x+c|e^{-|x|}dx = -2e + 2e^e, \\
\therefore \int_{-\infty}^{\infty} |x-a|e^{-|x-b|}dx &= 2|a-b| + 2e^{-|a-b|}, \\
D(f|g) &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\sigma_1}} e^{-\frac{x^2}{2\sigma_1^2}} \log \left(\frac{1}{\sqrt{2\sigma_1}} e^{-\frac{x^2}{2\sigma_1^2}} / \frac{1}{\sqrt{2\sigma_2}} e^{-\frac{x^2}{2\sigma_2^2}} \right) dx \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\sigma_1}} e^{-\frac{x^2}{2\sigma_1^2}} \left(\log \left(\frac{\sigma_2}{\sigma_1} \right) + \sqrt{2} \log e \left(\frac{|x-\mu_1|}{\sigma_1} - \frac{|x-\mu_2|}{\sigma_2} \right) \right) dx \\
&= \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\log e}{\sigma_1 \sigma_2} \left(\frac{\sigma_1^2}{2} (2 \frac{\sqrt{2}}{\sigma_1} |\mu_1 - \mu_2| + 2e^{-\frac{\sigma_2^2}{2\sigma_1^2}(\mu_1 - \mu_2)} - \frac{\log e}{\sigma_1^2} \frac{\sigma_1^2}{2}) \right) \\
&= \log \left(\frac{\sigma_2}{\sigma_1} \right) + \frac{\sigma_1 \log e}{\sigma_2} \sqrt{\frac{2}{\sigma_1}} |\mu_1 - \mu_2| + e^{-\frac{\sigma_2^2}{2\sigma_1^2}(\mu_1 - \mu_2)} - \log e.
\end{aligned}$$

(b) The first KL divergence – the second KL divergence = $\frac{\log e}{2\sigma_2^2}(\sigma_1^2 - \sigma_2^2) - \frac{\sigma_1 \log e}{\sigma_2}$

$$\begin{aligned}
\log e &= \frac{\log e}{2} \left(\left(\frac{\sigma_1}{\sigma_2} \right)^2 - 2 \frac{\sigma_1}{\sigma_2} + 1 \right) = \frac{\log e}{2} \left(\frac{\sigma_1}{\sigma_2} - 1 \right)^2 \geq 0, \\
\therefore \text{the first KL divergence} &\geq \text{the second KL divergence, the equation holds} \\
&\iff \sigma_1 = \sigma_2.
\end{aligned}$$

(c) Let $x := |\mu_1 - \mu_2|$.

$$\begin{aligned}
\text{The first KL divergence} &- \text{the second KL divergence} = \frac{\log e}{2} (\mu_1 - \mu_2)^2 - \\
&\log e \left(\frac{\sqrt{2}}{\sigma_1} |\mu_1 - \mu_2| + e^{-\frac{\sigma_2^2}{2\sigma_1^2}(\mu_1 - \mu_2)} \right) + \log e \\
&= \frac{\log e}{2} x^2 - \log e \left(\frac{\sqrt{2}}{\sigma_1} x + e^{-\frac{\sigma_2^2}{2\sigma_1^2}x} \right) + \log e \\
&= \log e \left(\frac{1}{2} x^2 - \frac{\sqrt{2}}{\sigma_1} x - e^{-\frac{\sigma_2^2}{2\sigma_1^2}x} + 1 \right), \\
\therefore \text{the first KL divergence is the larger} &\iff \frac{1}{2} x^2 - \frac{\sqrt{2}}{\sigma_1} x - e^{-\frac{\sigma_2^2}{2\sigma_1^2}x} + 1 \geq 0.
\end{aligned}$$

Problem 2.

$$\begin{aligned}
\text{(a) } h(X) &= E_{\mathbf{q} \sim J_X} \left(\log \frac{1}{\int_X \mathbf{q}(x)} \right) = \int_{-\infty}^{\infty} \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} (\log(2b) + \log e \left(\frac{|x-\mu|}{b} \right)) dx = \log(2b) + \\
&\log e \int_{-\infty}^{\infty} \frac{1}{b} e^{-\frac{|x-\mu|}{b}} \frac{2-x-\mu}{b} dx = \log(2b) + \log e = \log(2be).
\end{aligned}$$

Author: 许得刚

2

$$\begin{aligned}
\text{(b) From Problem 1 (a)(2), we know that } \int_{-\infty}^{\infty} |x-a|e^{-|x-b|}dx &= 2|a-b| + 2e^{-|a-b|}, \\
E(|X|) &= \int_{-\infty}^{\infty} |x| \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} dx = \frac{1}{2b} b^2 (2|\mu| + 2e^{-b|\mu|}) = b(|\mu| + e^{-b|\mu|}). \\
\text{Let } g(y) &:= y + e^{-y}, \\
&\Rightarrow g'(y) = 1 - e^{-y} > 0 \text{ when } y > 0, \\
&\Rightarrow g(y) \text{ is strictly increasing on } (0, \infty), \\
&\Rightarrow b(|\mu| + e^{-b|\mu|}) \stackrel{(1)}{\geq} bg(0) = 2b, \\
&\Rightarrow 2b \leq E(|X|) \leq B, \\
&\Rightarrow b \leq \frac{B}{2}, \\
&\Rightarrow h(X) = \log(2be) \leq \log Be, \text{ and when the equation holds, the distribution of} \\
&X \text{ is } \text{Lap}(0, \frac{B}{2}) \text{ since the equation in (1) holds} \iff \mu = 0, \text{ and the equation} \\
&\text{in (2) holds.}
\end{aligned}$$

Problem 3.

(a) Consider $\tilde{b}(x) := E[b(x, Y)] = E_{P_{Y|X}}[b(x, Y)]$.

Since $\tilde{b}(x) = \sum_{y \in \mathcal{Y}} P_{Y|X}(y|x) b(x, y)$ is a deterministic function of x , $\tilde{b}(x)$ is an input-only cost function.

$$\therefore \frac{1}{n} \sum_{i=1}^n E_{Y_i}[b(x_i, Y_i)] = \frac{1}{n} \sum_{i=1}^n \tilde{b}(x_i).$$

$$\therefore \text{the cost constraint becomes: } \frac{1}{n} \sum_{i=1}^n \tilde{b}(x_i) \leq B.$$

Therefore, this problem is equivalent to the channel coding problem with input-cost only function \tilde{b} .

$$\begin{aligned}
\text{From Theorem 1 in Lecture 5, } C(B) &= \max_{P_X, B_{P_X}[\tilde{b}(X)] \leq B} I(X; Y) \\
&\stackrel{(1)}{=} \max_{P_X, B_{P_X}[\tilde{b}(X)] \leq B} \max_{P_{Y|X}[\tilde{b}(X, Y)] \leq B} I(X; Y).
\end{aligned}$$

$$\text{(b) First, } P_{Y|X}(y|x) = P_X(y-x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y-x)^2}.$$

$$\text{Let } b(x, y) := y^2.$$

$$\text{The cost constraint is } \frac{1}{n} \sum_{i=1}^n E_{Y_i}[b(x_i, Y_i)] = \frac{1}{n} \sum_{i=1}^n E_{Y_i}[Y_i^2] \leq B.$$

From the formula in Problem 1 (a)(1):

$$\begin{aligned}
\tilde{b}(x) &:= E[b(x, Y)] = \int_{-\infty}^{\infty} P_{Y|X}(y|x) b(x, y) dy = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y-x)^2} y^2 dy \\
&= \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}y^2} ((y-x)^2 + 2(y-x)x + x^2) dy
\end{aligned}$$

Author: 许得刚

3

$$\begin{aligned}
&= \sigma^2 + 0 + x^2 = \sigma^2 + x^2, \\
&\Rightarrow \text{the cost constraint becomes } \frac{1}{n} \sum_{i=1}^n (\sigma^2 + x_i^2) = \frac{1}{n} \sum_{i=1}^n \tilde{b}(x_i) \leq B, \text{ which is} \\
&\frac{1}{n} \sum_{i=1}^n |x_i|^2 \leq B - \sigma^2. \\
\text{From the example of Gaussian channel capacity in Lecture 5, we get that } C(B) &= \\
&\frac{1}{2} \log \left(1 + \frac{B - \sigma^2}{\sigma^2} \right) = \frac{1}{2} \log \left(\frac{B}{\sigma^2} \right).
\end{aligned}$$

Problem 4. In HW2, we know that if $\sum_i p_i = \sum_i q_i = 1$ where $p_i, q_i \geq 0$, then

$$\sum_i p_i \log \frac{1}{p_i} \leq \sum_i p_i \log \frac{1}{q_i} - (1)$$

$$\text{(a) } D_{\min} = \min_{\mathbf{q} \in \mathcal{Q}} E[d(S, \mathbf{q}(S))] = \min_{\mathbf{q} \in \mathcal{Q}} E \left[\log \frac{1}{\mathbf{q}(S)} \right] = 0 \text{ if } \mathbf{q}(s) = \mathbb{I}(S=s).$$

$$D_{\max} = \max_{\mathbf{q}} E[d(S, \mathbf{q})] = \min_{\mathbf{q}} E \left[\log \frac{1}{\mathbf{q}(S)} \right].$$

$$\therefore E \left[\log \frac{1}{\mathbf{q}(S)} \right] = \sum P_S(s) \log \frac{1}{\mathbf{q}(s)} \stackrel{(1)}{\geq} \sum P_S(s) \log \frac{1}{P_S(s)} = H(S) = H(\pi), \text{ and the equation holds when } \mathbf{q}(s) = P_S(s).$$

$$\therefore D_{\max} = H(\pi).$$

$$\text{(b) } H(S|\mathbf{Q}) = E_{(S, \mathbf{Q}) \sim P} \left[\log \frac{1}{P_{S|\mathbf{Q}}} \right] = \sum P_{\mathbf{Q}}(\mathbf{q}) \sum_s P_{S|\mathbf{Q}}(s|\mathbf{q}) \log \frac{1}{P_{S|\mathbf{Q}}(s|\mathbf{q})}$$

$$\stackrel{(1)}{\leq} \sum_{\mathbf{q}} P_{\mathbf{Q}}(\mathbf{q}) \sum_s P_{S|\mathbf{Q}}(s|\mathbf{q}) \log \frac{1}{\mathbf{q}(s)} = E_{(S, \mathbf{Q}) \sim P} \left[\log \frac{1}{\mathbf{q}(S)} \right].$$

$$\begin{aligned}
\text{(c) } R(D) &= \inf_{(S, \mathbf{Q})} \left\{ I(S; \mathbf{Q}) \mid E \left[\log \frac{1}{\mathbf{Q}(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
&= \inf_{(S, \mathbf{Q})} \left\{ I(S; \mathbf{Q}) \mid H(S|\mathbf{Q}) \leq E \left[\log \frac{1}{\mathbf{Q}(S)} \right] \leq D \text{ and } S \sim \pi \right\} \\
&\stackrel{(1)}{\leq} \inf_{(S, \mathbf{Q})} \left\{ I(S; \mathbf{Q}) \mid H(S|\mathbf{Q}) \leq D \text{ and } S \sim \pi \right\} \\
&\stackrel{(2)}{\leq} \inf_{(S, \mathbf{Q})} \left\{ I(S; \mathbf{Q}) \mid H(S|\mathbf{Q}) \leq D \text{ and } S \sim \pi \text{ and } \mathbf{Q}(s) = 1 \text{ for some } s \in S \right\} \\
&= \min_{(S, \mathbf{Q})} \left\{ I(S; \mathbf{Q}) \mid H(S|\mathbf{Q}) \leq D \text{ and } S \sim \pi \right\}.
\end{aligned}$$

(d) Let $\mathbf{q}_*(s) := \mathbb{I}(s=s)$.

Consider the distribution $\mathbf{Q} = \mathbf{q}_*$:

The equation in (2) holds \iff the equation in (1) holds $\iff \forall s, \mathbf{q}_* P_{S|\mathbf{Q}}(s|\mathbf{q}) = \mathbf{q}_*(s)$, which is true because \mathbf{q}_* with nonzero probability, $\mathbf{q} = \mathbf{q}_*$ for some s , and $\mathbf{q}_*(s) = \mathbb{I}(s=s) \stackrel{\mathbf{q}_*}{=} P_{S|\mathbf{Q}}(s|\mathbf{q}_*)$.

Author: 许得刚

4

$$\begin{aligned}
&\text{The equation in (3) holds since } \mathbf{q}_* = 1 \text{ for } s \in S. \\
&\therefore \text{with this distribution, } R(D) = \min_{(S, \hat{S})} \left\{ I(S; \hat{S}) \mid H(S|\hat{S}) \leq D \text{ and } S \sim \pi \right\} \\
&= \min_{(S, \hat{S})} \left\{ H(S) - H(S|\hat{S}) \mid H(S|\hat{S}) \leq D \text{ and } S \sim \pi \right\} \\
&= \min_{(S, \hat{S})} \left\{ H(S) - H(S|\hat{S}) \mid H(S|\hat{S}) \leq D \text{ and } S \sim \pi \right\} \\
&= H(\pi) - D^{\arg \min_{D \leq H(\pi) \leq D^{\max}} \max(0, H(\pi) - D)}.
\end{aligned}$$

Author: 许得刚

5