Sample Problems for Quiz 9

For the preparation of Quiz 9, you can focus on Problems 2 and 3.

1. (a) Suppose we intend to transmit a sequence of information bits through 4-ary ASK modulation taking values from $\{-3, -1, +1, +3\}$. In order to reduce the *word error* rate (WER), we design a channel code as follows.

info bits	binary codeword	channel codeword
$U_1U_2U_3$	$V_1V_2V_3$	(X_1, X_2, X_3)
000	000000	(-3, -3, -3)
001	001001	(-3, +3, -1)
010	010010	(-1, -3, +3)
011	011011	(-1, +3, +1)
100	100100	(+3, -1, -3)
101	101101	(+3, +1, -1)
110	110110	(+1, -1, +3)
111	111111	(+1, +1, +1)

What is the code rate measured in information bits per channel usage?

- (b) If the transmitter sends (-3, -3, -3), (-1, -3, +3), (+3, -1, -3), (+1, -1, +3), (-1, -3, +3), (-3, -3, -3), (+1, -1, +3), (+3, +1, -1), (-3, +3, -1), (-3, +3, -1), and the receiver, after performing detection based on thirty 4-ary receptions, outputs (+3, -1, -3), (-1, -3, +3), (+3, -1, -3), (+1, -1, +3), (-1, -3, +3), (-3, -3, -3), (+1, +1, +1), (+3, +1, -1), (-3, +3, -1), (-3, -3, -3), what is the word error rate (WER)? What is the (information) bit error rate (BER)?
- (c) Suppose the information sequence to be transmitted $\dots U_9U_8U_7 \ U_6U_5U_4 \ U_3U_2U_1$ is independent and satisfies $P_{U_i}(0) = P_{U_i}(1) = \frac{1}{2}$. What is the corresponding empirical P_X from the transmission of $\dots X_9X_8X_7 \ X_6X_5X_4 \ X_3X_2X_1$, where $X_{3k+3}X_{3k+2}X_{3k+1}$ is the channel codeword due to information block $U_{3k+3}U_{3k+2}U_{3k+1}$?
- (d) After performing measurement, the system designer found that the 4-ary ASK transmissions can be (approximately) modelled by a discrete memoryless channel (DMC), and the channel transition probability follows

$$\mathbb{Q} = \begin{bmatrix}
P_{Y|X}(-3|-3) & P_{Y|X}(-1|-3) & P_{Y|X}(1|-3) & P_{Y|X}(3|-3) \\
P_{Y|X}(-3|-1) & P_{Y|X}(-1|-1) & P_{Y|X}(1|-1) & P_{Y|X}(3|-1) \\
P_{Y|X}(-3|1) & P_{Y|X}(-1|1) & P_{Y|X}(1|1) & P_{Y|X}(3|1) \\
P_{Y|X}(-3|3) & P_{Y|X}(-1|3) & P_{Y|X}(1|3) & P_{Y|X}(3|3)
\end{bmatrix}$$

$$= \begin{bmatrix}
1 - 3\epsilon & \epsilon & \epsilon & \epsilon \\
\epsilon & 1 - 3\epsilon & \epsilon & \epsilon \\
\epsilon & \epsilon & 1 - 3\epsilon & \epsilon \\
\epsilon & \epsilon & \epsilon & 1 - 3\epsilon
\end{bmatrix}$$

Since it is a symmetric channel,¹ the optimal input distribution that achieves the channel capacity is uniform, i.e.,

$$P_X(-3) = P_X(-1) = P_X(1) = P_X(3) = \frac{1}{4}.$$

¹A DMC is said to be symmetric if the rows of the transition probability matrix \mathbb{Q} are permutations of each other and the columns of \mathbb{Q} are permutations of each other. It is known that the channel capacity of a symmetric DMC can be achieved by uniform input distribution.

Find the channel capacity.

Hint: Derive I(X; Y) using the uniform input distribution.

- (e) Is the rate in (a) a *reliable* transmission rate, give that ϵ in (d) is 0.1?
- (f) We redesign the transmission scheme as follows.

info bits	binary codeword	channel codeword
$U_1 U_2 U_3$	$V_1V_2V_3$	$(X_1, X_2, X_3, X_4, X_5, X_6)$
000	000000000000	(-3, -3, -3, -3, -3, -3)
001	001001001001	(-3, +3, -1, -3, +3, -1)
010	010010010010	(-1, -3, +3, -1, -3, +3)
011	011011011011	(-1, +3, +1, -1, +3, +1)
100	100100100100	(+3, -1, -3, +3, -1, -3)
101	101101101101	(+3, +1, -1, +3, +1, -1)
110	110110110110	(+1, -1, +3, +1, -1, +3)
111	1111111111111	(+1, +1, +1, +1, +1, +1)

Re-do (e).

Solution.

(a) We transmit three information bits using three channel usages. Hence, the code rate is $R = \frac{3}{3} = 1$ (information) bits per channel usages.

Note: What Shannon's capacity formula concerns is the *information bit per channel* usage, not code bit per channel usage, which is 2 in this example.

(b) There are ten channel codewords transmitted, among which three codewords are incorrectly detected. Hence, WER = $\frac{3}{10} = 0.3$.

For the three incorrectly detected codewords, the information sequences transmitted are 000, 110 and 001 but the detector at the receiver outputs 100, 111 and 000. Hence, among the thirty information bits transmitted, only three information bits are erroneously detected, yielding the (information) bit error rate (BER) $\frac{3}{30} = 0.1$.

Note: For the transmission of the first codeword (-3, -3, -3), the receiver may receive (+3, -3, -3) through a symbol-by-symbol decision maker, which is not a channel codeword. According a certain decision rule, this received word of three 4-ary symbols is regarded to be most likely the codeword (+3, -1, -3). The corresponding estimate of the information bits is thus 100. We therefore have 1 information bit error and 1 codeword error from the transmission of the first codeword.

(c) Since each of the eight channel codewords will appear 1/8 of the time, we have

$$P_X(-3) = P_X(-1) = P_X(1) = P_X(3) = \frac{6}{24} = \frac{1}{4}$$

(d) With
$$\mathcal{X} = \mathcal{Y} = \{-3, -1, 1, 3\}$$
, we derive

$$C = I(X; Y)$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_X(x) P_{Y|X}(y|x) \log_2 \frac{P_{Y|X}(y|x)}{\sum_{x' \in \mathcal{X}} P_X(x') P_{Y|X}(y|x')}$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{1}{4} P_{Y|X}(y|x) \log_2 \frac{P_{Y|X}(y|x)}{\sum_{x' \in \mathcal{X}} \frac{1}{4} P_{Y|X}(y|x')}$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{1}{4} P_{Y|X}(y|x) \log_2 \frac{P_{Y|X}(y|x)}{\frac{1}{4}} \quad \left(\sum_{x' \in \mathcal{X}} P_{Y|X}(y|x') = (1 - 3\epsilon) + \epsilon + \epsilon + \epsilon = 1\right)$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{1}{4} P_{Y|X}(y|x) \cdot 2 + \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{1}{4} P_{Y|X}(y|x) \log_2 P_{Y|X}(y|x)$$

$$= \frac{1}{2} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{Y|X}(y|x) + \frac{1}{4} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{Y|X}(y|x) \log_2 P_{Y|X}(y|x)$$

$$= 2 + (1 - 3\epsilon) \log_2(1 - 3\epsilon) + 3\epsilon \log_2 \epsilon.$$

(e) When $\epsilon = 0.1$, we have

$$C = 2 + 0.7 \log_2(0.7) + 0.3 \log_2(0.1) \approx 0.6432 \dots$$
 info bits/channel usage

Since R in (a) is also measured by the unit of info bits/channel usage, and since R = 1 > 0.6432, the rate in (a) is not a reliable transmission rate.

Note: Because the rate is not a reliable transmission rate, WER cannot be made arbitrarily small by extending the length of channel codeword (such as transmitting n info bits using n channel usages for n very large).

(f) The rate of the new design is $\frac{3}{6} = 0.5$ info bits/channel usage, which is a reliable transmission rate because it is less than 0.6432 info bits/channel usage.

Note: Even if 0.5 info bits/channel usage is less than C = 0.6432 info bits/channel usage, Shannon only proves the existence of a code for any desired WER, but did not tell us how such code can be designed.

The binary information sequence may sometimes be the output of a lossless data compression coding scheme such as Huffman code or Lempel-Ziv code. Shannon's source coding theorem then gives the minimum source coding rate attainable by a lossless data compression scheme is

$$R = \frac{H(S^n)}{n} = H(S) \text{ info bits/source symbol.}$$

For example, the source entropy in Slide IDC 6-29 is given by

 $H(S) \approx 2.38$ info bits/source symbol (See Problem 2(c)),

and we can design a Huffman code with average codeword length (or average source coding rate) R arbitrary close to 2.38 info bits/source symbol. Whether this source can be *reliably* transmitted over the channel in (d) should be answered according to

$$\frac{R = 2.38 \text{ info bits/source symbols}}{T_s \text{ seconds/source symbol}} \gtrsim \frac{C = 0.6432 \text{ info bits/channel usage}}{T_c \text{ seconds/channel usage}}$$

where T_s is the average time between two source symbols, while the channel transmits one channel symbol in every T_c seconds.

- 2. Sub-problem (a) is for your reference. Not a part of the exam. You can see how Huffman proves optimality of his code. However, sub-problems (b) and (c) will be a part of Quiz 9 and final exam.
 - (a) Prove that Huffman code has the minimum average codeword length among all uniquely decodable codes.

Hint: From Slide IDC 6-28, Huffman proposes to combine the two least probable source symbols into a new single symbol, whose probability is equal to the sum of the probabilities of the original two (See the example on Slide IDC 6-29). This is called the *reduced* source \mathcal{X}' , which contains one less symbol of the source \mathcal{X} before combining.

Hint: Prove by contradiction that if a code \mathcal{C}' for the reduced source \mathcal{X}' is optimal, then the code \mathcal{C} for its immediate extended source \mathcal{X} is optimal.

- (b) Find a Huffman code for the source with probabilities {0.25, 0.25, 0.25, 0.1, 0.1, 0.05} and compute their average codeword length.
- (c) Why doesn't the average codeword length of the optimal Huffman code in (b) equal the source entropy?

Solution.

(a) Let the size of the source \mathcal{X} be M. Assume without loss of generality that

$$p_1 \ge p_2 \ge \cdots \ge p_M,$$

where $p_i = \Pr[X = a_i]$. Suppose \mathcal{C}' is optimal (in the sense of having the minimum average codeword length) for the reduced source \mathcal{X}' but the code \mathcal{C} for the extended source \mathcal{X} that must satisfy

$$ACL(\mathcal{C}) = ACL(\mathcal{C}') + p_{M-1} + p_M \tag{1}$$

is not optimal, where ACL is a shorthand for average codeword length. Then, there must exist an optimal code \mathcal{D} with ACL(\mathcal{D}) < ACL(\mathcal{C}) and with a_{M-1} and a_M as siblings (See Property 3 in Slide IDC 6-27).

As suggested by Huffman, combine a_{M-1} and a_M into one new symbol $a_{M-1,M}$ and assign its probability as $p_{M-1} + p_M$, which immediately gives a code \mathcal{D}' for the reduced source \mathcal{X}' , whose average codeword length satisfies

$$ACL(\mathcal{D}) = ACL(\mathcal{D}') + p_{M-1} + p_M.$$
(2)

The two equations (1) and (2) indicate $ACL(\mathcal{D}') < ACL(\mathcal{C}')$; a contradiction to the optimality of \mathcal{C}' for the reduced source \mathcal{X}' is obtained.

(b) Slide IDC 6-30 gives a Huffman code as 00,01,10,110,1110,1111, of which the ACL is

$$0.25 \times 2 + 0.25 \times 2 + 0.25 \times 2 + 0.1 \times 3 + 0.1 \times 4 + 0.05 \times 4 = 2.4.$$

(c) From Slide IDC 6-26, the average codeword length of an i.i.d. source $S_1, S_2, S_2, \ldots, S_n$ satisfies

$$H(S) = \frac{1}{n}H(S^n) \le \frac{1}{n}\bar{L}_n \le \frac{1}{n}H(S^n) + \frac{1}{n} = H(S) + \frac{1}{n}.$$

As the code in (b) only uses n = 1. Thus, the optimal lossless data compression yields

$$H(S) \le \bar{L}_1^* = 2.4 \le H(S) + 1,$$

where

$$H(S) = \frac{1}{4}\log_2\frac{1}{\frac{1}{4}} + \frac{1}{4}\log_2\frac{1}{\frac{1}{4}} + \frac{1}{4}\log_2\frac{1}{\frac{1}{4}} + 0.1\log_2\frac{1}{0.1} + 0.1\log_2\frac{1}{0.1} + 0.05\log_2\frac{1}{0.05}$$

= 1.8 + 0.25 log₂(5)
= 2.38...

Note: When taking n = 1, one can show that if all the probabilities p_1, p_2, \ldots, p_M are in the form of $\frac{1}{2^k}$ for some k, then the ACL of Huffman code is equal to H(S). For example, given

$$p_1 = \frac{1}{2}, \quad p_2 = p_3 = \frac{1}{4},$$

we have H(S) = 1.5 and the Huffman code $\{0, 10, 11\}$ also has the ACL equal to 1.5 bits.

- 3. (a) Compress 10101101010010100101001 by the Lempel-Ziv coding algorithm in Slide IDC 6-33.
 - (b) Decompress 010001111000101011011000 by the Lempel-Ziv coding algorithm if the decompressor knows the index is fixed as 3 bits in length.

Solution.

(a) We first parse the input stream as

and build the below table:

decimal index	1	2	3	4	5	6	7
binary index	001	010	011	100	101	110	111
string	0	1	10	101	1010	10100	101001

where the first two indices are default. The Lempel-Ziv coding algorithm then outputs:

(010, 0)(011, 1)(100, 0)(101, 0)(110, 1)(100, 0)

i.e.,

010001111000101011011000.

(b) The decompressor will treat the sequence as a sequential concatenation of (index, last bit):

 $(010,0)(011,1)(100,0)(101,0)(110,1)(100,0) \equiv (2,0)(3,1)(4,0)(5,0)(6,1)(4,0)$

Then, each step takes in one (index, last bit). Table can be built on the fly.

• $(4,0) \rightarrow 10$	$(4, 0) \rightarrow 1010$ and renew table as	inde	x	1	2	3	4	5					
	$(4,0) \rightarrow 1010$ and renew table as	strin	g	0	1	10	101	l 101	0				
• $(5,0) \rightarrow 101$	$(5,0)$ \rightarrow 10100 and renow table	ind	$\mathbf{e}\mathbf{x}$	1	2	: 3	4	1	5	6			
	$(5,0) \rightarrow 10100$ and renew table (as stri	ng	0	1	. 1() 1()1 10)10	1010	00		
• $(6,1) \rightarrow 1010$	$(6,1)$ \rightarrow 101001 and renew table	ind	ex	1	2	3	4	1 I	5	6		7	
	$(0,1) \rightarrow 101001$ and renew table	as stri	ng	0	1	1() 1()1 10	10	1010	00	10100)1
• $(4, 0) \rightarrow 101$	$(4, 0) \rightarrow 1010$ and table remains	index	1	2		3	4	5		6		7	
	$(4,0) \rightarrow 1010$ and table remains -	string	0	1		10	101	1010	1(0100	10)1001	
	Thus the decompressed sequence is												

1010110101001010011010.

Note: We start the index by 1 by following the textbook. See page 580 in textbook. It should be okay to start the index by 0.

4. From Slide IDC 6-60, the capacity for the continuous-input AWGN channel is given by

$$C = \frac{1}{2}\log_2\left(1 + \frac{P}{\sigma^2}\right)$$
 bits per channel usage,

where $\sigma^2 = \frac{N_0}{2}$, and P (which sets the power constraint on the system, i.e., $\frac{1}{n} \sum_{i=1}^{n} E[|X_i|^2] \leq P$) is measured in **Joule per channel usage**. What is the minimum E_b/N_0 required for reliable transmission, subject to $R = k/n = \frac{1}{2}$ bits/channel usage?

Solution. (You shall be careful not to compare two quantities of different units.)

The energy of n transmissions (equivalently, n channel usages) is kE_b . Thus, in average, we have

$$P = \frac{kE_b}{n} = RE_b \text{ Joule/channel usage.}$$

The noise power experienced in each transmission is $\sigma^2 = \frac{N_0}{2}$. Shannon then said that reliable transmission is possible only when

R bits/channel usage < C bits/channel usage.

By Shannon's formula, we know reliable transmission is possible if

$$R = \frac{k}{n} = \frac{1}{2} \text{ bits/channel usage} < \frac{1}{2} \log_2 \left(1 + \frac{RE_b}{\frac{N_0}{2}} \right) = \frac{1}{2} \log_2 \left(1 + \frac{E_b}{N_0} \right) \text{ bits/channel usage}$$

which is equivalent to

$$1 < \frac{E_b}{N_0}.$$

Hence, the minimum E_b/N_0 required for reliable transmission over the continuous-input AWGN channel, subject to $R = k/n = \frac{1}{2}$ bits/channel usage, is $10 \log_{10}(1) = 0$ dB.

Note: From Slide IDC 6-77, the minimum E_b/N_0 required for reliable transmission over the binary-input AWGN channel, subject to $R = k/n = \frac{1}{2}$ bits/channel usage, is 0.186 dB, which is only slightly larger than 0 dB. This confirms the effectiveness of digital communications. 5. (Just for your reference. Not a part of the quiz or exam. Determination of Channel Capacity by Lagrange Multipliers Technique) In order to simplify the notation, we denote

$$p_i = P_X(i), \quad p = (p_1, \dots, p_I), \text{ and } q_{j|i} = P_{Y|X}(j|i)$$

and let the ranges of i and j be $1, \ldots, I$ and $1, \ldots, J$, respectively. Then, the mutual information between channel input X and output Y is equal to

$$I(X;Y) = f(\mathbf{p}) = \sum_{i=1}^{I} \sum_{j=1}^{J} p_i q_{j|i} \log_2 \frac{q_{j|i}}{\sum_{i'=1}^{I} p_{i'} q_{j|i'}}$$

The channel capacity is then given by

$$C \triangleq \max_{\boldsymbol{p} = (p_1, \dots, p_{\mathrm{I}}) \in [0, 1]^{\mathrm{I}} : \sum_{i=1}^{\mathrm{I}} p_i = 1} f(\boldsymbol{p}) = \max_{\boldsymbol{p} = (p_1, \dots, p_{\mathrm{I}}) \in \mathcal{Q}} f(\boldsymbol{p}),$$
(3)

where

$$\mathcal{Q} = \left\{ \boldsymbol{p} = (p_1, \dots, p_{\mathtt{I}}) \in [0, 1)^{\mathtt{I}} : \sum_{i=1}^{\mathtt{I}} p_i = 1 \right\}.$$

Note that here we set each $0 \le p_i < 1$ because when $p_{i''} = 1$ for some i'', we have $p_i = 0$ for $i \ne i''$, implying

$$I(X;Y) = f(\mathbf{p}) = \sum_{i=1}^{\mathsf{I}} \sum_{j=1}^{\mathsf{J}} p_i q_{j|i} \log_2 \frac{q_{j|i}}{\sum_{i'=1}^{\mathsf{I}} p_{i'} q_{j|i'}} = \sum_{j=1}^{\mathsf{J}} q_{j|i''} \log_2 \frac{q_{j|i''}}{q_{j|i''}} = 0.$$

Thus, such an input distribution can be excluded in our determination of channel capacity.

(a) Using the Lagrange multipliers technique, we rewrite (3) as

$$C = \max_{\boldsymbol{p} \in \mathcal{Q}} f(\boldsymbol{p}) = \max_{\boldsymbol{p} \in \mathcal{Q}} f_{\lambda}(\boldsymbol{p}) \leq \underbrace{\max_{\boldsymbol{p} \in [0,1]^{\mathrm{I}}} f_{\lambda}(\boldsymbol{p})}_{\text{unconstrained}},$$

where

$$f_{\lambda}(\boldsymbol{p}) \triangleq f(\boldsymbol{p}) + \lambda \left(\sum_{i=1}^{\mathtt{I}} p_i - 1\right).$$

Prove that $\boldsymbol{p}_{\lambda}^{*}$, which maximizes $f_{\lambda}(\boldsymbol{p}, \boldsymbol{Q})$ over $\boldsymbol{p} \in [0, 1)^{\mathrm{I}}$, satisfies

$$\begin{cases} d_i(\boldsymbol{p}^*_{\lambda}) = \log_2(e) - \lambda & \text{if } p^*_{\lambda,i} > 0\\ d_i(\boldsymbol{p}^*_{\lambda}) \le \log_2(e) - \lambda & \text{if } p^*_{\lambda,i} = 0 \end{cases}$$

$$\tag{4}$$

where

$$d_i(\boldsymbol{p}) \triangleq \sum_{j=1}^{\mathsf{J}} q_{j|i} \log_2 \frac{q_{j|i}}{\sum_{i'=1}^{\mathsf{T}} p_{i'} q_{j|i'}}.$$

Note that $d_i(\mathbf{p})$ is usually denoted as I(i; Y) in the literature.

Hint: Because

$$\frac{\partial^2 f_{\lambda}(\boldsymbol{p})}{\partial^2 p_{i''}} < 0,$$

 $f_{\lambda}(\mathbf{p})$ is concave in $p_{i''} \in [0, 1)$ and hence the optimizer satisfies

$$\begin{cases} \left. \frac{\partial f_{\lambda}(\boldsymbol{p})}{\partial p_{i''}} \right|_{\boldsymbol{p}=\boldsymbol{p}_{\lambda}^{*}} = 0 \quad \text{if } p_{\lambda,i''}^{*} > 0\\ \left. \frac{\partial f_{\lambda}(\boldsymbol{p})}{\partial p_{i''}} \right|_{\boldsymbol{p}=\boldsymbol{p}_{\lambda}^{*}} \leq 0 \quad \text{if } p_{\lambda,i''}^{*} = 0 \end{cases}$$

(b) For binary symmetric channel with $0 < \epsilon < \frac{1}{2}$, we have I = J = 2, and

$$q_{j|i} = \begin{cases} \epsilon, & j \neq i \\ 1 - \epsilon, & j = i \end{cases}$$

Hence,

$$d_1(p_1, p_2) \triangleq \sum_{j=1}^2 q_{j|1} \log_2 \frac{q_{j|1}}{\sum_{i'=1}^2 p_{i'} q_{j|i'}} = q_{1|1} \log_2 \frac{q_{1|1}}{\sum_{i'=1}^2 p_{i'} q_{1|i'}} + q_{2|1} \log_2 \frac{q_{2|1}}{\sum_{i'=1}^2 p_{i'} q_{2|i'}}$$
$$= (1 - \epsilon) \log_2 \frac{(1 - \epsilon)}{p_1(1 - \epsilon) + p_2\epsilon} + \epsilon \log_2 \frac{\epsilon}{p_1\epsilon + p_2(1 - \epsilon)}$$

and

$$d_2(p_1, p_2) \triangleq \epsilon \log_2 \frac{\epsilon}{p_1(1-\epsilon) + p_2\epsilon} + (1-\epsilon) \log_2 \frac{(1-\epsilon)}{p_1\epsilon + p_2(1-\epsilon)}$$

Using (4) to argue that $p^*_{\lambda,1} = p^*_{\lambda,2}$ for any λ .

(c) Continue from (b). $\boldsymbol{p}_{\lambda}^*$ in (b) is a function of λ . Can we find a λ such that $p_{\lambda,1}^* + p_{\lambda,2}^* = 1$ and therefore $\boldsymbol{p}_{\lambda}^* \in \mathcal{Q}$? Note: If such λ exists, then $C = f(\boldsymbol{p}_{\lambda}^*(\lambda))$.

Solution.

(a) We first derive

$$\begin{split} \frac{\partial f_{\lambda}(\boldsymbol{p})}{\partial p_{i''}} &= \frac{\partial}{\partial p_{i''}} \left\{ \sum_{i=1}^{I} \sum_{j=1}^{J} p_i q_{j|i} \log_2 q_{j|i} - \sum_{i=1}^{I} \sum_{j=1}^{J} p_i q_{j|i} \log_2 \left(\sum_{i'=1}^{I} p_{i'} q_{j|i'} \right) + \lambda \left(\sum_{i=1}^{I} p_i - 1 \right) \right\} \\ &= \sum_{j=1}^{J} q_{j|i''} \log_2 q_{j|i''} - \left(\sum_{j=1}^{J} q_{j|i''} \log_2 \left[\sum_{i'=1}^{I} p_{i'} q_{j|i'} \right] + \log_2(e) \sum_{i=1}^{I} \sum_{j=1}^{J} p_i q_{j|i} \frac{q_{j|i''}}{\sum_{i'=1}^{I} p_{i'} q_{j|i'}} \right) + \lambda \\ &= \left(\sum_{j=1}^{J} q_{j|i''} \log_2 q_{j|i''} - \sum_{j=1}^{J} q_{j|i''} \log_2 \left[\sum_{i'=1}^{I} p_{i'} q_{j|i'} \right] \right) - \log_2(e) \sum_{i=1}^{I} \sum_{j=1}^{J} p_i q_{j|i} \frac{q_{j|i''}}{\sum_{i'=1}^{I} p_{i'} q_{j|i'}} + \lambda \\ &= d_{i''}(\boldsymbol{p}) - \log_2(e) \sum_{j=1}^{J} \left[\sum_{i'=1}^{I} p_{i} q_{j|i} \right] \frac{q_{j|i''}}{\sum_{i'=1}^{I} p_{i'} q_{j|i'}} + \lambda \\ &= d_{i''}(\boldsymbol{p}) - \log_2(e) \sum_{j=1}^{J} q_{j|i''} + \lambda \\ &= d_{i''}(\boldsymbol{p}) - \log_2(e) + \lambda \end{split}$$

and

$$\frac{\partial^2 f_{\lambda}(\boldsymbol{p})}{\partial^2 p_{i''}} = \frac{\partial d_{i''}(\boldsymbol{p})}{\partial p_{i''}} \\
= \frac{\partial}{\partial p_{i''}} \left(\sum_{j=1}^{J} q_{j|i''} \log_2 q_{j|i''} - \sum_{j=1}^{J} q_{j|i''} \log_2 \left[\sum_{i'=1}^{I} p_{i'} q_{j|i'} \right] \right) \\
= -\log_2(e) \sum_{j=1}^{J} q_{j|i''} \frac{q_{j|i''}}{\sum_{i'=1}^{I} p_{i'} q_{j|i'}} < 0.$$

Hence, the optimizer satisfies

$$\begin{cases} \frac{\partial f_{\lambda}(\boldsymbol{p})}{\partial p_{i''}} \Big|_{\boldsymbol{p}=\boldsymbol{p}_{\lambda}^{*}} = 0 \quad \text{if } p_{\lambda,i''}^{*} > 0 \\ \frac{\partial f_{\lambda}(\boldsymbol{p})}{\partial p_{i''}} \Big|_{\boldsymbol{p}=\boldsymbol{p}_{\lambda}^{*}} \leq 0 \quad \text{if } p_{\lambda,i''}^{*} = 0 \end{cases} = \begin{cases} d_{i''}(\boldsymbol{p}_{\lambda}^{*}) - \log_{2}(e) + \lambda = 0 \quad \text{if } p_{\lambda,i''}^{*} > 0 \\ d_{i''}(\boldsymbol{p}_{\lambda}^{*}) - \log_{2}(e) + \lambda \leq 0 \quad \text{if } p_{\lambda,i''}^{*} = 0 \end{cases}$$
$$= \begin{cases} d_{i''}(\boldsymbol{p}_{\lambda}^{*}) = \log_{2}(e) - \lambda \quad \text{if } p_{\lambda,i''}^{*} > 0 \\ d_{i''}(\boldsymbol{p}_{\lambda}^{*}) \leq \log_{2}(e) - \lambda \quad \text{if } p_{\lambda,i''}^{*} = 0 \end{cases}$$

(b) We first note that $p_{\lambda,1}^* < 1$ implies $p_{\lambda,2}^* > 0$, and vice versus. Hence, (4) indicates $d_1(p_{\lambda,1}^*, p_{\lambda,2}^*) = d_2(p_{\lambda,1}^*, p_{\lambda,2}^*) = \log_2(e) - \lambda,$

$$d_1(p_{\lambda,1}^*, p_{\lambda,2}^*) = d_2(p_{\lambda,1}^*, p_{\lambda,2}^*) = \log_2(e) - \lambda,$$

which implies

$$(1-\epsilon)\log_2\frac{(1-\epsilon)}{p_{\lambda,1}^*(1-\epsilon)+p_{\lambda,2}^*\epsilon}+\epsilon\log_2\frac{\epsilon}{p_{\lambda,1}^*\epsilon+p_{\lambda,2}^*(1-\epsilon)}$$

= $\epsilon\log_2\frac{\epsilon}{p_{\lambda,1}^*(1-\epsilon)+p_{\lambda,2}^*\epsilon}+(1-\epsilon)\log_2\frac{(1-\epsilon)}{p_{\lambda,1}^*\epsilon+p_{\lambda,2}^*(1-\epsilon)}$

equivalently,

$$\begin{aligned} (1-2\epsilon)\log_2(p_{\lambda,1}^*(1-\epsilon)+p_{\lambda,2}^*\epsilon) &= (1-2\epsilon)\log_2(p_{\lambda,1}^*\epsilon+p_{\lambda,2}^*(1-\epsilon))\\ \Leftrightarrow \quad (p_{\lambda,1}^*(1-\epsilon)+p_{\lambda,2}^*\epsilon) &= (p_{\lambda,1}^*\epsilon+p_{\lambda,2}^*(1-\epsilon))\\ \Leftrightarrow \quad p_{\lambda,1}^* &= p_{\lambda,2}^*. \end{aligned}$$

(c) Last, we solve ${\pmb p}^*_\lambda = (p^*_\lambda, p^*_\lambda)$ via

$$d_1(p_{\lambda}^*, p_{\lambda}^*) = (1-\epsilon) \log_2 \frac{(1-\epsilon)}{p_{\lambda}^*(1-\epsilon) + p_{\lambda}^*\epsilon} + \epsilon \log_2 \frac{\epsilon}{p_{\lambda}^*\epsilon + p_{\lambda}^*(1-\epsilon)}$$
$$= (1-\epsilon) \log_2 \frac{(1-\epsilon)}{p_{\lambda}^*} + \epsilon \log_2 \frac{\epsilon}{p_{\lambda}^*}$$
$$= (1-\epsilon) \log_2(1-\epsilon) + \epsilon \log_2(\epsilon) - \log_2(p_{\lambda}^*) = \log_2(\epsilon) - \lambda$$

and obtain

$$(1-\epsilon)\log_2(1-\epsilon) + \epsilon\log_2(\epsilon) - \log_2(p_\lambda^*) = \log_2(e) - \lambda$$

Since we must have $p_{\lambda}^* = \frac{1}{2}$ (because $p_{\lambda,1}^* + p_{\lambda,2}^* = 1$ and $p_{\lambda,1}^* = p_{\lambda,2}^*$) in order to have $p_{\lambda}^* \in \mathcal{Q}$, setting

$$\lambda = -(1 - \epsilon) \log_2(1 - \epsilon) - \epsilon \log_2(\epsilon) - 1 + \log_2(\epsilon)$$

will make $p_{\lambda}^* \in Q$. In such case, $\log_2(e) - \lambda = 1 + (1 - \epsilon) \log_2(1 - \epsilon) + \epsilon \log_2(\epsilon)$ and

$$C = f\left(\frac{1}{2}, \frac{1}{2}\right) = \sum_{i=1}^{2} p_{\lambda,i}^{*} d_{i}(\boldsymbol{p}_{\lambda}^{*})$$
$$= \sum_{i=1}^{2} p_{\lambda}^{*}(\log_{2}(e) - \lambda)$$
$$= \log_{2}(e) - \lambda$$
$$= 1 - \underbrace{\left[(1-\epsilon)\log_{2}\frac{1}{(1-\epsilon)} + \epsilon\log_{2}\frac{1}{\epsilon}\right]}_{H(\epsilon)}$$
$$= 1 - H(\epsilon) \text{ info bits/channel usage}$$