

# Homework 3 Solution

## Homework Policy: (READ BEFORE YOU START TO WORK)

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

## 1. (Linear coding achieves the BSC capacity) [14]

In the lecture, we show that for the channel coding problem over a binary symmetric channel  $\text{BSC}(p)$ ,  $\forall \delta > 0$  and  $\forall \epsilon \in (0, 1)$ , there exists a codebook  $\mathcal{C}$  of size  $2^k$  such that

$$k > n \left( d_b(p \parallel \tfrac{1}{2}) - \delta \right) \quad \text{and} \quad P_{\mathbf{e}, \text{ML}}^{(n)} \leq \epsilon.$$

The key to the proof is a random coding argument, where the random ensemble given in the lecture is such that

$$X_i(w) \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(1/2) \quad \forall i = 1, 2, \dots, n, \quad \forall w = 1, 2, \dots, 2^k,$$

and as a result, a key inequality is established as follows:

$$\Pr\{\mathbf{w}(\mathbf{X}(w) \oplus \mathbf{Z} \oplus \mathbf{X}(\tilde{w})) \leq n(p + \epsilon)\} \leq 2^{-n d_b(p + \epsilon \parallel \frac{1}{2})} \quad (1)$$

In this problem, let us consider an alternative random ensemble that comprises *linear codes*, that is, the codeword of a message bit vector  $w \equiv \mathbf{b} = [b_1 \ b_2 \ \dots \ b_k]$  is

$$\mathbf{x}(\mathbf{b}) = \mathbf{b}\mathbf{g},$$

for some binary matrix  $\mathbf{g} \in \{0, 1\}^{k \times n}$ , and the above arithmetic is in the binary field. In other words, the encoding function is a linear transform governed by the matrix  $\mathbf{g}$ , and the

codebook  $\mathcal{C}_{\mathbf{g}} = \{\mathbf{b}\mathbf{g} \mid \mathbf{b} \in \{0, 1\}^k\}$  can be viewed as a subspace in the vector space  $\{0, 1\}^n$  with  $k$  dimensions. The subscript associated to  $\mathcal{C}$  is to emphasize its dependency with the *generator* matrix  $\mathbf{g}$ .

The alternative random ensemble is  $\mathcal{C}_{\mathbf{G}}$ , where the matrix  $\mathbf{G}$  is random with

$$\mathbf{G}_{j,i} \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(1/2) \quad \forall j = 1, 2, \dots, k, \quad \forall i = 1, 2, \dots, n.$$

a) Show that the random codewords in  $\mathcal{C}_{\mathbf{G}}$  are pairwise independent, that is,

$$\mathbf{X}(\mathbf{b}) \perp\!\!\!\perp \mathbf{X}(\tilde{\mathbf{b}}) \quad \forall \mathbf{b} \neq \tilde{\mathbf{b}}. \quad [6]$$

b) Show that the random codewords in  $\mathcal{C}_{\mathbf{G}}$  are NOT mutually independent. [2]

c) Show that (1) holds and conclude that linear coding achieves the BSC capacity. [6]

### Solution:

a) By definition of the random codebook, any random codeword can be represented by  $\mathbf{b}\mathbf{G}$ , where  $\mathbf{b}$  is the message corresponding to it. Let  $\mathcal{I}$  be the set that collects the indices where  $\mathbf{b}$  has entry 1. Hence, we have for  $\mathbf{b} \neq 0$ ,

$$\begin{aligned} \Pr\{\mathbf{b}\mathbf{G} = \mathbf{c}\} &= \Pr\{\mathbf{b}\mathbf{G}_1 = c_1, \dots, \mathbf{b}\mathbf{G}_n = c_n\} \quad , \text{ where } \mathbf{G} = (\mathbf{G}_1 \dots \mathbf{G}_n) \text{ and } \mathbf{c} = (c_1 \dots c_n) \\ &= \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,1} = c_1, \dots, \bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,n} = c_n\right\} \\ &= \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,1} = c_1\right\} \dots \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,n} = c_n\right\} \quad (\text{by independence of } \mathbf{G}) \\ &= \frac{1}{2^n} \end{aligned}$$

Now, consider another codeword  $\mathbf{b}'\mathbf{G}$  and let its corresponding index set be  $\mathcal{I}'$ .

$$\begin{aligned} \Pr\{\mathbf{b}\mathbf{G} = \mathbf{c}, \mathbf{b}'\mathbf{G} = \mathbf{c}'\} &= \Pr\{\mathbf{b}\mathbf{G}_1 = c_1, \dots, \mathbf{b}\mathbf{G}_n = c_n, \mathbf{b}'\mathbf{G}_1 = c'_1, \dots, \mathbf{b}'\mathbf{G}_n = c'_n\} \\ &= \Pr\{\mathbf{b}\mathbf{G}_1 = c_1, \mathbf{b}'\mathbf{G}_1 = c'_1\} \dots \Pr\{\mathbf{b}\mathbf{G}_n = c_n, \mathbf{b}'\mathbf{G}_n = c'_n\} \end{aligned}$$

$$\begin{aligned} &\Pr\{\mathbf{b}\mathbf{G}_j = c_j, \mathbf{b}'\mathbf{G}_j = c'_j\} \\ &= \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,j} = c_j, \bigoplus_{i \in \mathcal{I}'} \mathbf{G}_{i,j} = c'_j\right\} \\ &= \Pr\left\{\bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 0\right\} \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,j} = c_j, \bigoplus_{i \in \mathcal{I}'} \mathbf{G}_{i,j} = c'_j \mid \bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 0\right\} \\ &\quad + \Pr\left\{\bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 1\right\} \Pr\left\{\bigoplus_{i \in \mathcal{I}} \mathbf{G}_{i,j} = c_j, \bigoplus_{i \in \mathcal{I}'} \mathbf{G}_{i,j} = c'_j \mid \bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 1\right\} \end{aligned}$$

$$\begin{aligned}
&= \Pr \left\{ \bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 0 \right\} \Pr \left\{ \bigoplus_{i \in \mathcal{I} \setminus \mathcal{I}'} \mathbf{G}_{i,j} = c_j, \bigoplus_{i \in \mathcal{I}' \setminus \mathcal{I}} \mathbf{G}_{i,j} = c'_j \right\} \\
&\quad + \Pr \left\{ \bigoplus_{i \in \mathcal{I} \cap \mathcal{I}'} \mathbf{G}_{i,j} = 1 \right\} \Pr \left\{ \bigoplus_{i \in \mathcal{I} \setminus \mathcal{I}'} \mathbf{G}_{i,j} = \bar{c}_j, \bigoplus_{i \in \mathcal{I}' \setminus \mathcal{I}} \mathbf{G}_{i,j} = \bar{c}'_j \right\} \\
&= \frac{1}{4} \quad (\text{by discussing whether } \mathcal{I} \setminus \mathcal{I}' \text{ or } \mathcal{I}' \setminus \mathcal{I} \text{ is empty})
\end{aligned}$$

As a result,  $\Pr \{ \mathbf{b}\mathbf{G} = \mathbf{c}, \mathbf{b}'\mathbf{G} = \mathbf{c}' \} = \frac{1}{4^n} = \Pr \{ \mathbf{b}\mathbf{G} = \mathbf{c} \} \Pr \{ \mathbf{b}'\mathbf{G} = \mathbf{c}' \}$

(The above still holds for  $\mathbf{b} = 0$ )

This proves the pairwise independence.

b) We can easily find the dependence by

$$\mathbf{c}_{(0,\dots,0,1)} + \mathbf{c}_{(0,\dots,1,0)} = (0, \dots, 0, 1) \mathbf{G} + (0, \dots, 1, 0) \mathbf{G} = (0, \dots, 1, 1) \mathbf{G} = \mathbf{c}_{(0,\dots,1,1)}$$

c) Basically, all steps in the slide can be followed. The key here is to mention that pairwise independence is sufficient to conclude that  $w(X(b) + X(\tilde{b}) + Z)$  follows  $\text{Binom}(n, \frac{1}{2})$ .

### Grading Policy

- a) Conceptual explanation [4] Mathematical details [2]
- b) Clear counter example [2]
- c) Point out where pairwise independence suffices [6]

## 2. (Data processing) [10]

Recall the data processing inequality for information divergence in L2:

$$D(\mathbf{P}_X \| \mathbf{Q}_X) \geq D(\mathbf{P}_Y \| \mathbf{Q}_Y) \quad (2)$$

where  $\mathbf{P}_Y$  and  $\mathbf{Q}_Y$  are the output distributions of a data processing block  $\mathbf{W}_{Y|X}$  when the input  $X$  follows  $\mathbf{P}_X$  and  $\mathbf{Q}_X$  respectively.

Suppose now for some special  $\mathbf{W}_{Y|X}$ , the above inequality (2) can be *strengthened* to

$$D(\mathbf{P}_X \| \mathbf{Q}_X) \geq \eta D(\mathbf{P}_Y \| \mathbf{Q}_Y) \quad \forall \mathbf{P}_X, \mathbf{Q}_X, \quad (3)$$

where  $\eta$  is some constant and  $\eta > 1$ .

Based on (3), prove the following *strengthened* data processing inequality for mutual information: if  $U - X - Y$  forms a Markov chain and the conditional law of  $Y$  given  $X$  is  $\mathbf{W}_{Y|X}$ , then

$$I(U; X) \geq \eta I(U; Y).$$

**Solution:**

$$\begin{aligned}
I(U; X) &= D(P_{X|U} \| P_X | P_U) && \text{(by p.27 in slide L3_NCC_ho)} \\
&= E_U[D(P_{X|U} \| P_X)] \\
&\geq E_U[\eta D(P_{Y|U} \| P_Y)] && \text{(by strengthened data processing inequality)} \\
&= \eta D(P_{Y|U} \| P_Y | P_U) \\
&= \eta I(U; Y).
\end{aligned}$$

**Grading Policy**

Correctly relate mutual information and KL divergence [3]

Correctly apply strengthened DPI [3]

Correct proof [4]

**3. (Capacity of the permutation channel) [12]**

A channel model in neural communication is the following:

- Input/output alphabet:  $\mathcal{X} = \mathcal{Y} = \{0, 1\}^d$
- Channel law:

$$P_{Y|X}(\mathbf{y}|\mathbf{x}) = \begin{cases} 1/\binom{d}{\|\mathbf{x}\|_1}, & \text{if } \|\mathbf{y}\|_1 = \|\mathbf{x}\|_1 \\ 0, & \text{otherwise} \end{cases}$$

Note that for a  $d$ -dimensional *binary* vector  $\mathbf{x}$ , its  $\ell_1$ -norm is the number of 1's in  $\mathbf{x}$ :

$$\|\mathbf{x}\|_1 = \sum_{i=1}^d \mathbb{1}\{x_i = 1\}.$$

In words, the channel permutes the length- $d$  binary vector uniformly at random. In this problem, let us compute the capacity of this channel, namely, find

$$C = \max_{P_{\mathbf{X}}} I(\mathbf{X}; \mathbf{Y}).$$

a) (Warm-up) Let  $L := \|\mathbf{X}\|_1$ . Show that  $I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{X}; \mathbf{Y}|L) + H(L)$ . [2]

b) Show that

$$C = \max_{P_L} \left\{ H(L) + \max_{P_{\mathbf{X}|L}} I(\mathbf{X}; \mathbf{Y}|L) \right\}$$

and compute the channel capacity  $C$  accordingly. What is the capacity achieving input distribution? [6]

c) Let  $\alpha$  be a constant between 0 and 1, that is,  $0 < \alpha < 1$ . Now suppose the channel delivers  $\mathbf{x}$  noiselessly with probability  $(1 - \alpha)$ , and permutes  $\mathbf{x}$  uniformly at random with probability  $\alpha$  (note: keeping  $\mathbf{x}$  the same is also one possible permutation.).

Compute the channel capacity  $C$  of this channel. What is the capacity achieving input distribution? [4]

**Solution:**

a)

$$\begin{aligned}
I(\mathbf{X}; \mathbf{Y}|L) + H(L) &= H(\mathbf{X}|L) - H(\mathbf{X}|\mathbf{Y}, L) + H(L) \\
&= H(\mathbf{X}, L) - H(\mathbf{X}|\mathbf{Y}, L) \\
&= H(\mathbf{X}) - H(\mathbf{X}|\mathbf{Y}) = I(\mathbf{X}; \mathbf{Y}) \quad (\text{since } H(L|\mathbf{X}) = 0 \text{ and } H(L|\mathbf{Y}) = 0)
\end{aligned}$$

b)

$$\begin{aligned}
C &= \max_{\mathbf{P}_{\mathbf{X}}} I(\mathbf{X}; \mathbf{Y}) \\
&= \max_{\mathbf{P}_{\mathbf{X}, L}: P_{\mathbf{X}, L}(\mathbf{x}, l) = 0 \text{ if } \|\mathbf{x}\|_1 \neq l} H(L) + I(\mathbf{X}; \mathbf{Y}|L) \quad (\text{by (a)}) \\
&= \max_{\mathbf{P}_L} \left( H(L) + \max_{\mathbf{P}_{\mathbf{X}|L}: P_{\mathbf{X}|L}(\mathbf{x}|l) = 0 \text{ if } \|\mathbf{x}\|_1 \neq l} I(\mathbf{X}; \mathbf{Y}|L) \right)
\end{aligned}$$

Hence, by the above formula and that  $I(\mathbf{X}; \mathbf{Y}|L = l) = 0 \forall l = 0, \dots, d$ , we have that

$$C = \max_{\mathbf{P}_L} H(L) = \log(d+1),$$

and a capacity achieving input distribution is  $\mathbf{P}_X(x) = \frac{1}{d+1} \frac{1}{\binom{d}{\|x\|_1}}$ .

c) For the modified channel, we can see that the property in (a) still holds, and thus the capacity of it can still be calculated by (b). Therefore, by solving the maximization step by step, we can derive the capacity and the input distribution. Also, we drop the subscript for convenience.

$$\begin{aligned}
\max_{\mathbf{P}_{\mathbf{X}|L}} I(\mathbf{X}; \mathbf{Y}|L) &= \sum_{l=0}^d P_L(l) \max_{\mathbf{P}_{X|L=l}} I(X; Y|L = l) \\
&= \sum_{l=0}^d P_L(l) \max_{\mathbf{P}_{X|L=l}} H(Y|L = l) - H(Y|X, L = l)
\end{aligned}$$

Note that the "sub-channel" when  $L = l$  is given is symmetric, so the maximum can be achieved by uniform conditional distribution i.e.  $\mathbf{P}_{X|L=l}(x) = \frac{1}{\binom{d}{l}}$ ,  $\forall \|x\|_1 = l$ .

In that case, we define

$$\begin{aligned}
C_l &:= \max_{\mathbf{P}_{X|L=l}} H(Y|L = l) - H(Y|X, L = l) \\
&= \log \binom{d}{l} + \left( 1 - \alpha + \frac{\alpha}{\binom{d}{l}} \right) \log \left( 1 - \alpha + \frac{\alpha}{\binom{d}{l}} \right) + \left( \binom{d}{l} - 1 \right) \frac{\alpha}{\binom{d}{l}} \log \frac{\alpha}{\binom{d}{l}}
\end{aligned}$$

Therefore,

$$\begin{aligned}
 C &= \max_{P_L} \left( H(L) + \max_{P_{\mathbf{X}|L}} I(\mathbf{X}; \mathbf{Y}|L) \right) \\
 &= \max_{P_L} \sum_{l=0}^d P_L(l) (-\log P_L(l) + C_l) \\
 &= \max_{P_L} \sum_{l=0}^d P_L(l) \log \frac{2^{C_l}}{P_L(l)} \leq \left( \log \sum_{l=0}^d 2^{C_l} \right) \quad (\text{by Jensen's inequality})
 \end{aligned}$$

Also, the upper bound can be achieved by  $P_L(l) = \frac{2^{C_l}}{\sum_{i=0}^d 2^{C_i}}$ .

Thus  $C = \log \left( \sum_{l=0}^d 2^{C_l} \right)$ .

### Grading Policy

- a) Correctness [2]
- b) Proof of the formula [2] Capacity [2] Achieving input distribution [2]
- c) Correctly Apply (b) [2] Calculation [2]

## 4. (List codes) [14]

In this problem, let us consider a variant of the channel coding problem over a DMC  $P_{Y|X}$  with Shannon capacity  $C$ .

Recall that in the formulation of channel coding, the decoder aims to uniquely decode the message. In practice, such an aim might be too stringent, and the decoder may just want to determine a *list* of plausible codewords from the channel output.

A  $(n, \lceil nR \rceil, \lceil nL \rceil)$  *list code* consists of

- An encoding function that maps each message  $w$  to a length- $n$  codeword  $x^n(w)$ , for each  $w \in \{1, \dots, 2^{\lceil nR \rceil}\}$ .
- A decoding function that maps a channel output length- $n$  sequence  $y^n$  to a list of messages  $\mathcal{L}(y^n) \subseteq \{1, \dots, 2^{\lceil nR \rceil}\}$  of size  $|\mathcal{L}| \leq 2^{\lceil nL \rceil}$ . An error occurs if the transmitted message is not contained in this list.

Hence, the definition of the *probability of error* becomes

$$P_e^{(n)} := \Pr\{W \notin \mathcal{L}(Y^n)\}.$$

A tuple  $(R, L)$  is said to be *achievable* if there exists a sequence of  $(n, \lceil nR \rceil, \lceil nL \rceil)$  list codes with

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

- a) Since the definition of the probability of error is changed, to prove the converse part of the coding theorem, a new Fano-type inequality is needed.

Consider a random variable  $U \in \mathcal{U} = \{1, 2, \dots, m\}$  and a random set  $\mathcal{V} \subseteq \mathcal{U}$ . For example, if  $\mathcal{V} \sim \text{Unif}(2^{\mathcal{U}})$  where  $2^{\mathcal{U}}$  denotes the collection of all subsets of  $\mathcal{U}$ , then  $H(\mathcal{V}) = \log_2(2^{|\mathcal{U}|}) = |\mathcal{U}| = m$ .

Suppose  $|\mathcal{V}| \leq v$  with probability 1. Let  $P_e := \Pr\{U \notin \mathcal{V}\}$ . Show that

$$H(U|\mathcal{V}) \leq H_b(P_e) + P_e \log(m) + (1 - P_e) \log(v). \quad [4]$$

- b) Show that for every sequence of  $(n, \lceil nR \rceil, \lceil nL \rceil)$  list codes with vanishing  $P_e^{(n)}$  as  $n \rightarrow \infty$ ,  $(R, L)$  must satisfy

$$R - L \leq C. \quad [4]$$

- c) Show that if  $R - L < C$ , then  $(R, L)$  is achievable. [6]

*Hint: You may design the encoder so that roughly  $2^{n(R-L)}$  distinct codewords need to be transmitted reliably over the channel.*

### Solution:

- a) Similar to the proof of Fano's inequality, let  $E = \mathbb{1}\{U \notin \mathcal{V}\}$ . Then:

$$\begin{aligned} H(U|\mathcal{V}) &\leq H(U, E|\mathcal{V}) = H(E|\mathcal{V}) + H(U|E, \mathcal{V}) \\ &\leq H(E) + \Pr\{E = 1\}H(U|E = 1, \mathcal{V}) + \Pr\{E = 0\}H(U|E = 0, \mathcal{V}) \\ &\leq H_b(P_e) + P_e \log m + (1 - P_e) \log v. \end{aligned}$$

- b) Follow similar arguments of the converse proof of the channel coding theorem:  
Let  $\mathcal{L}(Y^n)$  be the decoding list results from the channel output  $Y^n$ .

$$\begin{aligned} R &\leq \frac{k}{n} = \frac{H(W)}{n} = \frac{1}{n} (I(W; \mathcal{L}(Y^n)) + H(W|\mathcal{L}(Y^n))) \\ &\leq \frac{1}{n} I(W; Y^n) + \frac{1}{n} (1 + P_e^{(n)} \log 2^{\lceil nR \rceil} + (1 - P_e^{(n)})nL) \\ &=: \frac{1}{n} I(W; Y^n) + \epsilon_n. \end{aligned}$$

The second inequality follows from the data processing inequality and part (a), and  $\lim_{n \rightarrow \infty} \epsilon_n = 0$ .

Also,

$$\frac{1}{n} I(W; Y^n) \leq \frac{1}{n} \sum_{i=1}^n I(W; Y_i | Y^{i-1}).$$

Note that  $\forall i$ ,

$$I(W; Y_i | Y^{i-1}) \leq I(W, Y^{i-1}; Y_i) \leq I(W, Y^{i-1}, X_i; Y_i) = I(X_i; Y_i) \leq \max_{P_X} I(X; Y) = C.$$

Thus,  $\forall N \in \mathbb{N}$ ,

$$R \leq C + \epsilon_n.$$

Taking  $n \rightarrow \infty$ , we conclude that for every sequence of  $(2^{nR}, 2^{nL}, C)$  list codes with  $\lim_{n \rightarrow \infty} P_e^{(n)} = 0$ , we have  $R \leq C + L$ .

- c) By the hint, we group every  $2^{\lfloor nL \rfloor}$  messages into a set, with  $2^{\lceil nR \rceil - \lfloor nL \rfloor}$  sets. From each set, choose one representative message. Denote representatives as  $\{w'_1, \dots, w'_{2^{\lceil nR \rceil - \lfloor nL \rfloor}}\}$ , and their sets as  $S_{w'}$ .

To send message  $w$ , if  $w \in S_{w'}$ , send the representative  $w'$ . At the decoder, upon receiving  $w'$ , output  $S_{w'}$  as the list of size  $\leq 2^{nL}$ .

This reduces to a point-to-point channel coding problem with  $k = \lceil nR \rceil - \lfloor nL \rfloor$ . Applying noisy channel coding theorem, if  $R - L < C$ , exist a coding scheme such that the error probability  $\Pr(W \notin \mathcal{L}(Y^n))$  vanishes. Thus,  $(R, L)$  is achievable.

### Grading Policy

- a) Define  $E := \mathbb{1}\{U \notin \mathcal{V}\}$  [2] Correctness [2]
- b) Correctly apply (a) [2] Correctness [2]
- c) Encoding process [2] Decoding process [2] Error probability analysis [2]