

Lecture 7: Channel coding theorem for discrete-time continuous memoryless channel

Lectured by Dr. Saif K. Mohammed

Scribed by Mirsad Čirkić

Information Theory for Wireless Communication (ITWC) Spring 2012

Let us first define, for the random sequences $\mathbf{X} = [X_1, \dots, X_n]$ and $\mathbf{Y} = [Y_1, \dots, Y_n]$ and their corresponding sequence realizations $\mathbf{x} = [x_1, \dots, x_n]$ and $\mathbf{y} = [y_1, \dots, y_n]$ where $x_k, y_k \in \mathbb{R}$, the following probability density functions (pdfs): $f_{\mathbf{X}}(\mathbf{x}) \triangleq \prod_{k=1}^n f_X(x_k)$ as the input pdf, $f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^n f_{Y|X}(y_k|x_k)$ as the memoryless channel (transition) pdf, and $f_{\mathbf{Y}}(\mathbf{y})$ as the output pdf. Each element X_k (or just X without the index k) are independent and identically distributed i.i.d. Further, due to the i.i.d. property of X and the specified memoryless property imposed on the channel transition pdf, we have $f_{\mathbf{Y}}(\mathbf{y}) = \prod_{k=1}^n f_Y(y_k)$.

Define the so called and the mutual information quantity

Let

$$I(X; Y) \triangleq \int_{x,y} f_{X,Y}(x, y) \log \left(\frac{f_{X,Y}(x, y)}{f_X(x)f_Y(y)} \right) dx dy \quad (1)$$

define the mutual information between the random variables X and Y . Further, let $g : \mathbb{R} \rightarrow \mathbb{R}$ be the capacity cost function (mapping). Then, we define (n, N, λ) code with a capacity cost $S \in \mathbb{R}$ as a system

- $\{(\mathbf{u}_1, A_1), \dots, (\mathbf{u}_N, A_N)\}$
- $\mathbf{u}_i = [u_{i,1} \dots u_{i,n}] \in \mathbb{R}^n$, $A_i \subset \mathbb{R}^n$, $A_i \cap A_j = \emptyset \forall i \neq j$
- $\sum_{k=1}^n \frac{g(u_{i,k})}{n} \leq S$, $\forall i$
- $\bar{\lambda} \triangleq \frac{1}{N} \sum_{i=1}^N \lambda_i$, $\lambda_M \triangleq \max_{i=1, \dots, N} \lambda_i$, $\lambda_i \triangleq P(\mathbf{Y} \notin A_i | \mathbf{X} = \mathbf{u}_i)$, and $\lambda_M \leq \lambda$.

Theorem 1. Random Coding Theorem

We define the capacity with cost constraint S as

$$C(S) \triangleq \max_{f_X(\cdot)} I(X; Y) \quad (2)$$

s.t. $\int_x f_X(x)g(x) \leq S$.

This document is a property of Communication Systems Division, Department of Electrical Engineering, Linköping University, Sweden. Copyright must be obtained by writing to saif@isy.liu.se, erik.larsson@isy.liu.se prior to usage.

The theorem states that, For any “rate” R which satisfies $0 \leq R < C(S)$, there exists a sequence of $(n, 2^{nR}, \lambda(n))$ codes for $n = 1, 2, \dots$ such that $\lambda(n) \rightarrow 0$ as $n \rightarrow \infty$.

In what follows, a proof is presented and for that reason, we need to define

$$A_\epsilon(n) \triangleq \left\{ (\mathbf{x}, \mathbf{y}) : \begin{aligned} & \left| -\frac{\log f_{\mathbf{X}}(\mathbf{x})}{n} - h(X) \right| < \epsilon, \\ & \left| -\frac{\log f_{\mathbf{Y}}(\mathbf{y})}{n} - h(Y) \right| < \epsilon, \\ & \left| -\frac{\log f_{\mathbf{Y}, \mathbf{X}}(\mathbf{y}, \mathbf{x})}{n} - h(Y, X) \right| < \epsilon \end{aligned} \right\},$$

where $h(X) \triangleq -\int_{\mathbf{x}} f_X(\mathbf{x}) \log f_X(\mathbf{x}) d\mathbf{x}$. Further, we need to introduce the following two lemmas.

Lemma 1.

$$\boxed{\begin{array}{c} MS \\ f_X(x) \end{array}} \rightarrow \mathbf{X} \rightarrow \boxed{\begin{array}{c} MC \\ f_{Y|X}(y|x) \end{array}} \rightarrow \mathbf{Y}$$

- $P((\mathbf{X}, \mathbf{Y}) \in A_\epsilon(n)) > 1 - \epsilon$ for sufficiently large n
- $\text{Vol}(A_\epsilon(n)) \triangleq \int_{(\mathbf{x}, \mathbf{y}) \in A_\epsilon(n)} d\mathbf{x}d\mathbf{y} < 2^{n(h(Y, X) + \epsilon)}$
- $\text{Vol}(A_\epsilon(n)) > (1 - \epsilon)2^{n(h(Y, X) - \epsilon)}$ for sufficiently large n

We omit the explicit proofs since they become trivial from the discrete case proofs in previous lectures.

Lemma 2. Known as “the Packing Lemma”

For independent $\tilde{\mathbf{X}}$ and \mathbf{Y} , i.e.,

$$\boxed{\begin{array}{c} MS \\ f_X(x) \end{array}} \rightarrow \tilde{\mathbf{X}}$$

$$\boxed{\begin{array}{c} MS \\ f_X(x) \end{array}} \rightarrow \mathbf{X} \rightarrow \boxed{\begin{array}{c} MC \\ f_{Y|X}(y|x) \end{array}} \rightarrow \mathbf{Y}$$

the following inequalities hold for sufficiently large n ,

$$(1 - \epsilon)2^{-n(I(X; Y) + 3\epsilon)} < P\left((\tilde{\mathbf{X}}, \mathbf{Y}) \in A_\epsilon(n)\right) < 2^{-n(I(X; Y) - 3\epsilon)}.$$

$$P\left((\tilde{\mathbf{X}}, \mathbf{Y}) \in A_\epsilon(n)\right) = \int_{(\tilde{\mathbf{x}}, \mathbf{y}) \in A_\epsilon(n)} \underbrace{f_{\mathbf{X}}(\tilde{\mathbf{x}}) f_{\mathbf{Y}}(\mathbf{y})}_{f_{\tilde{\mathbf{X}}, \mathbf{Y}}(\tilde{\mathbf{x}}, \mathbf{y})} d\tilde{\mathbf{x}}d\mathbf{y}$$

From the definition of $A_\epsilon(n)$, we have $2^{-n(h(X) + \epsilon)} < f_{\mathbf{X}}(\tilde{\mathbf{x}}) < 2^{-n(h(X) - \epsilon)}$, $2^{-n(h(Y) + \epsilon)} < f_{\mathbf{Y}}(\mathbf{y}) < 2^{-n(h(Y) - \epsilon)}$. This gives us the upper bound

$$\begin{aligned} P\left((\tilde{\mathbf{X}}, \mathbf{Y}) \in A_\epsilon(n)\right) &< 2^{-n(h(X) + h(Y) - 2\epsilon)} \int_{(\tilde{\mathbf{x}}, \mathbf{y}) \in A_\epsilon(n)} d\tilde{\mathbf{x}}d\mathbf{y} \\ &< 2^{-n(h(X) + h(Y) - 2\epsilon)} \underbrace{2^{n(h(Y, X) + \epsilon)}}_{\text{lemma 1}} \\ &= 2^{-n(I(X; Y) - 3\epsilon)}, \end{aligned}$$

and analogously the lower bound $P\left((\tilde{\mathbf{X}}, \mathbf{Y}) \in A_\epsilon(n)\right) > 2^{-n(I(X:Y)+3\epsilon)}$, for sufficiently large n .

The proof of Theorem 1. relies on Shannon's random coding argument, which utilizes the fact that if the error probability averaged over all possible codes goes to zero as $n \rightarrow \infty$, then there must exist at least one code that provides an error probability smaller than the average.

1) *Random Code Generation:*

- Choose an arbitrary $\epsilon > 0$
- Choose a pdf $f_X(x)$ s.t. $\int_x f_X(x)g(x)dx = S - \epsilon < S$
- Generate $u_{i,k}$ for all $i = 1, \dots, N$ and $k = 1, \dots, n$ independently using $f_X(x)$.

2) *Encoding:* Let us label the possible codewords $\mathbf{u}_1, \dots, \mathbf{u}_N$ as messages $w = 1, \dots, N$. When the transmitter wants to communicate message k , it will transmit \mathbf{u}_k .

3) *Decoding:* Let us use the following typical set decoder (for the code \mathcal{C}) that, for a given channel output \mathbf{y} , estimates the transmitted message as or be in error if

$$\hat{w}_{\mathcal{C}}(\mathbf{y}) = \begin{cases} i, & \text{if } \mathbf{y} \in A_i(\mathcal{C}) \text{ and } \mathbf{y} \notin A_j(\mathcal{C}) \forall j \neq i \\ & \text{and } \sum_{k=1}^n \frac{g(u_{i,k})}{n} \leq S \\ 0, & \text{otherwise (i.e. decoding error)} \end{cases},$$

where $A_i(\mathcal{C}) \triangleq \{\mathbf{y} \in \mathbb{R}^n | (\mathbf{u}_i, \mathbf{y}) \in A_\epsilon(n)\}$.

The average error probability $\bar{\lambda}(\mathbf{u}_1, \dots, \mathbf{u}_N) \triangleq \frac{1}{N} \sum_{i=1}^N \lambda_i(\mathbf{u}_1, \dots, \mathbf{u}_N)$ and the error probability $\lambda_i(\mathbf{u}_1, \dots, \mathbf{u}_N) \triangleq P(\hat{w}(\mathbf{Y}) \neq i | \mathbf{X} = \mathbf{u}_i)$ are both deterministic for a particular realization of the code book $\mathcal{C} = (\mathbf{u}_1, \dots, \mathbf{u}_N)$. Each realization \mathbf{u}_i is drawn from the input pdf $f_X(x)$, and this is done independently over i . If we define $\mathbf{C} = (\mathbf{X}_1, \dots, \mathbf{X}_i, \dots, \mathbf{X}_N)$, where \mathbf{X}_i are i.i.d. with the pdf f_X , we can say that $f_{\mathcal{C}}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{i=1}^N f_X(\mathbf{x}_i)$ and that the code book \mathcal{C} is drawn from $f_{\mathcal{C}}(\mathbf{x}_1, \dots, \mathbf{x}_N)$. Now, the expectation of $\bar{\lambda}$ (taken over \mathcal{C} , i.e., the whole ensemble of codes) is

$$\mathbb{E}_{\mathcal{C}}\{\bar{\lambda}(\mathcal{C})\} = \mathbb{E}_{\mathcal{C}}\left\{\frac{1}{N} \sum_{i=1}^N \lambda_i(\mathcal{C})\right\} = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{C}}\{\lambda_i(\mathbf{X}_1, \dots, \mathbf{X}_N)\} = \mathbb{E}_{\mathcal{C}}\{\lambda_1(\mathbf{X}_1, \dots, \mathbf{X}_N)\},$$

where the last equality utilizes the fact that the expectation value is invariant to any reordering of $\mathbf{X}_1, \dots, \mathbf{X}_N$ due to the i.i.d. property of \mathbf{X}_i over i . Before we continue further, let us introduce a compact notation $A_0(\mathcal{C})$ for the set $\{\mathbf{y} | \sum_{k=1}^n \frac{g(u_{1,k})}{n} \leq S\}$, which may seem strange since the condition does not depend on \mathbf{y} . Essentially, the set $A_0(\mathcal{C})$ will either contain the whole \mathbf{y} -space if the transmitted codeword \mathbf{u}_1 fulfills the constraint on $g(\cdot)$, or it will be the null set when the

constraint is not met. Hence,

$$\begin{aligned}
\mathbb{E}_{\mathcal{C}}\{\lambda_1(\mathcal{C})\} &= \int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \left(\int_{\mathbf{y}; \hat{w}_{\mathcal{C}}(\mathbf{y}) \neq 1} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} \right) d\mathcal{C} \\
&= \int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \left(\int_{A_1^c(\mathcal{C}) \cup A_2(\mathcal{C}) \cdots \cup A_N(\mathcal{C}) \cup A_0^c(\mathcal{C})} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} \right) d\mathcal{C} \\
&\leq \underbrace{\int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \int_{A_1^c(\mathcal{C})} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C}}_{T_1} \\
&\quad + \sum_{i=2}^N \underbrace{\int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \int_{A_i(\mathcal{C})} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C}}_{T_i} \\
&\quad + \underbrace{\int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \int_{A_0^c(\mathcal{C})} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C}}_{T_0}.
\end{aligned}$$

We can continue from this point on by evaluating the integrals $T_0, T_1, T_2, \dots, T_N$ separately. Starting with T_0 , note that the condition defining the set A_0 does not depend on \mathbf{y} , which gives

$$\begin{aligned}
T_0 &= \int_{\mathcal{C} \in \mathbb{R}^{Nn}; \sum_{k=1}^n \frac{g(u_{1,k})}{n} \geq S} f_{\mathcal{C}}(\mathcal{C}) \int_{\mathbb{R}^n} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C} \\
&= \int_{\mathbf{u}_1 \in \mathbb{R}^n; \sum_{k=1}^n \frac{g(u_{1,k})}{n} \geq S} f_{\mathbf{X}}(\mathbf{u}_1) d\mathbf{u}_1 < \epsilon,
\end{aligned}$$

for sufficiently large n where the last inequality follows from the weak law of large numbers. That is, $\frac{1}{n} \sum_{k=1}^n g(u_{1,k}) \rightarrow \mathbb{E}_X\{g(X)\} = S - \epsilon$ in probability as $n \rightarrow \infty$. Continuing with T_1 ,

$$\begin{aligned}
T_1 &= \int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \int_{\mathbf{y}; (\mathbf{u}_1, \mathbf{y}) \notin A_{\epsilon}(n)} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C} \\
&= \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{u}_1) \int_{\mathbf{y}; (\mathbf{u}_1, \mathbf{y}) \notin A_{\epsilon}(n)} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathbf{u}_1 \\
&= \int_{A_{\epsilon}^c(n)} f_{\mathbf{Y}, \mathbf{X}}(\mathbf{y}, \mathbf{u}_1) d\mathbf{y} d\mathbf{u}_1 = P((\mathbf{X}, \mathbf{Y}) \notin A_{\epsilon}(n) | \mathbf{X} = \mathbf{u}_1) < \epsilon,
\end{aligned}$$

for sufficiently large n where the last inequality follows from lemma 1. Lastly T_i for $i = 2, \dots, N$,

$$\begin{aligned}
T_i &= \int_{\mathbb{R}^{Nn}} f_{\mathcal{C}}(\mathcal{C}) \int_{\mathbf{y}; (\mathbf{u}_i, \mathbf{y}) \in A_{\epsilon}(n)} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{u}_1) d\mathbf{y} d\mathcal{C} \\
&= \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{u}_i) \int_{\mathbf{y}; (\mathbf{u}_i, \mathbf{y}) \in A_{\epsilon}(n)} f_{\mathbf{Y}}(\mathbf{y}) d\mathbf{y} d\mathbf{u}_i \\
&= \int_{A_{\epsilon}(n)} f_{\mathbf{X}}(\mathbf{u}_i) f_{\mathbf{Y}}(\mathbf{y}) d\mathbf{y} d\mathbf{u}_i \leq 2^{-n(I(\mathbf{Y}; \mathbf{X}) - 3\epsilon)},
\end{aligned}$$

where the last inequality follows from lemma 2. These upper bounds result in

$$\mathbb{E}_{\mathcal{C}}\{\bar{\lambda}(\mathcal{C})\} < 2\epsilon + N2^{-n(I(\mathbf{Y}; \mathbf{X}) - 3\epsilon)} < 2\epsilon + 2^n R 2^{-n(I(\mathbf{Y}; \mathbf{X}) - 3\epsilon)} = 2\epsilon + 2^{-n(I(\mathbf{Y}; \mathbf{X}) - R - 3\epsilon)},$$

where $R \triangleq \frac{\log(N)}{n} \geq 0$ represents the ‘‘code rate’’. For $R < I(\mathbf{Y}; \mathbf{X}) - 3\epsilon$, we can choose sufficiently large n such that $\mathbb{E}_{\mathcal{C}}\{\bar{\lambda}(\mathcal{C})\} < 3\epsilon$.

Since $\mathbb{E}_C\{\bar{\lambda}(C)\} < 3\epsilon$, following Shannon's random coding argument, there exists a code $(n, 2^{nR}, \cdot)$ that has an average error probability that is smaller than 3ϵ . In essence, we will only pick a code book that will yield a small average error probability, which also means that only those codewords that satisfy the capacity cost constraint will be included. Note that the decoding regions A_i can be made disjoint without increasing the upper bound on the average error probability. So far, we have not said anything about the maximum error probability, which is the quantity of interest in order to communicate with a certain level of reliability. Nevertheless, given an upper bound on the average error probability, say 3ϵ , the best $N/2$ codewords (smallest error probabilities) have a maximum error probability that is smaller than 6ϵ . If it were not true, then the average error probability would not be smaller than 3ϵ due to the contribution of the worst $N/2$ codewords each having an error probability greater than 6ϵ , i.e., $\frac{1}{N} \frac{N}{2} 6\epsilon = 3\epsilon \not< 3\epsilon$. Using this fact, we can construct a new code consisting only of the $N/2$ best codewords, which would then reduce the new rate R' to

$$R' = \frac{\log(N/2)}{n} = R - \frac{1}{n}.$$

We can conclude now that since $S - \epsilon < S$, it follows that $R' < I(X; Y) < C(S)$. However by choosing $\epsilon > 0$ to be arbitrarily small and n sufficiently large, we can achieve any rate less than $C(S)$.

Capacity cost of the discrete-time continuous memoryless AWGN channel with an average input power constraint: Consider the discrete-time memoryless additive white Gaussian noise channel (AWGN), i.e., $Y = X + N$, with $X \in \mathbb{R}$ as the input, $Y \in \mathbb{R}$ as the output and $N \in \mathbb{R}$ as the AWGN with variance σ^2 . Further, let us consider an average power constraint on the input, i.e. $\mathbb{E}[X^2] \leq S$. The capacity cost for this channel can be computed using (2), with the cost function $g(x) = x^2$.

$$\begin{aligned} C(S) &= \max_{f_X(\cdot): \mathbb{E}_X\{X^2\} \leq S} I(X; Y) = \max_{f_X(\cdot): \mathbb{E}_X\{X^2\} \leq S} \underbrace{h(Y) - h(Y|X)}_{h(N)} \\ &= \max_{f_X(\cdot): \mathbb{E}_X\{X^2\} \leq S} h(Y) - \frac{1}{2} \log(2\pi e \sigma^2) = \max_{\text{Var}\{Y\}: \mathbb{E}_X\{X^2\} \leq S} \frac{1}{2} \log\left(\frac{\text{Var}\{Y\}}{\sigma^2}\right) \\ &= \frac{1}{2} \log\left(1 + \frac{S}{\sigma^2}\right). \end{aligned}$$

where the maximum is achieved when X is Gaussian distributed with mean 0 and variance S .

Later in Lecture 9, we will see that $\frac{1}{2} \log\left(1 + \frac{S}{\sigma^2}\right)$ is indeed the capacity of this channel.