

INFORMATION THEORY & CODING

Channel Code - 2

Dr. Rui Wang

Department of Electrical and Electronic Engineering
Southern Univ. of Science and Technology (SUSTech)

Email: wang.r@sustech.edu.cn

December 6, 2022



- **Channel capacity.** The logarithm of the number of distinguishable inputs is given by

$$C = \max_{p(x)} I(X; Y).$$

- **Examples**

- Binary symmetric channel: $C = 1 - H(p)$
- Binary erasure channel: $C = 1 - \alpha$
- Symmetric channel: $C = \log |\mathcal{Y}| - H$ (row of trans. matrix)

Definition

An (M, n) code for the channel $(\mathcal{X}, p(y|x), \mathcal{Y})$ consists of :

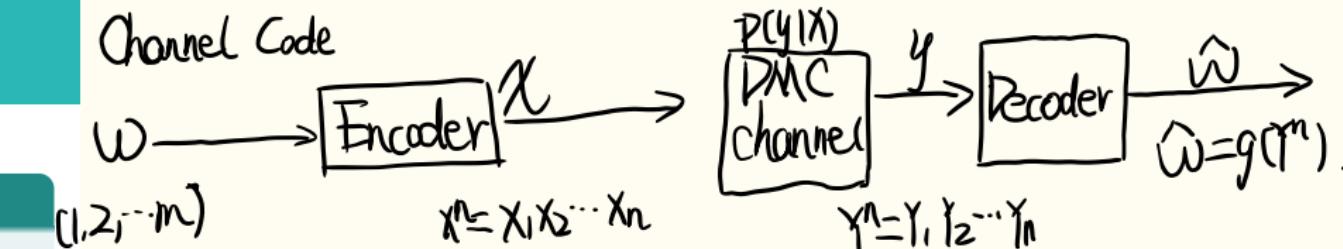
1. An index set $\{1, 2, \dots, M\}$ representing messages.
2. An encoding function $X^n : \{1, 2, \dots, M\} \rightarrow \mathcal{X}^n$, yielding codewords $x^n(1), x^n(2), \dots, x^n(M)$. The set of codewords is called **codebook**.
3. A decoding function $g : \mathcal{Y}^n \rightarrow \{1, 2, \dots, M\}$.

The rate R of an (M, n) code is

$$R = \frac{\log M}{n} \text{ bit per transmission}$$

On the other hand, we usually write

$$M = [2^{nR}]$$



An (M, n) code consists of

① An index set representing messages $\{1, 2, \dots, M\}$

② An encoder: $X^n(1), X^n(2), \dots, X^n(M)$

③ A decoder $g : \mathcal{Y}^n \rightarrow \{1, 2, 3, \dots, M\}$

Code Rate: $R = \frac{\log M}{n}$, 平均每 bit 中携带的信道信息量

$M = \lceil 2^{nR} \rceil \rightarrow$ 向上取整

R与C之间的关系 !!!

Definition

An (M, n) code for the channel $(\mathcal{X}, p(y|x), \mathcal{Y})$ consists of :

1. An index set $\{1, 2, \dots, M\}$ representing messages.
2. An encoding function $X^n : \{1, 2, \dots, M\} \rightarrow \mathcal{X}^n$, yielding codewords $x^n(1), x^n(2), \dots, x^n(M)$. The set of codewords is called **codebook**.
3. A decoding function $g : \mathcal{Y}^n \rightarrow \{1, 2, \dots, M\}$.

The **rate** R of an (M, n) code is

$$R = \frac{\log M}{n} \text{ bit per transmission}$$

On the other hand, we usually write

$$M = \lceil 2^{nR} \rceil$$

Performance Metric

- Conditional probability of error:

$$\lambda_i = \Pr[g(Y_n) \neq i | X^n = x^n(i)] = \sum_{y^n} p(y^n | x^n(i)) I(g(y^n) \neq i)$$

- Maximal probability of error: $\lambda^{(n)} = \max_{i \in \{1, 2, \dots, M\}} \lambda_i$
- Decoding error probability: $\Pr[W \neq g(Y^n)] = \sum_i \lambda_i \Pr[W = i]$
- Arithmetic average probability of error:

$$P_e^{(n)} = \frac{1}{M} \sum_{i=1}^M \lambda_i, \quad P_e^{(n)} \leq \lambda^{(n)}$$

If W is uniformly distributed:

$$P_e^{(n)} = \Pr[W \neq g(Y^n)] \text{ Decoding error probability}$$



Error free \rightarrow 极限情况下 Error probability $\rightarrow 0$.
$$\exists (\bar{T}^{nR}, n), \forall n > 0, \text{integer}$$
$$\lambda^{(n)} \downarrow$$
$$\text{when } n \rightarrow \infty, \lambda^{(n)} \rightarrow 0.$$

- A rate R is **achievable**,

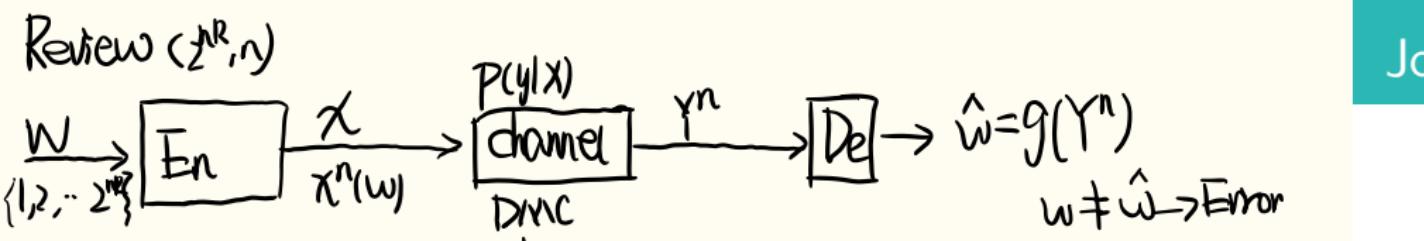
if there exists a sequence of codes with rate R and codeword length n , denoted as $(\lceil 2^{nR} \rceil, n)$, such that the maximal probability of error $\lambda^{(n)} \rightarrow 0$ as $n \rightarrow \infty$.

Recall that

The **rate** R of an (M, n) code is

$$R = \frac{\log M}{n} \text{ bit per transmission.}$$





codebook

$$C^n = \begin{pmatrix} x_c^n(1) \\ x_c^n(2) \\ \vdots \\ x_c^n(2^R) \end{pmatrix} = \begin{pmatrix} x_1(1) & x_2(1) & \cdots & \cdots & x_n(1) \\ x_1(2) & x_2(2) & \cdots & \cdots & x_n(2) \\ \vdots & \vdots & & & \vdots \\ x_1(2^R) & x_2(2^R) & \cdots & \cdots & x_n(2^R) \end{pmatrix}$$

$$\gamma_i = P_r[g(Y^n) = i | X_c^n(i)] \quad X^{(n)} = \max_i \gamma_i$$

R is achievable $\Rightarrow \exists C^1, C^{n+1}, C^{n+2}, \dots, C^{(n)} \xrightarrow{n \rightarrow \infty} 0$.

$$\begin{cases} R < C \Rightarrow R \text{ is achievable} \\ C^{(n)} \rightarrow 0 \Rightarrow R \leq C. \end{cases}$$

Joint Typical Set

- Joint typicality. Given two i.i.d. random variable sequences X^n and Y^n , the set of jointly typical sequences is

$$A_\epsilon^{(n)} = \left\{ (x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n : \begin{array}{l} \left| -\frac{1}{n} \log p(x^n) - H(X) \right| < \epsilon \\ \left| -\frac{1}{n} \log p(y^n) - H(Y) \right| < \epsilon \\ \left| -\frac{1}{n} \log p(x^n, y^n) - H(X, Y) \right| < \epsilon \end{array} \right\}$$

where $p(x^n, y^n) = \prod_{i=1}^n p(x_i, y_i)$.
 \uparrow x_1, x_2, \dots, x_n 之间相互独立
 y_1, y_2, \dots, y_n 之间相互独立
 y_i 由 x_i 决定, 所以 x_i 与 y_i 之间相互独立

Typical set.

$$X^n = x_1, x_2, \dots, x_n \sim P(X^n) = P(x_1)P(x_2) \cdots P(x_n)$$

$$A_\epsilon^n = \left\{ x^n \mid -\frac{1}{n} (\log P(x^n) - H(X)) < \epsilon \right\}$$

$$\Pr[X^n \in A_\epsilon^{(n)}] \rightarrow 1, n \rightarrow \infty$$

$$\text{Joint Typical Set } (x^n, y^n) \sim P(x^n, y^n) = P(x_1, y_1)P(x_2, y_2) \cdots P(x_n, y_n)$$

$$A_\epsilon^{(n)} = \left\{ (x^n, y^n) : \begin{array}{l} \left| -\frac{1}{n} \log p(x^n) - H(X) \right| < \epsilon, \left| -\frac{1}{n} \log p(y^n) - H(Y) \right| < \epsilon, \\ \left| -\frac{1}{n} \log p(x^n, y^n) - H(X, Y) \right| < \epsilon \end{array} \right\}$$



$$p(x^n) \rightarrow \boxed{ch} \rightarrow y^n$$

$$p(y|x)$$

$$p(y^n) = \sum_{y^n} p(x^n, y^n)$$

- Joint AEP Let (X^n, Y^n) be the sequences of length n drawn i.i.d. according to $p(x^n, y^n) = \prod_{i=1}^n p(x_i, y_i)$, then:

$$1. \Pr \left[(X^n, Y^n) \in A_\epsilon^{(n)} \right] \rightarrow 1 \text{ as } n \rightarrow \infty.$$

$$2. |A_\epsilon^{(n)}| \leq 2^{n(H(X,Y)+\epsilon)}$$

$$3. \text{If } (\tilde{X}^n, \tilde{Y}^n) \sim p(x^n)p(y^n), \text{ then}$$

$$\Pr \left[(\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)} \right] \leq 2^{-n(I(X;Y)-3\epsilon)} \xrightarrow{n \rightarrow \infty} \Pr[(\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)}]$$

Please refer to p196 for the proof (proof of Theorem 7.6.1) $\rightarrow \circ$.

$p(x)p(y|x) = p(x,y)$
 $(x^n, y^n) \sim p(x^n, y^n) = \prod_{i=1}^n p(x_i, y_i) \in A_\epsilon^{(n)}$

$p(x^n) \sim x_i^n$ $y_i^n \sim p(y^n)$
 $(x_i^n, y_i^n) \sim p(x^n, y^n)$

$p(x^n) \sim x_2^n$ $y_2^n \sim p(y^n)$
 $(x_2^n, y_2^n) \sim p(x^n, y^n)$

$\Pr[(x_2^n, y_2^n) \in A_\epsilon^{(n)}] \rightarrow 0, n \rightarrow \infty$

Joint Typical set, Joint Typical set
 一组数据与一组数据是否属于
 ①一次

$p(y|x) \Rightarrow C = \max_{p(x)} I(X;Y)$, let $p(x) = \arg \max I(X;Y)$ $p(x,y) = p(x)p(y|x)$
 $\Rightarrow p(x^n, y^n) = \prod_{i=1}^n p(x_i, y_i) \Rightarrow A_{\Sigma}^{(n)}$
Proof: Encoder $\rightarrow p(x)$. i.i.d generate C .
 $\Rightarrow |X|^{2^{nR} \times n}$ codebooks
 Transmit message w with C .
 ① $x_c^n(w)$ ② receive $y^n \sim p[y^n | x^n = x_c^n(w)]$
 ③ g: find the only \hat{w} , $(x_c^n(\hat{w}), y^n) \in A_{\Sigma}^{(n)}$
 Error: ① $(x_c^n(w), y^n) \notin A_{\Sigma}^{(n)}$ ② $(x_c^n(i), y^n) \in A_{\Sigma}^{(n)}$, $i \neq w$.
 let $\mathcal{E}_r = \{\hat{w}(y^n) \neq w\}$ $Pr(\mathcal{E}_r) = \sum_{C \in \mathcal{C}} \Pr(C) \cdot P_e^n(C) \rightarrow$ 使用 C 的概率
 $= \sum_{C \in \mathcal{C}} \Pr(C) \sum_{w=1}^{2^{nR}} \lambda_w(C) \frac{1}{2^{nR}} \rightarrow$ 使用 C 的概率
 $= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_{C \in \mathcal{C}} \Pr(C) \lambda_w(C) = \sum_{C \in \mathcal{C}} \Pr(C) \lambda(C)$

Channel Coding Theorem

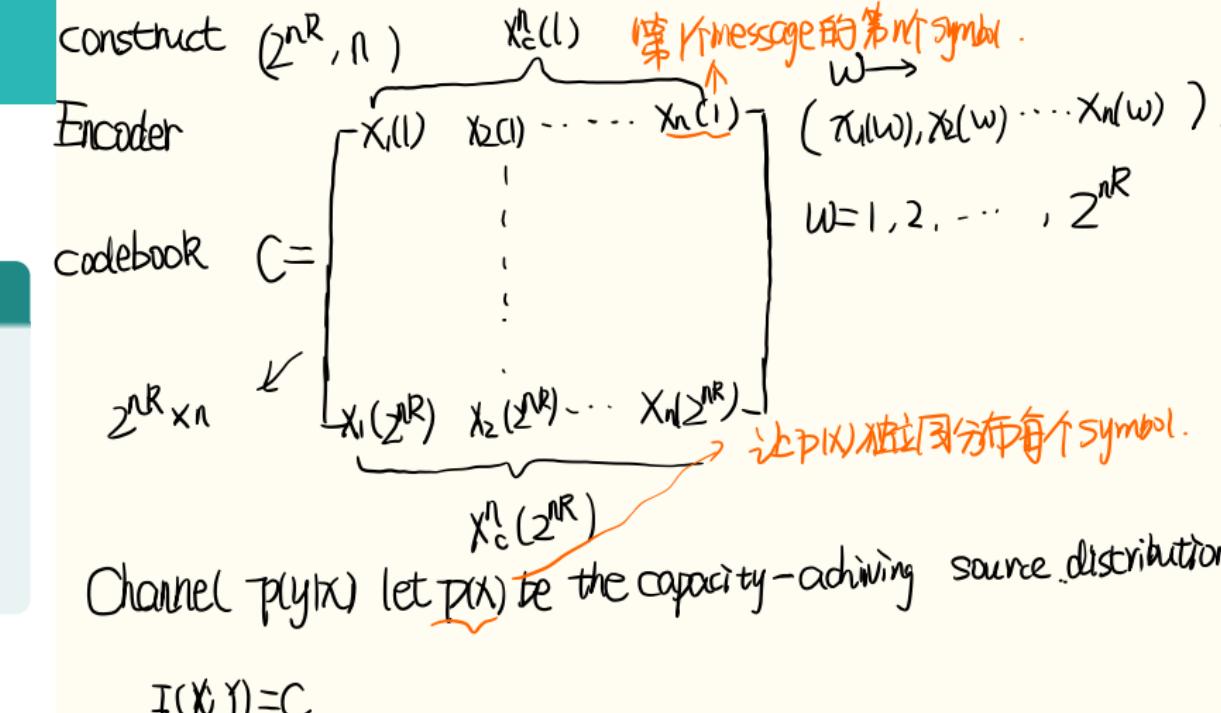
Theorem (Channel coding theorem)

For a discrete memoryless channel, **all rates below capacity C are achievable**. Specifically, for every rate $R < C$, there exists a sequence of $(2^{nR}, n)$ codes with maximum probability of error $\lambda^{(n)} \rightarrow 0$.

Conversely, any sequence of $(2^{nR}, n)$ codes with $\lambda^{(n)} \rightarrow 0$ must have $R \leq C$.
 最大的 Error probability.
 $n \rightarrow \infty, \lambda^{(n)} \rightarrow 0, P_e^{(n)} \rightarrow 0$

Achievability: when $R < C$, there exists zero-error code.

Converse: zero-error codes must have $R \leq C$.



w, i 都是在同一个 codebook 上, 书上会写 $\lambda(w|c) = \lambda_i(c)$.

$$\text{let } e_i(c) = \{(x_c^n(i), y_i^n) \in A_{\Sigma}^{(n)}\}$$

$$\lambda_i(c) = \Pr[\bar{e}_1(c) \cup e_2(c) \cup e_3(c) \dots e_{2^{nR}}(c) | w=i] \\ \leq \Pr[\bar{e}_1(c) | w=i] + \sum_{i=2}^{2^{nR}} \Pr[e_i(c) | w=i]$$

$$\Pr(\Sigma_r) \leq \sum_C \Pr[C] \Pr[\bar{e}_1(c) | w=i] + \sum_{i=2}^{2^{nR}} \sum_C \Pr[C] \Pr[e_i(c) | w=i]$$

$$\sum_C \Pr[C] \Pr[\bar{e}_1(c) | w=i] = \sum_C \prod_{i=1}^{2^{nR}} \Pr[x_c^n(i)] \Pr[\bar{e}_1(c)]$$

$$= \sum_{x_i^n} \Pr[x_i^n] \sum_{c: x_c^n(i) = x_i^n} \prod_{i=1}^{2^{nR}} \Pr[x_c^n(i)] \Pr[\bar{e}_1(c)]$$

$$= \sum_{x_i^n} \Pr[x_i^n] \Pr[\bar{e}_1(c)] \sum_{c: x_c^n(i) = x_i^n} \prod_{i=1}^{2^{nR}} \Pr[x_c^n(i)] = \sum_{x_i^n} \Pr[x_i^n] \Pr[(x_c^n(i), y_i^n) \notin A_{\Sigma}^{(n)}]$$

Random Codebook

- Generate a $(2^{nR}, n)$ code at random according to $p(x)$, where $p(x)$ is the **capacity achieving distribution**. The 2^{nR} are the rows of a matrix:

$$\mathcal{C} = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_n(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(2^{nR}) & x_2(2^{nR}) & \dots & x_n(2^{nR}) \end{bmatrix}.$$

Each entry is generated **i.i.d.** according to $p(x)$.

- Encoding:** map the message $w = \{1, 2, 3, \dots, 2^{nR}\}$ to codeword $[x_1(w), x_2(w), \dots, x_n(w)]$, i.e.

$$\mathcal{C} \rightarrow [x_1(w), x_2(w), \dots, x_n(w)] = x_{\mathcal{C}}^n(w), w = 1, 2, \dots, 2^{nR}$$

- We shall prove the average detection error probability (over all codebooks) tends to zero as n increase, which implies that there must exists one good codebook whose detection error probability tends to zero

Decoder: $w \rightarrow x_i^n(w) \xrightarrow{\text{channel}} y^n \rightarrow \hat{w} = g(y^n)$

$$\bullet P(y^n)$$

$$P(x, y) = P(x) P(y|x)$$

$$P(x^n, y^n) = \prod_{i=1}^n P(x_i, y_i)$$

$$A_{\Sigma}^{(n)}$$

g : find the only \hat{w}

such that $(x_{\hat{w}}^n(\hat{w}), y^n)$

is jointly typical.



$$\begin{aligned}
 &= \sum_{x_c^n(1)} \Pr[x_c^n(1)] \Pr[(x_c^n(1), y_i^n) \notin A_\varepsilon^{(n)}] = \Pr[(x_i^n, y_i^n) \notin A_\varepsilon^{(n)}] \rightarrow 0, n \rightarrow \infty. \\
 &= \Pr[(x_i^n, y_i^n) \notin A_\varepsilon^{(n)}] \leq \varepsilon \text{ for sufficiently large } n.
 \end{aligned}$$

$$\sum_c p_c \Pr[\tilde{e}_i(c) | w=1] \quad i=2 \dots 2^R$$

$$= \Pr[(x_i^n, y_i^n) \notin A_\varepsilon^{(n)} | w=1] \leq 2^{-n[I(X;Y)-3\varepsilon]}$$

$$\Pr[\Sigma_r] \leq \sum_{i=2}^{2^R} 2^{-n[I(X;Y)-3\varepsilon]} = \sum 2^{-n[I(X;Y)-3\varepsilon]} (2^R - 1)$$

$$< \sum 2^{-n[I(X;Y)-3\varepsilon]} 2^R = \sum 2^{-n[I(X;Y)-R-3\varepsilon]} = \sum 2^{-n[C-R-3\varepsilon]}$$

If $C-R-3\varepsilon > 0$. $\Pr[\Sigma_r] \rightarrow 0$. \Rightarrow If $R < C-3\varepsilon$, $\Pr[\Sigma_r] \rightarrow 0$

\Rightarrow If $R < C$, $\Pr[\Sigma_r] \rightarrow 0$.

Jointly Typical Decoding

- **Decoding:** finds the only \hat{w} such that $(x_c^n(\hat{w}), Y_c^n)$ is jointly typical.
- **Decoding error:** Suppose message 1 is sent to via codeword $x_c^n(1)$ and Y_c^n is the received signal, the possible decoding error events include:
 - $(x_c^n(1), Y_c^n)$ is not joint typical.
 - $(x_c^n(i), Y_c^n)$ is joint typical ($i = 2, 3, \dots, 2^R$).
- **Idea of proof:** According to joint AEP, since $x_c^n(1)$ and Y_c^n are generated according to joint distribution $p(x^n, y^n)$, the chance of the first event is small. Moreover, since Y_c^n is generated independently of $x_c^n(i)$, the total chance of the second event is also small.

(2^{nR}, n)

W → Xⁿ(W) → Yⁿ → \hat{W} 构成了 Markov chain

message. codeword. $\Pr[W \neq \hat{W}] = Pe^{(n)}$ 仍然是 Markov chain, 可用 Fano's inequality

$\{1, \dots, 2^{nR}\}$ 均匀分布时最大.

$H(W|\hat{W}) \leq H(Pe^{(n)}) + Pe^{(n)} \log(2^{nR}-1) \leq 1 + nRPe^{(n)}$

$H(W|\hat{W}) = H(W) - I(W; \hat{W}) = \log_2 2^{nR} - I(W; \hat{W})$

假设 W 的选择.

$I(X^n; Y^n) = H(Y^n) - H(Y^n|X^n) = H(Y^n) - \sum_{i=1}^n H(Y_i|X_i) \leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i|X_i)$

仔细理解

$= \sum_{i=1}^n H(Y_i) - H(X_i|Y_i) = \sum_{i=1}^n I(X_i; Y_i) \leq nC.$

$\rightarrow H(W|\hat{W}) \geq nR - I(X^n; Y^n) \geq n(R-C) \rightarrow n(R-C) \leq 1 + nRPe^{(n)} \rightarrow R \leq \frac{1}{n(1-Pe^{(n)})} + \frac{C}{nPe^{(n)}}$

$n \rightarrow \infty, \lambda^{(n)} \rightarrow 0, Pe^{(n)} \rightarrow 0, \frac{1}{n(1-Pe^{(n)})} + \frac{C}{nPe^{(n)}} \rightarrow C. R \leq C$

↓ Proof for achievability

- A message W is chosen according to a uniform distribution

$$\Pr[W = w] = 2^{-nR},$$
 for $w = 1, 2, \dots, 2^{nR}$. The w -th codeword $x_{\mathcal{C}}^n(w)$, corresponding to the w -th row of \mathcal{C} , is sent over the channel.
- The receiver receives a sequence $Y_{\mathcal{C}}^n$ according to the distribution according to the distribution

$$\Pr(y_{\mathcal{C}}^n | x_{\mathcal{C}}^n(w)) = \prod_{i=1}^n \Pr(y_{i,\mathcal{C}} | x_{i,\mathcal{C}}(w)),$$
 and guesses which message was sent using **jointly typical decoding**.

 南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

Proof for achievability

- Let $\varepsilon = \{\hat{W}(Y^n) \neq W\}$ denote the error event, $\lambda_w(\mathcal{C})$ be the error probability of the w -th codeword of code \mathcal{C} . The **average probability of error**, over all codewords and all codebooks, is:

$$\begin{aligned}\Pr(\varepsilon) &= \sum_{\mathcal{C}} \Pr(\mathcal{C}) P_e^{(n)}(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(\mathcal{C}) \\ &= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_w(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_1(\mathcal{C}),\end{aligned}$$

where $\sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_1(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_w(\mathcal{C})$, $\forall w \neq 1$.

Proof for achievability

- Let $Y_{\mathcal{C}}^n$ be the received signal for $x_{\mathcal{C}}^n(1)$

$$e_i(\mathcal{C}) = \{(x_{\mathcal{C}}^n(i), Y_{\mathcal{C}}^n) \in A_{\epsilon}^{(n)}\}, i \in \{1, 2, \dots, 2^{nR}\},$$

and $e_i^c(\mathcal{C}) = \neg e_i(\mathcal{C})$. Thus,

$$\begin{aligned}\Pr[\varepsilon] &= \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_1(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr \left[e_1^c(\mathcal{C}) \cup (\bigcup_{i=2}^{2^{nR}} e_i(\mathcal{C})) \middle| W = 1 \right] \\ &\leq \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_1^c(\mathcal{C}) | W = 1] + \sum_{\mathcal{C}} \Pr(\mathcal{C}) \sum_{i=2}^{2^{nR}} \Pr[e_i(\mathcal{C}) | W = 1] \\ &= \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_1^c(\mathcal{C}) | W = 1] + \sum_{i=2}^{2^{nR}} \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_i(\mathcal{C}) | W = 1]\end{aligned}$$



Proof for achievability

$$\begin{aligned} & \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_1^c(\mathcal{C})|W = 1] \\ &= \sum_{\mathcal{C}} \left(\prod_{i=1}^{2^{nR}} \Pr(x_{\mathcal{C}}^n(i)) \right) \Pr[e_1^c(\mathcal{C})|W = 1] \\ &= \sum_{x_1^n} \sum_{\mathcal{C}: x_{\mathcal{C}}^n(1) = x_1^n} \prod_{i=1}^{2^{nR}} \Pr(x_{\mathcal{C}}^n(i)) \Pr(x_1^n \text{ and } Y^n \text{ are not joint typical} | W = 1) \\ &= \sum_{x_1^n} \Pr(x_1^n) \Pr(x_1^n \text{ and } Y^n \text{ are not joint typical} | W = 1) \\ &\quad \times \sum_{\mathcal{C}: x_{\mathcal{C}}^n(1) = x_1^n} \prod_{i=2}^{2^{nR}} \Pr(x_{\mathcal{C}}^n(i)) \\ &= \sum_{x_1^n} \Pr(x_1^n) \Pr(x_1^n \text{ and } Y^n \text{ are not joint typical} | W = 1) \\ &= \Pr(X_1^n \text{ and } Y^n \text{ are not joint typical} | W = 1) = \Pr(E_1^c | W = 1) \end{aligned}$$



Proof for achievability

- Similarly,

$$\begin{aligned}\sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_1(\mathcal{C})|W = 1] &= \Pr(\textcolor{blue}{X_i^n \text{ and } Y^n \text{ are joint typical}}|W = 1) \\ &= \Pr(E_i|W = 1)\end{aligned}$$

- As a result,

$$\Pr[\varepsilon] \leq \Pr[E_1^c|W = 1] + \sum_{i=2}^{2^{nR}} \Pr[E_i|W = 1]$$

Proof for achievability

- By the joint AEP, $\Pr[E_1^c | W = 1] \leq \epsilon$ for n sufficiently large. By the code generation process, $X^n(1)$ and $X^n(i)$ are independent for $i \neq 1$, so are Y^n and $X^n(i)$. Hence the probability that $X^n(i)$ and Y^n are jointly typical is $\leq 2^{-n(I(X;Y)-3\epsilon)}$ by the joint AEP.

$$\begin{aligned}\Pr[\varepsilon] &\leq \epsilon + \sum_{i=2}^{2^{nR}} 2^{-n(I(X;Y)-3\epsilon)} \\ &= \epsilon + (2^{nR} - 1)2^{-n(I(X;Y)-3\epsilon)} \\ &\leq \epsilon + 2^{3n\epsilon} 2^{-n(I(X;Y)-R)} \\ &\leq 2\epsilon \quad \text{for } R \leq I(X;Y) - 4\epsilon \text{ and sufficiently large } n\end{aligned}$$

Hence, if $R < I(X;Y)$, we can choose ϵ and n so that the average probability of error, over codebooks and codewords, is less than 2ϵ .

- Since $p(x)$ is the capacity achieving distribution, $R < I(X;Y)$ becomes $R < C$.



Proof for achievability

- Get rid of the average over codebooks. Since the average probability of error is $\leq 2\epsilon$, there exists **at least one** codebook \mathcal{C}^* with a small average probability of error ($\Pr(\varepsilon|\mathcal{C}^*) \leq 2\epsilon$). Since we have chosen \hat{W} according to a uniform distribution, we have

$$\Pr(\varepsilon|\mathcal{C}^*) = \frac{1}{2^{nR}} \sum_{i=1}^{2^{nR}} \lambda_i(\mathcal{C}^*).$$

- Throw away the worst half of the codewords in the best codebook \mathcal{C}^* . We have $\Pr(\varepsilon|\mathcal{C}^*) \leq \frac{1}{2^{nR}} \sum \lambda_i(\mathcal{C}^*) \leq 2\epsilon$. This implies that **at least half** the indices i and their associated codewords $X^n(I)$ must have conditional probability of error $\lambda_i \leq 4\epsilon$. If we reindex the codewords, we have 2^{nR-1} codewords. The rate now is $R' = R - \frac{1}{n}$ with maximal probability of error $\lambda^{(n)} \leq 4\epsilon$.



Proof for the converse

- The index W is uniformly distributed on the set $\mathcal{W} = \{1, 2, \dots, 2^{nR}\}$, and the sequence Y^n is related to W . From Y^n , we estimate the index W as $\hat{W} = g(Y^n)$. Thus, $W \rightarrow X^n(W) \rightarrow Y^n \rightarrow \hat{W}$ forms a Markov chain.

Data processing inequality: $I(W; \hat{W}) \leq I(X^n(W); Y^n)$

Lemma (Fano's inequality)

For a discrete memoryless channel with a codebook \mathcal{C} and the input message W uniformly distributed over 2^{nR} , we have

$$H(W|\hat{W}) \leq 1 + P_e^{(n)}nR.$$



Proof for the converse

- The index W is uniformly distributed on the set $\mathcal{W} = \{1, 2, \dots, 2^{nR}\}$, and the sequence Y^n is related to W . From Y^n , we estimate the index W as $\hat{W} = g(Y^n)$. Thus, $W \rightarrow X^n(W) \rightarrow Y^n \rightarrow \hat{W}$ forms a Markov chain.

Data processing inequality: $I(W; \hat{W}) \leq I(X^n(W); Y^n)$

Lemma (Fano's inequality)

For a discrete memoryless channel with a codebook \mathcal{C} and the input message W uniformly distributed over 2^{nR} , we have

$$H(W|\hat{W}) \leq 1 + P_e^{(n)}nR.$$



Lemma

Let Y^n be the result of passing X^n through a discrete memoryless channel of capacity C . Then

$$I(X^n; Y^n) \leq nC, \quad \text{for all } p(x^n).$$

Proof.

$$\begin{aligned} I(X^n; Y^n) &= H(Y^n) - H(Y^n | X^n) = H(Y^n) - \sum_{i=1}^n H(Y_i | Y_1, \dots, Y_{i-1}, X^n) \\ &= H(Y^n) - \sum_{i=1}^n H(Y_i | X_i) \quad \text{memoryless} \\ &\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \quad \text{independence bound} \\ &= \sum_{i=1}^n I(X_i | Y_i) \leq nC \end{aligned}$$

Lemma

Let Y^n be the result of passing X^n through a discrete memoryless channel of capacity C . Then

$$I(X^n; Y^n) \leq nC, \quad \text{for all } p(x^n).$$

Proof.

$$\begin{aligned} I(X^n; Y^n) &= H(Y^n) - H(Y^n | X^n) = H(Y^n) - \sum_{i=1}^n H(Y_i | Y_1, \dots, Y_{i-1}, X^n) \\ &= H(Y^n) - \sum_{i=1}^n H(Y_i | X_i) \quad \text{memoryless} \\ &\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \quad \text{independence bound} \\ &= \sum_{i=1}^n I(X_i | Y_i) \leq nC \end{aligned}$$

Proof for the converse

Proof.

Converse to channel coding theorem: Since W has a uniform distribution, we have

$$\begin{aligned} nR &= H(W) = H(W|\hat{W}) + I(W;\hat{W}) \\ &\leq 1 + P_e^{(n)}nR + I(W;\hat{W}) \quad \text{Fano's inequality} \\ &\leq 1 + P_e^{(n)}nR + I(X^n;Y^n) \quad \text{data-processing inequality} \\ &\leq 1 + P_e^{(n)}nR + nC \quad \text{Lemma 7.9.2} \end{aligned}$$

We obtain $R \leq \frac{1}{n(1+P_e^{(n)})} + \frac{C}{1+P_e^{(n)}} \rightarrow \frac{1}{n} + C$.

Letting $n \rightarrow \infty$, we have $R \leq C$.



Reading & Homework

- **Reading:** Chapter 7: 7.6-7.10
- **Homework:** Problems 7.15, 7.31.