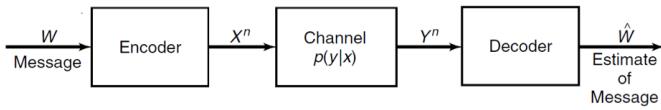
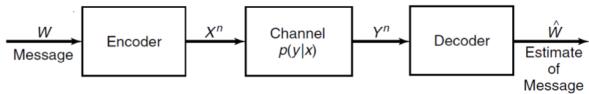


# Communication System Model



- $X^n = [X_1, X_2, \dots, X_n]$
- $Y^n = [Y_1, Y_2, \dots, Y_n]$
- Channel  $p(y^n|x^n)$ : probability of observing  $y^n$  given input sequence  $x^n$

## Discrete memoryless channel (DMC)



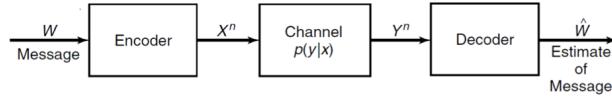
### Definition

A **discrete channel** consists of an input alphabet  $\mathcal{X}$  and output alphabet  $\mathcal{Y}$  and a probability transition matrix  $p(y^n|x^n)$  that expresses the probability of observing the output sequence  $y^n$  given that we send the sequence  $x^n$ .

### Definition

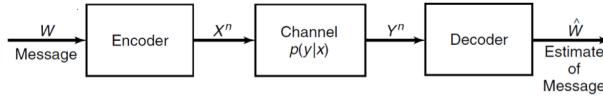
The channel is called **memoryless** if  $p(y^n|x^n) = \prod_{i=1}^n p(y_i|x_i)$ .

# Communication System Model

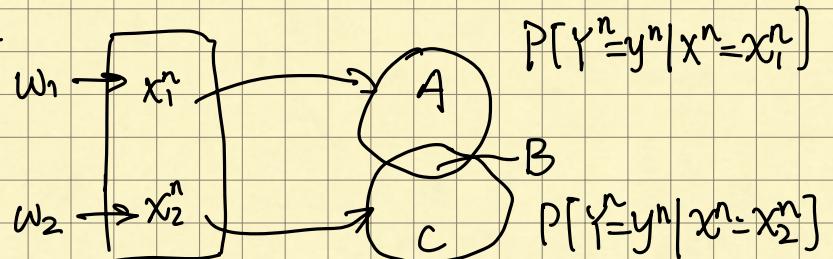


- $X^n = [X_1, X_2, \dots, X_n] \in \mathcal{X}^n$ ,  $Y^n = [Y_1, Y_2, \dots, Y_n] \in \mathcal{Y}^n$
- Channel  $p(y^n|x^n)$ : probability of observing  $y^n$  given input sequence  $x^n$
- **Memoryless**:  $p(y^n|x^n) = \prod_{i=1}^n p(y_i|x_i)$
- Messages are mapped into some sequence of the channel symbols. Output sequence is random but **has a distribution that depends on the input sequences**. Each possible input sequence may induce several possible outputs, and hence inputs are **confusable**. Can we choose a **non-confusable** subset of input sequences?

## Duality



- **Data compression**: we **remove** all the redundancy in the data to form the most compressed version possible.
- **Data transmission**: we **add** redundancy in a controlled manner to combat errors in the channel.



Region A:  $\hat{W}=W_1$ , Region C:  $\hat{W}=W_2$

Region B:  $\hat{W}=W_1 \cup W_2$ ?

要尽力使 region B 变小, 以减少错误。  
n 越大, region B 越小, 但 n 不能过大。

# "Survivor"

- You were deserted on a small island. You met a native and asked about the weather.

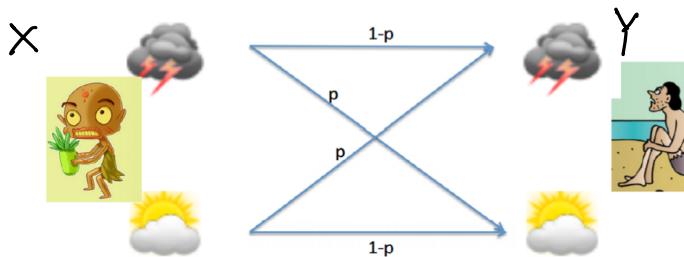
- True weather is a random variable  $X$

$$X = \begin{cases} \text{rain} & \text{w.p. } \alpha, \\ \text{sunny} & \text{w.p. } 1 - \alpha, \end{cases}$$

- Native knows tomorrow's weather perfectly, but only tells truth with probability  $1 - p$ .

- Native's answer is a random variable  $Y \in \{\text{rain, sunny}\}$ .

- How informative is the native's answer?



What is  $I(X; Y)$ ?

- $I(X; Y) = H(X) - H(X|Y)$
- $H(X) = H(\alpha) = -\alpha \log \alpha - (1 - \alpha) \log(1 - \alpha)$
- $H(X|Y) = H(X|Y = \text{rain})p(\text{rain}) + H(X|Y = \text{sunny})p(\text{sunny})$
- $H(X|Y = \text{rain})$  is equal to  
 $-\sum_{i \in \{\text{rain, sunny}\}} p(X = i|Y = \text{rain}) \log p(X = i|Y = \text{rain})$ . Note that

$$p(X = \text{rain}|Y = \text{rain}) = \frac{p(X = \text{rain}|Y = \text{rain})p(X = \text{rain})}{p(Y = \text{rain})} = \frac{(1-p)\alpha}{(1-p)\alpha + p(1-\alpha)}$$

$$\text{Thus, } H(X|Y) = \alpha H\left(\frac{(1-p)\alpha}{(1-p)\alpha + p(1-\alpha)}\right) + (1 - \alpha) H\left(\frac{p\alpha}{p\alpha + (1-p)(1-\alpha)}\right)$$

- $I(X; Y) = H(\alpha) - \alpha H\left(\frac{(1-p)\alpha}{(1-p)\alpha + p(1-\alpha)}\right) - (1 - \alpha) H\left(\frac{p\alpha}{p\alpha + (1-p)(1-\alpha)}\right)$

## Special Cases

- $I(X; Y) = H(\alpha) - \alpha H\left(\frac{(1-p)\alpha}{(1-p)\alpha + p(1-\alpha)}\right) - (1 - \alpha) H\left(\frac{p\alpha}{p\alpha + (1-p)(1-\alpha)}\right)$

- Always telling the truth:  $p = 0$

$$I(X; Y) = H(\alpha) - \alpha H(1) - (1 - \alpha) H(0) = H(\alpha) \leq 1 \text{ bit}$$

- Telling truth half of the time:  $p = 1/2$

$$I(X; Y) = H(\alpha) - \alpha H(\alpha) - (1 - \alpha) H(\alpha) = 0 \text{ bit}$$

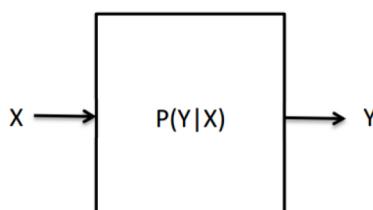
- Fix  $p$ , maximize with respect to  $\alpha$ , maximum achieved when  $\alpha = 1/2$

$$\max_{\alpha} I(X; Y) = H(1/2) - \frac{1}{2}H(1-p) - \frac{1}{2}H(p) = 1 - H(P)$$

## Information Channel Capacity

Definition ("Information" Channel Capacity)

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



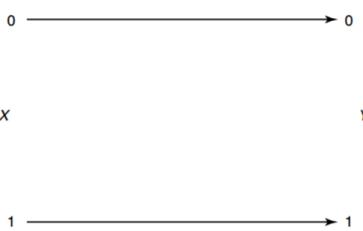
$$P(X, Y) = P(X)P(Y|X)$$

↓  
Source channel

$I(X; Y)$  表示  $X$  与  $Y$  之间传递的信息量。  
 $X$  是 source 的分布,  $P$  是 channel 的分布。

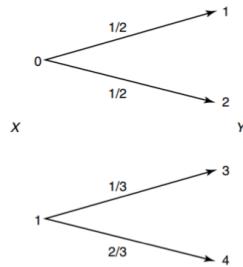
# Examples

- Binary noiseless channel



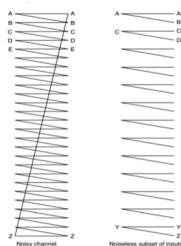
$$C = \max I(X; Y) = \log 2 = 1 \text{ bits} \quad (\text{with } p(x) = (\frac{1}{2}, \frac{1}{2}))$$

- Noisy channel with nonoverlapping outputs



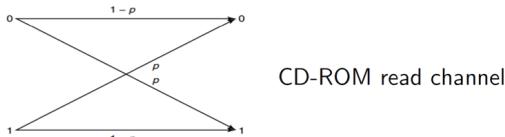
$$C = \max I(X; Y) = \log 2 = 1 \text{ bits} \quad (\text{with } p(x) = (\frac{1}{2}, \frac{1}{2}))$$

- Noisy typewriter



$$C = \max I(X; Y) = \log \frac{26}{2} = \log 13 \text{ bits} \quad (\text{with } p(x) \text{ uniformly distributed})$$

- Binary symmetric channel



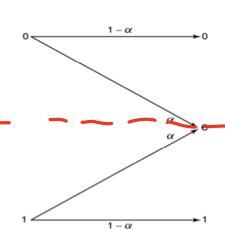
$$I(X; Y) = H(Y) - H(Y|X) = H(Y) - \sum_{x \in \{0,1\}} p(x)H(Y|X=x)$$

$$= H(Y) - \sum_{x \in \{0,1\}} p(x)H(p) = H(Y) - H(p) \leq 1 - H(p)$$

$$C = \max I(X; Y) = I - H(p) \text{ bits}$$

- Binary erasure channel

$$\begin{aligned} C &= \max_{p(x)} I(X; Y) \\ &= \max_{p(x)} (H(Y) - H(Y|X)) \\ &= \max_{p(x)} H(Y) - H(\alpha) \end{aligned}$$



Let  $\Pr[X = 1] = \pi$ , then  $(1-\pi)\alpha + \pi\alpha$

$$H(Y) = H((1-\pi)(1-\alpha), \alpha, \pi(1-\alpha)) \quad !! = H(\alpha) + (1-\alpha)H(\pi)$$

Thus,  $C = \max_{\pi} (1-\alpha)H(\pi) = 1 - \alpha$  (with  $\pi = \frac{1}{2}$ )

无随机性,  $H(Y|X) = 0$ .

$$I(X; Y) = H(Y) - H(Y|X) = H(Y)$$

$$C = \max I(X; Y) = \max H(Y) = 1$$

无重复输出, 与上者同.

$$I(X; Y) = H(X) - H(X|Y) \quad !!$$

$$C = \max H(X) = 1$$

$$\begin{aligned} C &= \max H(Y) - H(Y|X), \text{ 给定 } X, Y \text{ 的取值只有两个,} \\ &= \max H(Y) - \log 2, \text{ 且等可能.} \\ &= \log 26 - \log 2 \end{aligned}$$

$$(x \mid X=A \text{ 且 } Y=\begin{cases} A, & P=\frac{1}{2} \\ B, & P=\frac{1}{2} \end{cases})$$

# Symmetric channel

$$p(y|x) = \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.5 & 0.3 & 0.2 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}.$$

All the rows of the transition matrix are permutations of each other and so are the columns. Let  $r$  be a row of the transition matrix.

$$I(X;Y) = H(Y) - H(Y|X) = H(Y) - H(r) \leq \log |\mathcal{Y}| - H(r)$$

with equality if  $\mathcal{Y}$  is uniformly distributed. If  $p(x) = \frac{1}{|\mathcal{X}|}$ ,  $Y$  is also uniformly distributed:

$$p(y) = \sum_{x \in \mathcal{X}} p(y|x)p(x) = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} p(y|x) = \frac{c}{|\mathcal{X}|} = \frac{1}{|\mathcal{Y}|},$$

where  $c$  is the sum of the entries in one column.



## Fundamental question

- How fast can we transmit information over a channel?
- Suppose a source sends  $r$  messages per second, and the entropy of a message is  $H$  bits per message, information rate is  $R = rH$  bits/second.
- Intuition: as  $R$  increases, error will increase.
- Surprisingly, Shannon showed error can approach to zero, as long as

$$R < C$$

## Review

- **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)

## Channel Code

### 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

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

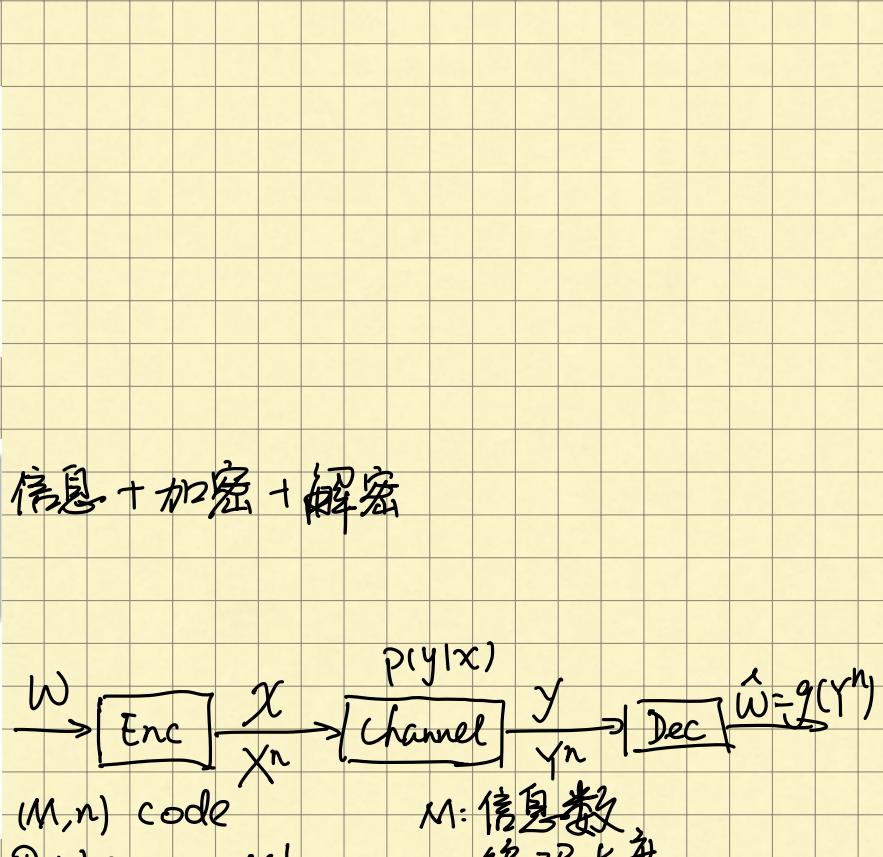
### Data Rate

发送一次能传输的bit数

$$R = \frac{\log M}{n} \rightarrow M \text{ 个可能的信息编码} \quad \text{D-ary码需要 } \log_2 M \text{ 个 codeword}$$

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

$$\begin{aligned} C &= \max H(Y) - H(Y|X) \\ &= \max H(Y) - H(0.2, 0.3, 0.5) \\ &= \log 3 - H(0.2, 0.3, 0.5) \end{aligned}$$



# Performance Metric

- Conditional probability of error:

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

给定第*i*个信息的编码  
解码出来不为第*i*个信息

- 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:

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

## Achievable Rate

- A rate  $R$  is achievable,

if there exists a sequence of codes with rate  $R$  and codeword length  $n$ , denoted as  $([2^{nR}], n)$ , 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

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

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

## Joint AEP

- 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:

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

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

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

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

Please refer to p196 for the proof (proof of Theorem 7.6.1)

$\lambda_i$  即为解码错误率的概率

$\lambda^{(n)}$  遍近于0，则出错概率低

允许出错，但 codeword 无限长时，最差出错概率为0。

typical Set:  $(x_1, \dots, x_n) \sim p(x) \triangleq X^n$

$$A_\epsilon^{(n)} = \{x^n : |-\frac{1}{n} \log p(x) - H(X)| < \epsilon\}$$

joint typical set:

$$((x_1, y_1), \dots, (x_n, y_n)) \sim p(x, y) \triangleq (X^n, Y^n)$$

$$A_\epsilon^{(n)} = \{(x^n, y^n) : |-\frac{1}{n} \log p(x^n) - H(X)| < \epsilon; |-\frac{1}{n} \log p(y^n) - H(Y)| < \epsilon\} \leftarrow X^n$$

$$|-\frac{1}{n} \log p(y^n) - H(Y)| < \epsilon; |-\frac{1}{n} \log p(x^n, y^n) - H(X, Y)| < \epsilon \leftarrow (X^n, Y^n)$$

$$p(x, y) \sim (X^n, Y^n) \quad X \text{与 } Y \text{ 有相关性}$$

$$p(x)p(y) \sim (\tilde{X}^n, \tilde{Y}^n) \quad \tilde{X} \text{与 } \tilde{Y} \text{ 无关, 但 } X \text{ 与 } \tilde{X} \text{ 同分布} \quad Y \text{ 与 } \tilde{Y} \text{ 同分布}$$

$n \rightarrow \infty$ , 按  $p(x, y)$  产生的  $X$  与  $Y$  一定在 typical set 上  
但按  $p(x)p(y)$  产生的  $\tilde{X}$  与  $\tilde{Y}$  不在 typical set 上

# 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  $R \leq C$ .

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

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

## 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:

$$\mathcal{C} = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_n(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(2^{nR}) & x_2(2^{nR}) & \dots & x_n(2^{nR}) \end{bmatrix} \xrightarrow{\text{rows}} X_C^n(1) \xrightarrow{\text{rows}} X_C^n(2^{nR})$$

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



## 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 ( $i = 2, 3, \dots, 2^{nR}$ ).
- 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.

## 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  $x_C^n(w)$ , corresponding to the  $w$ -th row of  $\mathcal{C}$ , is sent over the channel.

- The receiver receives a sequence  $Y_C^n$  according to the distribution according to the distribution

$$\Pr(Y_C^n | x_C^n(w)) = \prod_{i=1}^n \Pr(y_{i,C} | x_{i,C}(w)),$$

and guesses which message was sent using jointly typical decoding.

①  $R < C \Rightarrow R$  is achievable,  $\lim_{n \rightarrow \infty} \lambda^{(n)} \rightarrow 0$

②  $\lambda^{(n)} \rightarrow 0 \Rightarrow R \leq C$

证 ①  $R < C \Rightarrow R$  is achievable,  $\lim_{n \rightarrow \infty} \lambda^{(n)} \rightarrow 0$

$C$  是  $2^{nR} \times n$  的矩阵, 每行都是一个 message 的 Codeword

每一项都由  $p(x)$  独立同分布产生

$$W \rightarrow X_C^n(W) \rightarrow Y^n \rightarrow \hat{W} = g(Y^n)$$

$$\uparrow p(x) \quad \uparrow p(y|x) \quad \uparrow p(x)p(y|x) = p(x,y) \Rightarrow \text{产生 } A_\epsilon^{(n)}$$

因为收到  $Y^n$  后,  $g$  要寻找  $\hat{w}$ , 使得  $(X_C^n(\hat{w}), Y^n)$  在 joint typical set 中。

成功的情况:

$$(X_C^n(i), Y^n) \in A_\epsilon^{(n)} \text{ 且 } i = \hat{w}, \text{ 并且唯一}$$

出错的情况:

$$\textcircled{1} (X_C^n(\hat{w}), Y^n) \notin A_\epsilon^{(n)}$$

$$\textcircled{2} (X_C^n(i), Y^n) \in A_\epsilon^{(n)}, i \neq \hat{w}$$

由 joint AEP 的第三条性质, ② 出现的概率率  $\rightarrow 0$

假设有一种传输: virtual transmission

① uniformly pick up one message

② randomly generate a codebook  $C$

③ Joint typical decoding.  $A_\epsilon^{(n)} \sim p(x)p(y|x) = p(x,y)$

若这种传输能有当  $n \rightarrow \infty$ ,  $\lambda^{(n)} \rightarrow 0$ . 那么必然存在 codebook 满足  $\lambda^{(n)} \rightarrow 0$  当  $n \rightarrow \infty$ .

- 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:

$$\Pr(\varepsilon) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) P_e^{(n)}(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(\mathcal{C})$$

$$= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_w(\mathcal{C}) = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_1(\mathcal{C}),$$

w出错概率

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

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

$$\text{错误事件 } e_i(\mathcal{C}) = \{(x_{\mathcal{C}}^n(i), Y_{\mathcal{C}}^n) \in A_{\varepsilon}^{(n)}\}, i \in \{1, 2, \dots, 2^{nR}\},$$

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

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

需要证

→ 0 → 0

由于  $e_i(\mathcal{C})$  中每一行的生成都是  $i$  的，故

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

由上可知

- Similarly,

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

- As a result,

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

- By the joint AEP,  $\Pr[E_1^c | W=1] \leq \varepsilon$  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  $\leq 2^{-n(I(X;Y)-3\varepsilon)}$  by the joint AEP.

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

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

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

$$\Pr[\hat{W} \neq W] = \Pr[Er] = \sum_{\mathcal{C}} \Pr(\mathcal{C}) \lambda_1(\mathcal{C})$$

$$\begin{aligned} e_i(\mathcal{C}) &= \{(x_{\mathcal{C}}^n(i), Y_{\mathcal{C}}^n) \in A_{\varepsilon}^{(n)}\}, i=1, 2, \dots, 2^{nR}. \\ \lambda_1(\mathcal{C}) &= \Pr[\bar{e}_1(\mathcal{C}) \cup e_2(\mathcal{C}) \cup \dots \cup e_{2^{nR}}(\mathcal{C}) | W=1] \\ &\leq \Pr[\bar{e}_1(\mathcal{C}) | W=1] + \sum_{i=2}^{2^{nR}} \Pr[e_i(\mathcal{C}) | W=1] \\ \Pr[Er] &\leq \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[\bar{e}_1(\mathcal{C}) | W=1] + \sum_{i=2}^{2^{nR}} \sum_{\mathcal{C}} \Pr(\mathcal{C}) \Pr[e_i(\mathcal{C}) | W=1] \\ &\leq \varepsilon \end{aligned}$$

$\leq 2^{-n(I(X;Y)-3\varepsilon)}$

由于  $C$  中每一行的生成都是  $i$  的，故

$$\Pr[\mathcal{C}] = \Pr[x_{\mathcal{C}}^n(1)] \Pr[x_{\mathcal{C}}^n(2)] \dots \Pr[x_{\mathcal{C}}^n(2^{nR})]$$

$$\sum_{\mathcal{C}} = \sum_{x_1^n} \sum_{C: x_{\mathcal{C}}^n(1)=x_1^n} \Pr[C], \Pr[\mathcal{C}] = \prod_{i=1}^{2^{nR}} \Pr[x_{\mathcal{C}}^n(i)]$$

$$\sum_{\mathcal{C}} \Pr[\mathcal{C}] \Pr[\bar{e}_1(\mathcal{C}) | W=1]$$

$$= \sum_{\mathcal{C}} \prod_{i=1}^{2^{nR}} \Pr[x_{\mathcal{C}}^n(i)] \Pr[\bar{e}_1(\mathcal{C}) | W=1]$$

$$= \sum_{x_1^n} \sum_{C: x_{\mathcal{C}}^n(1)=x_1^n} \Pr[x_{\mathcal{C}}^n(1)] \prod_{i=2}^{2^{nR}} \Pr[x_{\mathcal{C}}^n(i)] \Pr[\bar{e}_1(\mathcal{C}) | W=1]$$

$$= \sum_{x_1^n} \Pr[x_{\mathcal{C}}^n(1)] \Pr[\bar{e}_1(\mathcal{C}) | W=1] \sum_{C: x_{\mathcal{C}}^n(1)=x_1^n} \prod_{i=2}^{2^{nR}} \Pr[x_{\mathcal{C}}^n(i)]$$

$$\sum_{C: x_{\mathcal{C}}^n(1)=x_1^n} = \sum_{x_2^n, x_3^n, \dots, x_{2^{nR}}^n} \text{只固定了第一行，后面均随机选取.}$$

$$\Rightarrow \sum_{C: x_{\mathcal{C}}^n(1)=x_1^n} \prod_{i=2}^{2^{nR}} \Pr[x_{\mathcal{C}}^n(i)] = \sum_{x_2^n} \Pr[x_{\mathcal{C}}^n(2)] \sum_{x_3^n} \Pr[x_{\mathcal{C}}^n(3)] \dots \sum_{x_{2^{nR}}^n} \Pr[x_{\mathcal{C}}^n(2^{nR})] = 1$$

$$\Pr[\bar{e}_1] = \sum_{x_1^n} \Pr[x_1^n] \Pr[(x_1^n, Y_1^n) \notin A_{\varepsilon}^{(n)} | W=1]$$

$$= \Pr[(x_1^n, Y_1^n) \notin A_{\varepsilon}^{(n)} | W=1] \leq \varepsilon \text{ for sufficiently large } n.$$

$$\Pr[er] \leq \varepsilon + \sum_{i=2}^{2^{nR}} 2^{-n(I(X;Y)-3\varepsilon)}$$

$$< \varepsilon + 2^{nR} \cdot 2^{-n(I(X;Y)-3\varepsilon)}$$

$$= \varepsilon + 2^{-n(I(X;Y)-R-3\varepsilon)} \geq \varepsilon$$

若  $R \leq I(X;Y) - 4\varepsilon$ , 则有

$$\leq \varepsilon + 2^{-n\varepsilon} \leq 2\varepsilon \text{ for sufficiently large } n.$$

- 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 ( $\Pr(\epsilon|\mathcal{C}^*) \leq 2\epsilon$ ). Since we have chosen  $\hat{W}$  according to a uniform distribution, we have

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

- Throw away the worst half of the codewords in the best codebook  $\mathcal{C}^*$ . We have  $\Pr(\epsilon|\mathcal{C}^*) \leq \frac{1}{2^{nR}} \sum \lambda_i(\mathcal{C}^*) \leq 2\epsilon$ . 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  $\lambda^{(n)} \leq 4\epsilon$ .

所以对于足够的  $n$ ,  $R < C = \max I(X; Y)$ , 因为

$$\Pr[\text{Er}] \leq 2\epsilon$$

也即  $\exists \mathcal{C}^*$ ,  $\Pr[\text{Er}(\mathcal{C}^*)] \leq \epsilon$ , for  $R < C$ , and sufficiently large  $n$ .

$$\Pr[\text{Er}(\mathcal{C}^*)] = \frac{1}{2^{nR}} \sum_{i=1}^{2^{nR}} \lambda_i(\mathcal{C}^*) \leq 2\epsilon$$

不失一般性, 设  $\lambda_1(\mathcal{C}^*) \leq \lambda_2(\mathcal{C}^*) \leq \dots \leq \lambda_{2^{nR}}(\mathcal{C}^*)$   
将  $i > 2^{nR-1}$  (即纠错概率较大的那一半) 破去  
对剩余的  $\tilde{\mathcal{C}}^*$ , 有 codeword  $2^{nR-1}$ . 此时的 rate 为

$$R' = \frac{\log 2^{nR-1}}{n} = R - \frac{1}{n} \rightarrow R, n \rightarrow \infty$$

且有  $\lambda^{(n)}(\tilde{\mathcal{C}}^*) \leq 4\epsilon$  (若不满足, 则必然有  
 $\frac{1}{2^{nR}} \sum_{i=1}^{2^{nR}} \lambda_i(\mathcal{C}^*) > 2\epsilon$ , 不满足前提)

证 ②  $\lambda^{(n)} \rightarrow 0 \Rightarrow R \leq C$

$$W \rightarrow X_C^n(W) \rightarrow Y^n \rightarrow \hat{W} = g(Y^n) \text{ 构成 Markov chain } (2^{nR}, n)$$

$$\text{令 } P_e^{(n)} = \Pr[W \neq \hat{W}]$$

对  $W \rightarrow Y^n \rightarrow \hat{W}$  使用 Fano 不等式

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

$$H(W|\hat{W}) = H(W) - I(W; \hat{W})$$

因为  $W$  是在  $2^{nR}$  的均匀分布上取值, 故  $H(W) = \log 2^{nR} = nR$

$$\text{所以 } nR - I(W; \hat{W}) < 1 + P_e^{(n)} nR$$

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

$$I(X^n; Y^n) \leq nC$$

$$\text{所以 } nR - nC < 1 + P_e^{(n)} nR$$

$$\Rightarrow R < \frac{1+nC}{n(1-P_e^{(n)})} = \frac{1}{n(1-P_e^{(n)})} + \frac{C}{1-P_e^{(n)}} \\ n \rightarrow \infty, \quad \Rightarrow \infty \quad \Rightarrow C$$

Proof.

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.

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  $R \leq C$ . □