

## Elements of Information Theory

### Chapter 5: Channel Capacity and Channel Coding Theorem

Bilingual course  
(Chinese taught course)

Information and Communication Eng. Dept.

Deng Ke

## Outline

- Sequence
- Channel Capacity
- Symmetric Channel
- Decoding Rule
- Joint Typical Set and Joint AEP
- Channel Coding Theorem

Xi'an Jiaotong University

## Sequence

- The channel input is  $x^N$  with length  $N$ , after  $N$  transmissions, the output is  $y^N$  also with length  $N$   

$$X^N = (x_1, x_2, \dots, x_N), x_i \in \{a_1, a_2, \dots, a_K\} \quad K^N$$

$$Y^N = (y_1, y_2, \dots, y_N), y_i \in \{b_1, b_2, \dots, b_J\} \quad J^N$$
- This can be regarded as once transmission whose channel model with  $K^N$  input and  $J^N$  output
- Probability  $p(y^N | x^N) = \prod_{i=1}^N p(y_i | x_i)$
- Condition entropy

Xi'an Jiaotong University

## Sequence

$$\begin{aligned} H(X^N | Y^N) &= \sum_{x^N, y^N} p(x^N, y^N) \log \frac{1}{p(x^N | y^N)} \\ &= \dots = \sum_{x_1, y_1} p(x_1, y_1) \log \frac{1}{p(x_1 | y_1)} + \dots + \sum_{x_N, y_N} p(x_N, y_N) \log \frac{1}{p(x_N | y_N)} \\ &= H(X_1 | Y_1) + \dots + H(X_N | Y_N) \\ &= NH(X | Y) \end{aligned}$$

- Then

$$\begin{aligned} I(X^N; Y^N) &= H(X^N) - H(X^N | Y^N) \\ &= NH(X) - NH(X | Y) \\ &= N[H(X) - H(X | Y)] \\ &= NI(X; Y) \end{aligned}$$

Xi'an Jiaotong University

## Channel Capacity

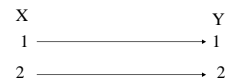
- **Definition** We define the "information" channel capacity of a discrete memoryless channel as  

$$C = \max_{p(x)} I(X; Y) \quad (7.1)$$
 where the maximum is taken over all possible input distributions  $p(x)$
- Interpretation
  - Mutual information: A measure of the amount of information that one random variable contains about another random variable.
  - Original uncertainty (information):  $H(X)$
  - Uncertainty at the receiver (Y observed):  $H(X|Y)$
  - Uncertainty reduction:  $H(X) - H(X|Y) = I(X; Y)$

Xi'an Jiaotong University

## Some Examples

- Noiseless Channel



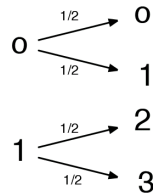
$$\begin{aligned} C &= \max I(X; Y) \\ &= \max [H(X) - H(X | Y)] \\ &= \max [H(X) - 0] \\ &= 1 \end{aligned}$$

Xi'an Jiaotong University

## Some Examples

- 2-input-4-output channel

$$\begin{aligned} C &= \max I(X;Y) \\ &= \max [H(X) - H(X|Y)] \\ &= \max [H(X) - 0] \\ &= 1 \end{aligned}$$

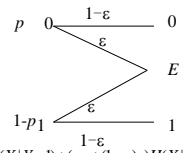


Xi'an Jiaotong University

## Some Examples

- Binary erasure channel

$$\begin{aligned} C &= \max I(X;Y) = \max H(X) - H(X|Y) \\ &= \max H(X) - \sum p(y) H(X|Y=y) \\ &= \max H(X) - [p(1-\epsilon)H(X|Y=0) + (1-p)(1-\epsilon)H(X|Y=1) + (p\epsilon + (1-p)\epsilon)H(X|Y=E)] \\ &= \max H(X) - [p(1-\epsilon)0 + (1-p)(1-\epsilon)0 + \epsilon H(X)] \\ &= \max H(X)(1-\epsilon) \\ &= (1-\epsilon) \end{aligned}$$

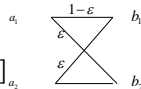


Xi'an Jiaotong University

## Some Examples

- BSC

$$\begin{aligned} C &= \max I(X;Y) \\ &= \max [H(Y) - H(Y|X)]_{a_i} \\ &= \max [H(Y) - \Omega(\epsilon)] \\ &= 1 - \Omega(\epsilon) \end{aligned}$$



$$\begin{aligned} I(X;Y) &= \Omega(\epsilon + p - 2\epsilon p) - \Omega(\epsilon) \quad \text{To achieve } C \\ \epsilon + p - 2\epsilon p &= 1/2 \Rightarrow (1-2\epsilon)p = (1-2\epsilon)/2 \Rightarrow p = 1/2 \end{aligned}$$

Xi'an Jiaotong University

## Channel Capacity

- Review
- A function  $f(x)$  is said to be *convex* over an interval  $(a,b)$  if for every  $x_1, x_2 \in (a,b)$  and  $0 \leq \lambda \leq 1$ 

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$
- (Jensen's inequality) If  $f$  is a convex function and  $X$  is a random variable,  $Ef(X) \geq f(EX)$
- The sum of cap(cup) functions is also a cap(cup) function.
- The mutual information  $I(X;Y)$  is a cap function of  $p(x)$

Xi'an Jiaotong University

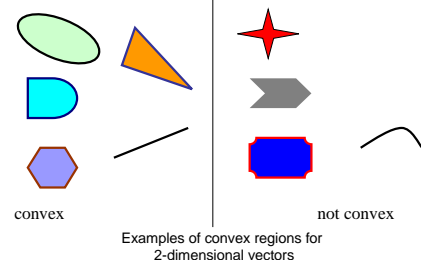
## Convex Region

- A region  $R$  is convex if for each vector  $\vec{\alpha}, \vec{\beta}$  in  $R$ , the vector  $\theta\vec{\alpha} + (1-\theta)\vec{\beta}$  is in  $R$  for  $0 \leq \theta \leq 1$ 

$$\begin{aligned} \theta\vec{\alpha} + (1-\theta)\vec{\beta} &= \theta\vec{\alpha} + (1-\theta)\vec{\alpha} - (1-\theta)\vec{\alpha} + (1-\theta)\vec{\beta} \\ &= \vec{\alpha} + (1-\theta)(\vec{\beta} - \vec{\alpha}) \end{aligned}$$
- For each pair of points in the region, the straight line between those points stays in the region
- A vector is defined to be a probability vector if its components are all nonnegative and sum to 1
- The region of probability vector is convex

Xi'an Jiaotong University

## Convex Region



Xi'an Jiaotong University

## Convex Region

- Let  $\bar{\alpha}$  and  $\bar{\beta}$  be probability vectors and let  $\bar{\gamma} = \theta \bar{\alpha} + (1-\theta) \bar{\beta}$
- Nonnegative  $\gamma_i = \theta \alpha_i + (1-\theta) \beta_i \geq 0$
- Sum  $\sum_i \gamma_i = \sum_i [\theta \alpha_i + (1-\theta) \beta_i]$   
 $= \theta \sum_i \alpha_i + (1-\theta) \sum_i \beta_i = \theta + (1-\theta) = 1$
- A real-valued function  $f$  of a vector is defined to be convex cap over a convex region  $R$  if for all  $\bar{\alpha}, \bar{\beta}$  in  $R$ , and  $\theta, 0 < \theta < 1$  the function satisfies  $\theta f(\bar{\alpha}) + (1-\theta) f(\bar{\beta}) \leq f[\theta \bar{\alpha} + (1-\theta) \bar{\beta}]$

Xi'an Jiaotong University

## Maximize of a convex cap function

- Let  $f(\bar{\alpha})$  be a convex cap function of  $\bar{\alpha}$ , where  $\bar{\alpha}$  is a probability vector. Assume that the partial derivatives,  $\partial f(\bar{\alpha}) / \partial \alpha_i$  are defined and continuous over the region  $R$  with possible exception that  $\lim_{\alpha_i \rightarrow 0} \frac{\partial f(\bar{\alpha})}{\partial \alpha_i} = \infty$ . Then the necessary and sufficient conditions on a probability vector  $\bar{\alpha}$  to maximize  $f$  over the region  $R$  is

$$\frac{\partial f(\bar{\alpha})}{\partial \alpha_i} = \lambda \quad \text{all } k \text{ such that } \alpha_i > 0$$

$$\frac{\partial f(\bar{\alpha})}{\partial \alpha_i} \leq \lambda \quad \text{all } k \text{ such that } \alpha_i = 0$$

Xi'an Jiaotong University

## Lagrange Multiplier

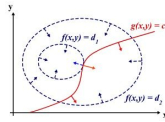
$$\begin{cases} \max & f(x, y) \\ g(x, y) = c \end{cases}$$

$$L(x, y, \lambda) = f(x, y) - \lambda [g(x, y) - c]$$

$$\frac{\partial L(x, y, \lambda)}{\partial \lambda} = 0 \Rightarrow g(x, y) = c$$

$$\frac{\partial L(x, y, \lambda)}{\partial x} = 0 \Rightarrow \frac{\partial f(x, y)}{\partial x} = \lambda \frac{\partial g(x, y)}{\partial x}$$

$$\frac{\partial L(x, y, \lambda)}{\partial y} = 0 \Rightarrow \frac{\partial f(x, y)}{\partial y} = \lambda \frac{\partial g(x, y)}{\partial y}$$



Xi'an Jiaotong University

## Finding Channel Capacity

$$\frac{\partial I(X;Y)}{\partial Q(x_n)} = \lambda \quad \text{all } k \text{ with } Q(x_n) > 0; \quad \frac{\partial I(X;Y)}{\partial Q(x_n)} \leq \lambda \quad \text{all } k \text{ with } Q(x_n) = 0$$

$$\begin{aligned} \frac{\partial I(X;Y)}{\partial Q(x_n)} &= \frac{\partial}{\partial Q(x_n)} \left[ \sum_j \sum_i Q(x_i) p(y_j | x_i) \log \frac{p(y_j | x_i)}{p(y_j)} \right] \\ &= \frac{\partial}{\partial Q(x_n)} \left[ \sum_j \sum_i \left[ Q(x_i) p(y_j | x_i) \log \frac{p(y_j | x_i)}{\sum_k Q(x_k) p(y_j | x_k)} \right] \right] \\ &= \sum_j p(y_j | x_n) \log \frac{p(y_j | x_n)}{p(y_j)} + \sum_j \sum_i \left[ Q(x_i) p(y_j | x_i) \frac{\partial}{\partial Q(x_n)} \log \frac{p(y_j | x_i)}{p(y_j)} \right] \\ &= \frac{\partial}{\partial Q(x_n)} \log \frac{p(y_j | x_n)}{p(y_j)} = \log e - \frac{p(y_j)}{p(y_j | x_n)} \left( p(y_j | x_n) \left( -\frac{1}{p^2(y_j)} \right) \right) + \frac{\partial}{\partial Q(x_n)} p(y_j) \\ &= -\frac{\log e}{p(y_j)} + \frac{\partial}{\partial Q(x_n)} \sum_i Q(x_i) p(y_j | x_i) \\ &= -\frac{\log e}{p(y_j)} + p(y_j | x_n) \end{aligned}$$

Xi'an Jiaotong University

## Finding Channel Capacity

$$\begin{aligned} & \sum_k \sum_j Q(x_k) p(y_j | x_k) - \frac{\log e}{p(y_j)} p(y_j | x_n) \\ &= -\log e \sum_k \sum_j p(x_k, y_j) \frac{p(y_j | x_n)}{p(y_j)} \\ &= -\log e \sum_j \frac{p(y_j | x_n)}{p(y_j)} \sum_k p(x_k, y_j) \\ &= -\log e \sum_j \frac{p(y_j | x_n)}{p(y_j)} p(y_j) \\ &= -\log e \sum_j p(y_j | x_n) \\ &= -\log e \end{aligned}$$

Xi'an Jiaotong University

## Finding Channel Capacity

$$\begin{aligned} \frac{\partial I(X;Y)}{\partial Q(x_n)} &= \sum_j p(y_j | x_n) \log \frac{p(y_j | x_n)}{p(y_j)} - \log e \\ &= I(X = x_n; Y) - \log e \end{aligned}$$

- Then

$$\frac{\partial I(X;Y)}{\partial Q(x_n)} = I(X = x_n; Y) - \log e = \lambda \quad Q(x_n) > 0,$$

$$\frac{\partial I(X;Y)}{\partial Q(x_n)} = I(X = x_n; Y) - \log e \leq \lambda \quad Q(x_n) = 0,$$

- A set of necessary and sufficient conditions on an input probability vector  $Q^{(n)} = [Q(x_1), Q(x_2), \dots, Q(x_n)]$  to achieve capacity on a DMC is

Xi'an Jiaotong University

## Symmetric Channel

$$I(X = x_s; Y) = C \quad Q(x_s) > 0,$$

$$I(X = x_s; Y) \leq C \quad Q(x_s) = 0.$$

- inputs as rows and outputs as columns
- The rows of the probability transition matrix are permutations of each other and so are the columns. *symmetric I*
- The rows of the probability transition matrix are permutations of each other and all the column sums are equal. *weakly symmetric*
- (The columns can be partitioned into subsets in such a way that in each subset, the rows are permutations of each other and so are the columns (if more than 1). *symmetric II*)

Xi'an Jiaotong University

## Symmetric Channel

$$\begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.2 & 0.5 & 0.3 \\ 0.5 & 0.3 & 0.2 \end{bmatrix} \quad \begin{bmatrix} 1/3 & 1/6 & 1/2 \\ 1/2 & 1/2 & 1/6 \end{bmatrix} \quad \begin{bmatrix} 0.7 & 0.2 & 0.1 \\ 0.1 & 0.2 & 0.7 \end{bmatrix}$$

- For a symmetric DMC, capacity is achieved by using the inputs with equal probability

Xi'an Jiaotong University

## Symmetric Channel

$$I(x = k; Y) = I_1(x = k; Y) + I_2(x = k; Y)$$

$$I(x = 1; Y) = \sum_{j=1}^K p(j|1) \log \frac{p(j|1)}{p(j)} + \sum_{j=N+1}^L p(j|1) \log \frac{p(j|1)}{p(j)}$$

$$= \sum_{j=1}^K p(j|1) \log \frac{p(j|1)}{\sum_i p(j|i)} + \sum_{j=N+1}^L p(j|1) \log \frac{p(j|1)}{\sum_i Q(i)p(j|i)}$$

$$\frac{Q(i) = 1/K}{\sum_{j=1}^K p(j|1) \log \frac{p(j|1)}{\frac{1}{K} \sum_i p(j|i)} + \sum_{j=N+1}^L p(j|1) \log \frac{p(j|1)}{\frac{1}{K} \sum_i p(j|i)}}$$

$$I(x = 2; Y) = \sum_{j=1}^K p(j|2) \log \frac{p(j|2)}{\frac{1}{K} \sum_i p(j|i)} + \sum_{j=N+1}^L p(j|2) \log \frac{p(j|2)}{\frac{1}{K} \sum_i p(j|i)}$$

⋮

$$I(x = K; Y) = \sum_{j=1}^K p(j|K) \log \frac{p(j|K)}{\frac{1}{K} \sum_i p(j|i)} + \sum_{j=N+1}^L p(j|K) \log \frac{p(j|K)}{\frac{1}{K} \sum_i p(j|i)}$$

Xi'an Jiaotong University

## Symmetric Channel

- Consider the column permutation

$$\frac{1}{K} \sum p(j|i) = \text{const}$$

- Consider the row permutation

$$I_1(x = k; Y) = \text{const}$$

- And then  $I_2(x = k; Y) = \text{const}$

- Hence

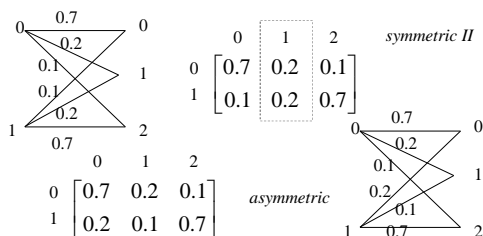
$$I(x = k; Y) = \text{const} \quad k = 1, 2, \dots, K$$

- Finally, we have

$$C = I(x = k; Y) \quad k = 1, 2, \dots, K$$

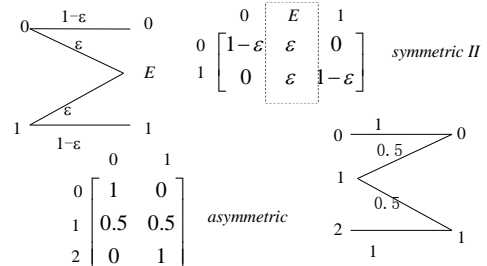
Xi'an Jiaotong University

## Symmetric Channel



Xi'an Jiaotong University

## Symmetric Channel



西安交通大学信通系

## Symmetric Channel

- Example: Binary Erasure Channel

$$C = I(X=0;Y) \big|_{p(x=0)=0.5}$$

$$= \sum_y p(y|x=1) \log \frac{p(y|x=1)}{p(y)}$$

$$= (1-\varepsilon) \log \frac{1-\varepsilon}{\frac{1}{2}(1-\varepsilon)} + \varepsilon \log \frac{\varepsilon}{\varepsilon}$$

$$= 1-\varepsilon$$

西安交通大学信通系

## Weakly Symmetric Channel

- Theorem 7.2.1 For a weakly symmetric channel  $C = \log |Y| - H(\text{row of a transition matrix})$
- and this is achieved by a uniform distribution on the input alphabet
- Proof: Letting  $r$  be a row of the transition matrix, we have  $I(X;Y) = H(Y) - H(Y|X) = H(Y) - H(r)$
- From before we know, when  $p(x) = 1/|X|$ , the capacity is achieved, and we have

$$p(y) = \sum_x p(x|y)p(x) = \frac{1}{|X|} \sum_x p(x|y) = \frac{c}{|X|} = \frac{1}{|Y|}$$

西安交通大学信通系

## Decoding Rule

- Minimum error decoding rule
  - After receiving  $y_j$ , for all a posteriori probability  $P(x_i|y_j)$ ,  $P(x_2|y_j), \dots, P(x_i|y_j), \dots$  Find maximum  $P(x^*|y_j)$
  - i.e. if  $x_i \neq x^*$  we have  $P(x^*|y_j) > P(x_i|y_j)$ , then  $g(y_j) = x^*$
  - This rule is also called MAP rule
- Example
  - Let the source probability and channel matrix is

$$P(x) = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$$

$$\begin{bmatrix} 1/2 & 1/3 & 1/6 \\ 1/6 & 1/2 & 1/3 \\ 1/3 & 1/6 & 1/2 \end{bmatrix}$$

西安交通大学信通系

## Decoding Rule

- Find the decoding scheme to minimize  $P_e$
- According to the MAP rule, first we get the MAP

$$P(x/y) = \frac{P(xy)}{P(y)}$$

$$P(xy) = \begin{bmatrix} 1/4 & 1/6 & 1/12 \\ 1/24 & 1/8 & 1/12 \\ 1/12 & 1/24 & 1/8 \end{bmatrix}$$

- and  $P(y) = [3/8 \ 1/3 \ 7/24]$ ,  $P(x/y) = x_2 \begin{bmatrix} 2/3 & 1/2 & 2/7 \\ 1/9 & 3/8 & 2/7 \\ 2/9 & 1/8 & 3/7 \end{bmatrix}$
- hence

西安交通大学信通系

## Decoding Rule

Then, the decoding scheme is

$$y_1 \rightarrow x_1, y_2 \rightarrow x_1, y_3 \rightarrow x_3$$

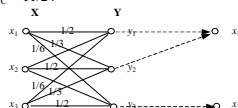
And the probability of correct transmission

$$P_C = P(x_1) \cdot P(y_1|x_1) + P(x_1) \cdot P(y_2|x_1) + P(x_3) \cdot P(y_3|x_3)$$

$$= 1/2 \times 1/2 + 1/2 \times 1/3 + 1/4 \times 1/2 = 13/24$$

The error probability is

$$P_E = 1 - P_C = 11/24$$



西安交通大学信通系

## Decoding Rule

- Maximum Likelihood Rule
  - MAP rule is the minimum error rule, but the computation is complex, we need other rule
  - When the source with  $M$  states is uniformly distributed
  - we have  $P(x_i) = \frac{1}{M}$
  - $P(x_i/y_j) = \frac{P(x_i y_j)}{P(y_j)} = \frac{1}{MP(y_j)} P(y_j/x_i)$
  - which means, we can decode only from the  $P(y_j/x_i)$
  - i.e. if  $x_i \neq x^*$  we have  $P(y_j/x^*) > P(y_j/x_i)$ , then  $g(y_j) = x^*$

西安交通大学信通系

## Decoding Rule

- Example  
Let a BSC with  $\epsilon=0.01$ , the source is uniformly distributed  
(1) Find minimum  $P_E$  (2) After the channel code "0"  $\rightarrow$  "000", "1"  $\rightarrow$  "111", find minimum  $P_E$
- Solution  
(1)  $\epsilon=P(1|0)=P(0|1)=0.01$ ,  $P(0)=P(1)=1/2$ , then  $g(0)=0$ ,  $g(1)=1$ , and  
$$P_E = \frac{1}{2} \cdot (0.01 + 0.01) = 10^{-2}$$
  
(2) Let the channel input  $\alpha_1=000$   $\alpha_2=001$ , the output is  
 $\beta_0=000$ ,  $\beta_1=001$ ,  $\beta_2=010$ ,  $\beta_3=100$ ,  $\beta_4=011$ ,  $\beta_5=101$ ,  $\beta_6=110$ ,  $\beta_7=111$   
and the channel matrix is

西安交通大学信通系

## Decoding Rule

$$\begin{array}{cccccccc} & 000 & 001 & 010 & 100 & 011 & 101 & 110 & 111 \\ \text{000} & (1-\epsilon)^3 & (1-\epsilon)^2\epsilon & (1-\epsilon)^2\epsilon & (1-\epsilon)^2\epsilon & (1-\epsilon)\epsilon^2 & (1-\epsilon)\epsilon^2 & (1-\epsilon)\epsilon^2 & \epsilon^3 \\ \text{111} & \epsilon^3 & (1-\epsilon)\epsilon^2 & (1-\epsilon)\epsilon^2 & (1-\epsilon)\epsilon^2 & (1-\epsilon)^2\epsilon & (1-\epsilon)^2\epsilon & (1-\epsilon)^2\epsilon & (1-\epsilon)^3 \end{array}$$

$g(\beta_0)=g(\beta_1)=g(\beta_2)=g(\beta_3)=\alpha_1=000 \rightarrow 0$   
 $g(\beta_4)=g(\beta_5)=g(\beta_6)=g(\beta_7)=\alpha_2=111 \rightarrow 1$   
 and the error probability  
 $P_E = 3 \cdot (1-\epsilon)\epsilon^2 + \epsilon^3 \approx 3\epsilon^2 =$

西安交通大学信通系

## Channel Coding Theorem

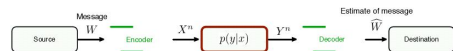
- Information channel capacity:
  - The maximum of the mutual information
- Operational channel capacity:
  - Highest rate (bits/channel use) that can communicate at reliably
- Channel coding theorem says: information capacity = operational capacity
- Definition: the rate  $R$  of an  $(M, n)$  code is

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

西安交通大学信通系

## Channel Coding Theorem

- Definition: Achievability. A rate  $R$  is called *achievable* if there exists a sequence of  $(\lceil 2^{nR} \rceil, n)$  codes such that the maximal probability of error  $\lambda^{(n)}$  (i.e. maximal  $\Pr\{\text{Error}\}$ ) tends to 0 as  $n \rightarrow \infty$ . Note  $(2^{nR}, n)$  codes means  $(\lfloor 2^{nR} \rfloor, n)$  codes
- Definition: Capacity. The *Capacity* of a channel is supremum of all achievable rates.



西安交通大学信通系

## Channel Coding Theorem

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

1. An index set  $\{1, 2, \dots, M\}$  over messages  $W$ .
2. An encoding function  $X^n: \{1, 2, \dots, M\} \rightarrow \mathcal{X}^n$ , yielding codewords  $x^n(1), x^n(2)$ . (This set is called the *codebook*  $\mathcal{C}$ .)  
 $x^n(W)$  passes through the channel and is received as a random sequence  $Y^n \sim p(y^n | x^n)$ .
3. A (deterministic) decoding function

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

which is an estimator  $\hat{W} = g(Y^n)$  of  $W \in \{1, 2, \dots, M\}$ . It declares an error if  $\hat{W} \neq W$ .

西安交通大学信通系

## Channel Coding Theorem

Definition: Conditional probability of error.

$$\text{Let } \lambda_i = \Pr\{g(Y^n) \neq i | X^n = x^n(i)\} = \sum_{y^n} p(y^n | x^n(i)) I_{\{g(y^n) \neq i\}}(y^n)$$

be the *conditional probability of error* given that index  $i$  was sent.  
Definition: Maximal probability of error. The maximal probability of error  $\lambda^{(n)}$  for an  $(M, n)$  code is defined as

$$\lambda^{(n)} = \max_{i \in \{1, 2, \dots, M\}} \lambda_i.$$

Definition: Average probability of error. The (arithmetic) average probability of error  $P_e^{(n)}$  for an  $(M, n)$  code is:

$$P_e^{(n)} = \frac{1}{M} \sum_{i=1}^M \lambda_i.$$

- Note that  $P_e^{(n)} = \Pr\{W \neq g(Y^n)\}$  if  $W$  is chosen uniformly.
- Also,  $P_e^{(n)} \leq \lambda^{(n)}$ ; i.e., the average probability of error is less than the maximal probability of error.

西安交通大学信通系

## Joint Typical Sequences

- Definition: The set  $A_\epsilon^{(n)}$  of jointly typical sequences  $\{(x^n, y^n)\}$  with respect to  $p(x, y)$  is the set of  $n$ -sequences with empirical entropies  $\epsilon$ -close to the true entropies

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

where

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

西安交通大学信通系

## Joint AEP

- Theorem 7.6.1 (Joint AEP) Let  $(X^n, Y^n)$  be sequences of length  $n$  drawn i.i.d according to  $p(x^*, y^*) = \prod_{i=1}^n p(x_i, y_i)$ , then

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

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

$$1 = \sum_{(x^n, y^n) \in A_\epsilon^{(n)}} p(x^n, y^n)$$

$$\geq \sum_{(x^n, y^n) \in A_\epsilon^{(n)}} p(x^n, y^n)$$

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

西安交通大学信通系

## Joint AEP

- If  $(\tilde{X}^n, \tilde{Y}^n) \sim p(x^n)p(y^n)$  so that  $\tilde{X}^n$  and  $\tilde{Y}^n$  are independent with the same marginals as  $p(x^n, y^n)$  then, we have

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

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

- for sufficient large  $n$

$$\Pr((\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)}) = \sum_{(x^n, y^n) \in A_\epsilon^{(n)}} p(x^n)p(y^n)$$

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

$$= 2^{-n(I(X; Y) - 3\epsilon)}.$$

西安交通大学信通系

## Joint AEP

For sufficiently large  $n$ ,  $\Pr(A_\epsilon^{(n)}) \geq 1 - \epsilon$ , and therefore

$$1 - \epsilon \leq \sum_{(x^n, y^n) \in A_\epsilon^{(n)}} p(x^n, y^n)$$

$$\leq |A_\epsilon^{(n)}| 2^{-n(H(X, Y) - \epsilon)} \quad |A_\epsilon^{(n)}| \geq (1 - \epsilon) 2^{n(H(X, Y) - \epsilon)}$$

$$\Pr((\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)}) = \sum_{(x^n, y^n) \in A_\epsilon^{(n)}} p(x^n)p(y^n)$$

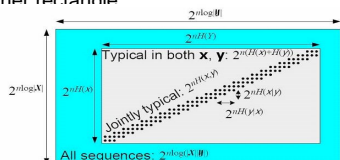
$$\geq (1 - \epsilon) 2^{n(H(X, Y) - \epsilon)} 2^{-n(H(X) + \epsilon)} 2^{-n(H(Y) + \epsilon)}$$

$$= (1 - \epsilon) 2^{-n(I(X, Y) + 3\epsilon)}. \quad \square$$

西安交通大学信通系

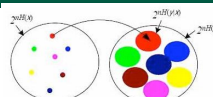
## Explanation

- There are about  $2^{nH(X)}$  typical  $X$  in all
- Each typical  $Y$  is jointly typical with about  $2^{nH(X|Y)}$  of those typical  $X$ 's
- The joint typical pairs are a fraction  $2^{-nI(X; Y)}$  of inner rectangle



西安交通大学信通系

## Explanation



- For each typical input sequence (how many?) there are about  $2^{nH(Y|X)}$  possible output sequences, all equally likely.
- Want to ensure that no two typical input sequences produce the same output sequences.

西安交通大学信通系

## Channel Coding Theorem

- Intuition

- Random choice of codeword
- Decoding rule: joint typicality decoding

The probability that any other codeword looks jointly typical with the received sequence is  $2^{-nI(X;Y)}$ . Hence, if we have fewer than  $2^{nI(X;Y)}$  codewords in the input side, then with high probability there will be no other codewords that can be confused with the transmitted codeword, and the probability of error is small.

西安交通大学信通系

## Channel Coding Theorem

- **Theorem 7.7.1 (Channel coding theorem)** For a DMC, 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  $R \leq C$

- Proof: All capacity proof can be divided into two parts:
  - i) Achievability: All  $R \leq C$  are achievable
  - ii) Converse: No rate  $R > C$  is achievable
- Only Achievability

西安交通大学信通系

## Channel Coding Theorem

- Generate a  $(2^{nR}, n)$  code at random according to distribution  $p(x)$

$$C = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_n(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(2^{nR}) & x_2(2^{nR}) & \dots & x_n(2^{nR}) \end{bmatrix} \quad \Pr(C) = \prod_{i=1}^{2^{nR}} \prod_{j=1}^n p(x_j(w))$$

- Each entry in this matrix is generated i.i.d.
- The code  $C$  is then revealed to both sender and receiver. Both sender and receiver are also assumed to know the channel transition matrix  $p(y|x)$  for the channel.
- A message  $W$  is chosen according to a uniform distribution  $\Pr(W = w) = 2^{-nR}$ ,  $w = 1, 2, \dots, 2^{nR}$

西安交通大学信通系

## Channel Coding Theorem

- The codeword corresponding to the  $w$ -throw of  $C$ , is sent over the channel.
- The receiver receives a sequence  $Y^n$  according to the distribution  $p(Y^n | X^n(w)) = \prod_{i=1}^n p(y_i | x_i(w))$
- The receiver guesses which message was sent use jointly typical decoding
  - 1.  $(X^n(\hat{w}), Y^n)$  is jointly typical
  - 2. There is NO other index  $w' \neq \hat{w}$  such that  $(X^n(w'), Y^n) \in A_\epsilon^{(n)}$ . Otherwise it declare an error
- Decoding error if  $w \neq \hat{w}$ , Let  $\epsilon$  be the event  $\{w \neq \hat{w}\}$

西安交通大学信通系

## Channel Coding Theorem

- Instead of calculating the probability of error for a single code, we calculate the average over all codes generated at random according to the distribution.
- For a typical codeword, two kinds of error :
  - Either the output  $Y^n$  is not jointly typical with the transmitted codeword
  - Or there is some other codeword that is jointly typical with  $Y^n$ .
- For any rival codeword, the probability that it is jointly typical with the received sequence is approximately  $2^{-nI(X;Y)}$

西安交通大学信通系

## Channel Coding Theorem

$$\Pr(\epsilon) = \sum_C \Pr(C) P_\epsilon^{(n)}(C) = \sum_C \Pr(C) \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(C) = \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_C \Pr(C) \lambda_w(C)$$

- By the symmetry of the code construction, the average probability of error averaged over all codes does not depend on the particular index that was sent.
- Thus, we can assume without loss of generality that the message  $W = 1$  was sent

$$\Pr(\epsilon) = \sum_C \Pr(C) \lambda_1(C) = \Pr(\epsilon | W = 1)$$

西安交通大学信通系



## Channel Coding Theorem

- Let  $E_i = \{(X^n(i), Y^n) \text{ is in } A_\phi^{(n)}\}$ ,  $i = 1, 2, \dots, 2^{nR}$
  - Error scenario:
    - 1) When the transmitted codeword and the received sequence are not jointly typical
    - 2) When a wrong codeword is jointly typical with the received sequence
- $$\Pr(\mathcal{E} | W = 1) = P(E_1^c \cup E_2 \cup \dots \cup E_{2^{nR}} | W = 1)$$
- By joint AEP  $P(E_1^c | W = 1) \rightarrow 0 \Rightarrow P(E_1^c | W = 1) \leq \epsilon$
  - By the code generation process,  $X^n(1)$  and  $X^n(i)$  are independent for  $i \neq 1$ . 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

西安交通大学信通系

## Channel Coding Theorem

$$\begin{aligned} \Pr(\mathcal{E}) &= \Pr(\mathcal{E} | W = 1) \leq P(E_1^c | W = 1) + \sum_{i=2}^{2^{nR}} P(E_i | W = 1) \\ &\leq \epsilon + \sum_{i=2}^{2^{nR}} 2^{-n(I(X;Y) - 3\epsilon)} \\ &= \epsilon + (2^{nR} - 1) 2^{-n(I(X;Y) - 3\epsilon)} \\ &\leq \epsilon + 2^{nR} 2^{-n(I(X;Y) - 3\epsilon)} \leq 2\epsilon \quad R < I(X;Y) - 3\epsilon \\ \Pr(\mathcal{E}) &\xrightarrow{n \rightarrow \infty} 0 \quad \text{if } R < I(X;Y) \end{aligned}$$

- Choose  $p(x) \rightarrow p^*(x)$  that maximize  $I(X;Y)$

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

西安交通大学信通系

## Channel Coding Theorem

- Although the theorem shows that there exist good codes with arbitrarily small probability of error for long block lengths, it does not provide a way of constructing the best codes.
- However, without some structure in the code, it is very difficult to decode (the simple scheme of table lookup requires an exponentially large table). Hence the theorem does not provide a practical coding scheme.

西安交通大学信通系