

Definition 2.3.5

Given a channel with input alphabet \mathcal{X} output alphabet \mathcal{Y} and transition probabilities

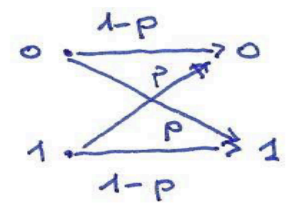
$P(y|x) \geq 0$; $\sum_{y \in \mathcal{Y}} P(y|x) = 1$, its capacity

is given by

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

Example 2.3.2.

Binary symmetric channel



0 ≤ p ≤ 1

Transition matrix

$P(y|x) = \begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$
↑ column
↑ row

$I(X; Y) = H(Y) - H(Y|X) = H(Y) + \sum_{x=0}^1 P(x) \sum_{y=0}^1 P(y|x) \log P(x|x)$
 $= H(Y) + p \log p + (1-p) \log(1-p) = H(Y) - h(p)$

$h(p) := -p \log p - (1-p) \log(1-p)$

$$H(Y) = - \sum_{y=0}^1 P_Y(y) \log P_Y(y) ; \quad P_Y(y) = \sum_{x=0}^1 P(y|x) P_X(x)$$

$$P_Y(0) = P(0|0) P_X(0) + P(0|1) P_X(1) = (1-p) P_X(0) + p P_X(1)$$

$$P_Y(1) = P(1|0) P_X(0) + P(1|1) P_X(1) = p P_X(0) + (1-p) P_X(1)$$

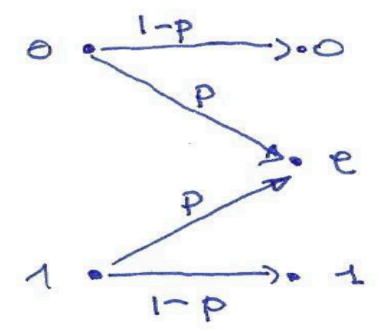
$$H(Y) \leq 1 ; \quad H(Y) = 1 \iff P_Y(0) = \frac{1}{2} = P_Y(1) \iff P_X(0) = P_X(1) = \frac{1}{2}$$

$$C = \max_{\pi_X} H(Y) - h(p) = 1 - h(p) \iff \pi_X = \left(\frac{1}{2}, \frac{1}{2}\right)$$

Example 2.3.3

Binary erasure channel

$$P(y|x) = \begin{pmatrix} 1-p & p & 0 \\ 0 & p & 1-p \end{pmatrix}$$



$$H(Y|X=0) = H(Y|X=1) = h(p) \implies$$

$$I(X;Y) = H(Y) - h(p)$$

$$H(Y) = -P_Y(0) \log P_Y(0) - P_Y(e) \log P_Y(e) - P_Y(1) \log P_Y(1)$$

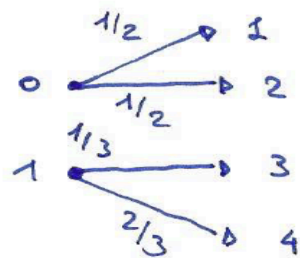
$$P_Y(0) = (1-p)P_X(0) ; P_Y(e) = pP_X(0) + pP_X(1) = p ; P_Y(1) = (1-p)P_X(1)$$

$$\begin{aligned} H(Y) &= - (1-p)P_X(0) (\log(1-p) + \log P_X(0)) - p \log p - (1-p)P_X(1) (\log(1-p) + \log P_X(1)) \\ &= H(p) + (1-p) H(\pi) \quad \pi := P_X(0) \end{aligned}$$

$$C = \max_{\pi_X} H(Y) - H(p) = (1-p) \max_{\pi_X} H(\pi) = 1-p \quad \leftarrow \pi_X = \left(\frac{1}{2}, \frac{1}{2}\right)$$

Example 2.3.4

Noisy channel with non-overlapping outputs



$$P(Y|X) = \begin{pmatrix} 1/2 & 1/2 & 0 & 0 \\ 0 & 0 & 1/3 & 2/3 \end{pmatrix}$$

$$H(Y|X=0) = 1$$

$$H(Y|X=1) = -\frac{1}{3} + \log 3$$

$$H(Y|X) = P_X(0) + P_X(1) \left(\log 3 - \frac{2}{3} \right)$$

$$P_Y(1) = \frac{P_X(0)}{2} ; P_Y(2) = \frac{P_X(0)}{2} ; P_Y(3) = \frac{P_X(1)}{3} ; P_Y(4) = \frac{2P_X(1)}{3}$$

$$\begin{aligned}
 H(Y) &= - P_X(0) \log \frac{P_X(0)}{2} - \frac{P_X(1)}{3} \log \frac{P_X(1)}{3} - \frac{2}{3} P_X(1) \log \left(\frac{2}{3} P_X(1) \right) \\
 &= - P_X(0) \log P_X(0) + P_X(0) - P_X(1) \log P_X(1) + P_X(1) \log 3 - \frac{2}{3} P_X(1)
 \end{aligned}$$

$$\begin{aligned}
 I(X; Y) &= h(\pi) + P_X(1) \log 3 - \frac{2}{3} P_X(1) + P_X(0) - P_X(0) - P_X(1) \log 3 + \frac{2}{3} P_X(1) \\
 &= h(\pi) \leq 1 \quad \pi := P_X(0)
 \end{aligned}$$

$$C = \max_{\pi_X} h(\pi) = 1 \iff \pi_X = \{1, 1\}$$

Definition 2.3.6

132

An (M, n) code for the channel $(\mathcal{X}, P(Y|X), \mathcal{Y})$

consists of

1. an index set $\{1, 2, \dots, M\}$;

2. an encoding function $X^{(n)}: \{1, 2, \dots, M\} \rightarrow \mathcal{X}^{(n)}$

yielding codewords $X^{(n)}(w) = X_1(w)X_2(w)\dots X_n(w) \in \mathcal{X}^{(n)}$

where $w \in \{1, 2, \dots, M\}$;

3. a decoding function $g: \mathcal{Y}^{(n)} \rightarrow \{1, 2, \dots, M\}$

which assigns a guess to each received $\gamma^{(n)} \in \mathcal{Y}^{(n)}$.

Definition 2.3.7

Let $d_i = \text{Prob}(g(\gamma^{(n)}) \neq w \mid X^{(n)} = X^{(n)}(w))$

$$= \sum_{\gamma^{(n)} \in \mathcal{Y}^{(n)}: g(\gamma^{(n)}) \neq w} P(\gamma^{(n)} \mid X^{(n)}(w))$$

be the conditional probability of error upon sending the i message w .

The maximal probability of error for an (M, n) -code is

$$d^{(n)} = \max_{w \in \{1, 2, \dots, M\}} d_w$$

Definition 2.3.8

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

$$R := \frac{\log M}{n} \text{ bits per channel use}$$

A rate R is achievable if there exists a sequence of

$(2^{Rn}, n)$ codes such that the maximal probability

of error $d^{(n)} \xrightarrow[n]{} 0$.

Definition 2.3.5

The capacity of a discrete memoryless channel is the supremum of all achievable rates:

$$C = \sup_R \{ R \text{ achievable} \}$$

Remark: The capacity of a discrete memoryless channel is the highest number of bits that one can send per use of the noisy channel.

Remark : given codewords $x^{(n)}$ chosen i.i.d with single choice Shannon entropy $H(X)$ there will be roughly $\boxed{2^{nH(X)}}$ typical ^{input} words (see AEP in Example 2-1.3).

By sending the codewords $x^{(n)}$ through the discrete memoryless channel, one obtains output words $y^{(n)}$ that are also i.i.d. Indeed,

$$P(y^{(n)}) = \sum_{x^{(n)}} P(y^{(n)} | x^{(n)}) P(x^{(n)}) = \prod_{j=1}^n [P(y_j | x_j) P(x_j)] = \prod_{j=1}^n P(y_j)$$

Therefore there will be $\boxed{2^{nH(Y)}}$ typical output words.

Since $P(x^{(n)}, y^{(n)}) = P(y^{(n)} | x^{(n)}) P(x^{(n)}) = \prod_{j=1}^n P(y_j, x_j)$,

there will be $\boxed{2^{nH(X,Y)}}$ jointly typical pair of words.

By decoding $y^{(n)}$ with w if $y^{(n)}$ is jointly typical with $x^{(n)}(w)$, an error will occur with probability

$$\boxed{\frac{2^{nH(X,Y)}}{2^{n(H(X)+H(Y))}} = \frac{1}{2}^{-nI(X;Y)}}$$

Definition 2.3.10

The set $A_\epsilon^{(n)}$ of jointly typical sequences $\{(x^{(n)}, y^{(n)})\}$

with respect to the distribution $p(x, y)$ is the set of sequences such that

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

$$p(x^{(n)}, y^{(n)}) = \prod_{j=1}^n p(x_j, y_j).$$

Theorem 2.3.1

joint AEP

Let $(x^{(n)}, y^{(n)})$ be sequences of length n drawn i.i.d. according to

$$p(x^{(n)}, y^{(n)}) = \prod_{j=1}^n p(x_j, y_j). \text{ Then, } 1) \text{ Prob}((x^{(n)}, y^{(n)}) \in A_\epsilon^{(n)}) \xrightarrow{n} 1; \quad 2) \#(A_\epsilon^{(n)}) \leq 2^{n(H(X, Y) + \epsilon)}$$

3) if $\tilde{x}^{(n)}, \tilde{y}^{(n)}$ are independent with the marginals of $p(x^{(n)}, y^{(n)})$, then

$$\text{Prob}((\tilde{x}^{(n)}, \tilde{y}^{(n)}) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X; Y) - 3\epsilon)}; \quad (1 - \epsilon) 2^{-n(I(X; Y) + 3\epsilon)} \leq \text{Prob}((\tilde{x}^{(n)}, \tilde{y}^{(n)}) \in A_\epsilon^{(n)}).$$

Proof: from the weak law of large numbers

$$\begin{cases} -\frac{1}{n} \log p(x^{(n)}) \xrightarrow{P} H(X) \\ -\frac{1}{n} \log p(y^{(n)}) \xrightarrow{P} H(Y) \\ -\frac{1}{n} \log p(x^{(n)}, y^{(n)}) \xrightarrow{P} H(X, Y) \end{cases}$$

Then, given the sets

$$B_{1, \varepsilon}^{(n)} = \{x^{(n)} \in \mathcal{X}^{(n)} : \left| -\frac{1}{n} \log p(x^{(n)}) - H(X) \right| > \varepsilon\}$$

$$B_{2, \varepsilon}^{(n)} = \{y^{(n)} \in \mathcal{Y}^{(n)} : \left| -\frac{1}{n} \log p(y^{(n)}) - H(Y) \right| > \varepsilon\}$$

$$B_{3, \varepsilon}^{(n)} = \{(x^{(n)}, y^{(n)}) \in \mathcal{X}^{(n)} \times \mathcal{Y}^{(n)} : \left| -\frac{1}{n} \log p(x^{(n)}, y^{(n)}) - H(X, Y) \right| > \varepsilon\}$$

there exists \bar{n} such that, for $n > \bar{n}$,

$$\text{Prob}(B_{1, \varepsilon}^{(n)}) < \frac{\varepsilon}{3}, \text{Prob}(B_{2, \varepsilon}^{(n)}) < \frac{\varepsilon}{3}, \text{Prob}(B_{3, \varepsilon}^{(n)}) < \frac{\varepsilon}{3}$$

Notice that

$$A_{\varepsilon}^{(n)} = \mathcal{X}^{(n)} \times \mathcal{Y}^{(n)} \setminus (B_{1, \varepsilon}^{(n)} \cup B_{2, \varepsilon}^{(n)} \cup B_{3, \varepsilon}^{(n)})$$

Then,

$$\text{Prob}(A_{\varepsilon}^{(n)}) = 1 - \text{Prob}\left(\bigcup_{i=1}^3 B_{i, \varepsilon}^{(n)}\right) \geq 1 - \sum_{i=1}^3 \text{Prob}(B_{i, \varepsilon}^{(n)}) \geq 1 - \varepsilon$$

Notice that if $(x^{(n)}, y^{(n)}) \in A_{\epsilon}^{(n)}$

$$\begin{aligned} 2^{-n(H(X)+\epsilon)} &\leq p(x^{(n)}) \leq 2^{-n(H(X)-\epsilon)} \\ 2^{-n(H(Y)+\epsilon)} &\leq p(y^{(n)}) \leq 2^{-n(H(Y)-\epsilon)} \\ 2^{-n(H(X,Y)+\epsilon)} &\leq p(x^{(n)}, y^{(n)}) \leq 2^{-n(H(X,Y)-\epsilon)} \end{aligned}$$

$$\bullet 1 = \sum_{(x^{(n)}, y^{(n)}) \in \mathcal{X}^{(n)} \times \mathcal{Y}^{(n)}} p(x^{(n)}, y^{(n)}) \geq \sum_{(x^{(n)}, y^{(n)}) \in A_{\epsilon}^{(n)}} p(x^{(n)}, y^{(n)}) \geq \#(A_{\epsilon}^{(n)}) 2^{-n(H(X,Y)+\epsilon)}$$

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

$$\bullet 1 - \epsilon \leq \sum_{(x^{(n)}, y^{(n)}) \in A_{\epsilon}^{(n)}} p(x^{(n)}, y^{(n)}) \leq \#(A_{\epsilon}^{(n)}) 2^{-n(H(X,Y)-\epsilon)}$$

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

$$\bullet \dots \text{Prob}(\tilde{X}^{(n)}, \tilde{Y}^{(n)} \in A_{\epsilon}^{(n)}) = \sum_{(x^{(n)}, y^{(n)}) \in A_{\epsilon}^{(n)}} p(x^{(n)}) p(y^{(n)}) \leq \#(A_{\epsilon}^{(n)}) 2^{-n(H(X)+H(Y)-2\epsilon)} \leq 2^{-n(I(X;Y)-3\epsilon)}$$

$$\bullet \dots \text{Prob}(\tilde{X}^{(n)}, \tilde{Y}^{(n)} \in A_{\epsilon}^{(n)}) \geq \#(A_{\epsilon}^{(n)}) 2^{-n(H(X)+H(Y)+2\epsilon)} \geq (1 - \epsilon) 2^{-n(I(X;Y)+3\epsilon)}$$

Remark: Shannon's noisy channel theorem asserts that the channel capacity can be achieved so that, whatever the ~~noise~~ memoryless discrete noisy channel, one could send C bits per channel use in the limit of a large number of channel uses.

The ingredients in Shannon's proof are:

1. Allowing an arbitrarily small but non-zero probability of error.
2. Using the channel many times in order to get the law of large numbers into effect.
3. Calculating the average of the probability of error over a random choice of codebooks.

Remark: in the proof given below, one decodes by joint typicality and then looks for a code-word jointly typical with the received sequence. If there is only one that is decided to be the transmitted code-word. The probability that any other code-word looks jointly typical with the ~~received~~ received sequence is 2^{-nI} so the errors can be kept low if we do with fewer than 2^{nI} codewords.

Theorem 2.3.2 : Channel Coding Theorem

- 1) All rates below capacity C are achievable. Namely, if $R < C$, there exists a sequence of $(2^{nR}, n)$ codes with maximum probability of error $d^{(n)} \xrightarrow{n} 0$.
- 2) Conversely, any sequence of $(2^{nR}, n)$ codes with $d^{(n)} \xrightarrow{n} 0$ must have $R \leq C$.

Proof of part 1).

Each code \mathcal{E} is identified by a matrix $[x_i(w)]_{\substack{i=1,2,\dots,n \\ w=1,2,\dots,2^{nR}}}$

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

- Assume each matrix entry to be generated according to the probability $p(x)$

and independently :
$$p(x^{(n)}(w)) = \prod_{j=1}^n p(x_j(w))$$

- ~~the~~ The probability of a given code \mathcal{E} is then

$$\text{Prob}(\mathcal{E}) = \prod_{w=1}^{2^{nR}} \prod_{j=1}^m p(x_j | w)$$

(140)

- A word w is chosen from $\mathcal{W} = \{1, 2, \dots, 2^{nR}\}$ according to the uniform

distribution:
$$\text{Prob}(w) = 2^{-nR} \quad \forall w \in \mathcal{W}$$

- The code-word $x^{(n)}(w)$ from the randomly generated code \mathcal{E} is sent over the channel.

- The receiver receives a sequence $y^{(n)}$ with probability

$$P(y^{(n)} | x^{(n)}(w)) = \prod_{j=1}^m P(y_j | x_j(w))$$

- The receiver guesses the word sent by typical set decoding:

\hat{w} was sent if $\left\{ \begin{array}{l} x^{(n)}(\hat{w}) \text{ and } y^{(n)} \text{ are jointly typical} \\ \text{there is no other } \bar{w} \in \mathcal{W} \text{ such that } x^{(n)}(\bar{w}) \text{ and } y^{(n)} \\ \text{are jointly typical} \end{array} \right.$

- An error has occurred if either $x^{(n)}(\hat{w})$ and $y^{(n)}$ are not jointly typical or more $x^{(n)}(w)$ jointly typical with $y^{(n)}$ exist.

- Average probability of error: let \mathcal{E} denote the error occurring by wrong decoding. The average probability of error is:

$$\text{Prob}^{\text{av}}(\mathcal{E}) = \sum_{\mathcal{E}} \text{Prob}(\mathcal{E}) \frac{1}{2^{nR}} \sum_{w \in \mathcal{W}} d_w^{(n)}(\mathcal{E})$$

$$d_w^{(n)}(\mathcal{E}) = \sum_{y^{(n)}: g(y^{(n)}) \neq w} p(y^{(n)} | x_{\mathcal{E}}^{(n)}(w))$$

- Simple example of random coding: $\mathcal{W} = \{0, 1\}$; $\mathcal{C} = \{0, 1\}$; $m = 2$

$$\mathcal{E}_1 = \begin{cases} 0 \rightarrow 00 \\ 1 \rightarrow 00 \end{cases}; \quad \mathcal{E}_2 = \begin{cases} 0 \rightarrow 00 \\ 1 \rightarrow 01 \end{cases}; \quad \mathcal{E}_3 = \begin{cases} 0 \rightarrow 00 \\ 1 \rightarrow 10 \end{cases}; \quad \mathcal{E}_4 = \begin{cases} 0 \rightarrow 00 \\ 1 \rightarrow 11 \end{cases}$$

$$\mathcal{E}_5 = \begin{cases} 0 \rightarrow 01 \\ 1 \rightarrow 00 \end{cases}; \quad \mathcal{E}_6 = \begin{cases} 0 \rightarrow 01 \\ 1 \rightarrow 01 \end{cases}; \quad \mathcal{E}_7 = \begin{cases} 0 \rightarrow 01 \\ 1 \rightarrow 10 \end{cases}; \quad \mathcal{E}_8 = \begin{cases} 0 \rightarrow 01 \\ 1 \rightarrow 11 \end{cases}$$

$$\mathcal{E}_9 = \begin{cases} 0 \rightarrow 10 \\ 1 \rightarrow 00 \end{cases}; \quad \mathcal{E}_{10} = \begin{cases} 0 \rightarrow 10 \\ 1 \rightarrow 01 \end{cases}; \quad \mathcal{E}_{11} = \begin{cases} 0 \rightarrow 10 \\ 1 \rightarrow 10 \end{cases}; \quad \mathcal{E}_{12} = \begin{cases} 0 \rightarrow 10 \\ 1 \rightarrow 11 \end{cases}$$

$$\mathcal{E}_{13} = \begin{cases} 0 \rightarrow 11 \\ 1 \rightarrow 00 \end{cases}; \quad \mathcal{E}_{14} = \begin{cases} 0 \rightarrow 11 \\ 1 \rightarrow 01 \end{cases}; \quad \mathcal{E}_{15} = \begin{cases} 0 \rightarrow 11 \\ 1 \rightarrow 10 \end{cases}; \quad \mathcal{E}_{16} = \begin{cases} 0 \rightarrow 11 \\ 1 \rightarrow 11 \end{cases}$$

- Because of the random generation of the code-words, to each $w \in \mathcal{W}$ we associates all $x^{(n)} \in \mathcal{X}^{(n)}$. Then,

$$\sum_{\mathcal{E}} \text{Prob}(\mathcal{E}) d_w^{(n)}(\mathcal{E}) \quad \boxed{\text{does not}} \quad \text{depend on } w \quad \text{and}$$

$$\boxed{\text{Prob}^w(\mathcal{E}) = \sum_{\mathcal{E}} \text{Prob}(\mathcal{E}) \frac{1}{2^{nR}} \sum_{w \in \mathcal{W}} d_w^{(n)}(\mathcal{E}) = \sum_{\mathcal{E}} \text{Prob}(\mathcal{E}) d_1^{(n)}(\mathcal{E})}$$

- Let $E_w = \{ (x^{(n)}(w), \gamma^{(n)}) \in A_{\mathcal{E}}^{(n)} \}$ be the subset of pairs

where $\gamma^{(n)}$ is the string received when $x^{(n)}(w)$ is sent through the channel.

Then, $E_1^c = \{ (x^{(n)}(1), \gamma^{(n)}) \notin A_{\mathcal{E}}^{(n)} \}$ and $E_{w \neq 1}$ are the subsets from each

errors originate. Then,

$$\text{Prob}^w(\mathcal{E}) = \text{Prob}^w(E_1^c \cup \bigcup_{w \neq 1} E_w) \leq \text{Prob}^w(E_1^c) + \sum_{w=2}^{2^{nR}} \text{Prob}^w(E_w)$$

- By the random generation of the code words and the joint AEP

$$\boxed{\text{Prob}^w(E_1^c) \leq \epsilon}$$