

Channel coding theorem

Hangjun Cho.



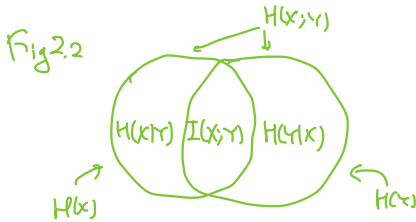
Claude Elwood Shannon
(April 30, 1916 – February 24, 2001)

7	Channel Capacity	183
7.1	Examples of Channel Capacity	184
7.1.1	Noiseless Binary Channel	184
7.1.2	Noisy Channel with Nonoverlapping Outputs	185
7.1.3	Noisy Typewriter	186
7.1.4	Binary Symmetric Channel	187
7.1.5	Binary Erasure Channel	188
7.2	Symmetric Channels	189
7.3	Properties of Channel Capacity	191
7.4	Preview of the Channel Coding Theorem	191
7.5	Definitions	192
7.6	Jointly Typical Sequences	195
7.7	Channel Coding Theorem	199
7.8	Zero-Error Codes	205
7.9	Fano's Inequality and the Converse to the Coding Theorem	206
7.10	Equality in the Converse to the Channel Coding Theorem	208
7.11	Hamming Codes	210
7.12	Feedback Capacity	216
7.13	Source-Channel Separation Theorem	218
	Summary	222
	Problems	223
	Historical Notes	240



Def: A discrete channel is a system consisting of an input X , output Y , $p(y|x)$.

Information channel capacity: $C = \max_{p(x)} I(X;Y)$

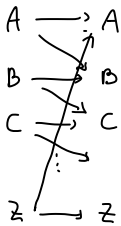


Mutual information.

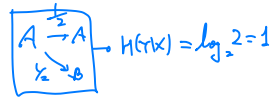
Thm 2.4.1.

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

Ex) D.1.3. Noisy Typewriter.



$$\begin{aligned} C &= \max [H(Y) - H(Y|X)] \\ &= \max H(Y) - 1 \\ &= \log 26 - 1 \\ &= \log 3 \end{aligned}$$



Thm 2.6.4.

Thm 2.6.4. $H(X) \leq \log |X|$. equality holds iff X has a uniform disc. over X .

§ 1.5 Definition for theorem.

Def. • A discrete channel: $(\mathcal{X}, p(y|x), \mathcal{Y})$

finite set \rightarrow prob. mass fun.

• The n th extension of DMC: \leftarrow discrete memoryless channel

$(\mathcal{X}^n, p(y^n|x^n), \mathcal{Y}^n)$ where $p(y_k|x^k, y^{k-1}) = p(y_k|x_k)$

\downarrow
 $y^n = (y_1, y_2, \dots, y_n)$: vector notation.

• A channel without feedback if

$$p(x_k|x^{k+1}, y^{k+1}) = p(x_k|x^{k+1})$$

Def. An (M, n) code for $(\mathcal{X}, p(y|x), \mathcal{Y})$:

1. An index set: $\{1, 2, \dots, M\} =: I^M$

codewords \in codebook.

2. encoding fun: $x^n: I^M \rightarrow \mathcal{X}^n$ $i \mapsto x^n(i)$

3. decoding fun.: $g: \mathcal{Y}^n \rightarrow I^M$

Def. • Conditional prob. error. $\lambda_i = P(g(Y^n) \neq i | X^n = x^n(i))$

• The maximal prob. of error. $\lambda^{(n)} = \max_i \lambda_i$

• The average prob. of error. $P_e^{(n)} = \frac{1}{M} \sum_{i=1}^M \lambda_i$.

- If the index W is chosen uniformly over I^M ,

$$P_e^{(n)} = P(W \neq g(Y^n))$$

Def. The rate R of (M, n) Code: $R = \frac{\log M}{n}$

• A rate R is achievable $\rightarrow 2^{nR} \in \mathbb{Z}$

if \exists a seq. of $(2^{nR}, n)$ Code, s.t. $\lambda^{(n)} \xrightarrow{n \rightarrow \infty} 0$

• The capacity of a DMC is the supremum of all achievable rates.
(achievable capacity)

typical set

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

$$A_\epsilon^{(n)} := \{(x^n, y^n) \in X^n \times Y^n :$$

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

Where

$$p(x^n, y^n) = \prod_{z=1}^n p(x_z, y_z)$$

$$\quad \quad \quad \uparrow$$

$$\quad \quad \quad p(x_i) p(y_i)$$

Thm 7.6.1 (X^n, Y^n) : sequences drawn i.i.d. ad $p(x^n, y^n) = \prod_{z=1}^n p(x_z, y_z)$

1. $P((X^n, Y^n) \in A_\epsilon^{(n)}) \rightarrow 1$ 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 $P((\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X; Y) - 3\epsilon)}$

For $n \gg 1$, $P((\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon^{(n)}) \geq (1 - \epsilon) 2^{-n(I(X; Y) + 3\epsilon)}$

Thm. (Channel Coding Theorem.) For a DMC,

" $R < C \Rightarrow \exists (2^{nR}, n)$ code with $d^{(n)} \rightarrow 0$ as $n \rightarrow \infty$ "

Conversely, if \exists a seq. of $(2^{nR}, n)$ codes with $d^{(n)} \rightarrow 0$, $R \leq C$

pt. (Achievability) Fix $p(x)$. ^{distribution.} generate $(2^{nR}, n)$ code at random and

$$p(x^n) = \prod_{i=1}^n p(x_i)$$

Which gives a codebook.

$$C = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(2^{nR})} & x_2^{(2^{nR})} & \dots & x_n^{(2^{nR})} \end{bmatrix} \rightarrow \text{the 1st merge goes to } n\text{-vector.}$$

Note that $P(C) = \prod_{\omega=1}^{2^{nR}} p(x^n(\omega)) = \prod_{\omega=1}^{2^{nR}} \prod_{i=1}^n p(x_i(\omega))$

seq. of events:

1. C is generated and $p(x)$.

2. C is revealed to encoder and decoder
 \downarrow
 know $p(y|x)$

3. W is chosen and a uniform dist: $P(W=\omega) = 2^{-nR}$, $\omega=1, 2, \dots, 2^{nR}$.

4. $X^n(\omega)$ is sent.

5. decoder receives a seq. Y^n and $P(y^n|x^n(\omega)) = \prod_{i=1}^n p(y_i|x_i(\omega))$.

b. The receiver decodes the index:

$$\hat{w} = g(\gamma^n) \text{ if } \begin{cases} (x^n(\hat{w}), \gamma^n) \text{ is jointly typical.} \\ \nexists w' (\neq \hat{w}) \text{ s.t. } (x^n(w'), \gamma^n) \in A_{\epsilon}^{(n)}. \end{cases}$$

error, when no such \hat{w} exists or there is more than one such.

? Decoding error if $\hat{w} \neq w$. $\mathcal{E} = \{\hat{w} \neq w\}$

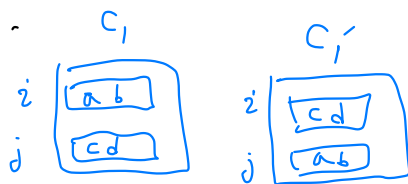
(Study of prob. of error)

Calculate the ave. prob. of error, averaged over all codes.

$$\begin{aligned} P(\mathcal{E}) &= \sum_C P(C) P_{\mathcal{E}}^{(n)}(C) && \underline{P_{\mathcal{E}}^{(n)} = P(w \neq g(\gamma^n))} \\ &= \sum_C P(C) \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(C) \\ &= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_C P(C) \lambda_w(C). \end{aligned}$$

By symmetry of code construction, $\sum_C P(C) \lambda_w(C)$ does not

depend on w .



WOLG ass, $w=1$ was sent.

$$P(\mathcal{E}) = \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_C P(c) \lambda_1(c)$$

$$= \sum_C P(c) \lambda_1(c) = P(\mathcal{E} | W=1)$$

Define $E_i = \{(X^n(i), Y^n) \in A_\epsilon^n\}$ for $i = \{1, 2, \dots, 2^{nR}\}$

$W=1$

- E_1^c $(X^n(1), Y^n)$ is not jointly typical.

An error occurs:

- $E_2 \cup \dots \cup E_{2^{nR}}$ Wrong codeword is jointly typical.

Y^n is the result of sending the first code word $X^n(1)$.

$$P(\mathcal{E} | W=1) = P(E_1^c \cup E_2 \cup \dots \cup E_{2^{nR}})$$

$$\leq P(E_1^c) + \sum_{i=2}^{2^{nR}} P(E_i)$$

$$P(A \cup B) \leq P(A) + P(B)$$

$$\leq \underbrace{\epsilon}_{(1)} + \underbrace{\sum_{i=2}^{2^{nR}} 2^{-n(I(X;Y) - 3\epsilon)}}_{(2)} \quad \text{by joint AEP.}$$

(1)

(2) $X^n(1)$ and $X^n(i)$: indep. for $i \neq 1$.

$$= \epsilon + (2^{nR} - 1) \cdot 2^{-n(I(X;Y) - 3\epsilon)}$$

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

(if $R < I(X;Y) - 3\epsilon$ and $n \gg 1$.)

\therefore For $R < I(X;Y)$, $\forall \epsilon, \exists n$ s.t. $P(\mathcal{E}) \leq 2\epsilon$.

(Finding a codebook we want)

1. Choose $p(x)$ that maximize $I(X;Y)$, and call it $p^*(x)$.

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

2. $\exists C^*$ s.t. $P_e^{(n)}(C^*) \leq 2\epsilon$ ↔ $P(\epsilon) = \sum_c P(c) P_e^{(n)}$

3. Aim: $\lambda^{(n)}(C^*) \rightarrow 0$

$$2\epsilon \geq \frac{1}{2^{nR}} \underbrace{\sum_{i=1}^{2^{nR}} \lambda_i(C^*)}_{\lambda^*} = \frac{1}{2^{nR}} \left(\underbrace{\sum \lambda_i}_{\substack{\lambda^* > \lambda_i \\ \text{best half}}} + \sum \lambda_i \right)$$

$\lambda^* = \lambda_{i_2}$

$\lambda^* < \lambda_i$
worst half.

For best half, $\lambda_i \leq 4\epsilon$ for all i in best half.

If not $\lambda^* \geq 4\epsilon$. $\frac{1}{2^{nR}} \sum_{\lambda^* < \lambda_i} \lambda_i > \frac{1}{2} \cdot 4\epsilon = 2\epsilon$ *

Throw away the worst half of codeword in C^* .

And take 2^{nR-1} codewords (best half of C^*)

New Rate $\frac{1}{n} \cdot \log(2^{nR-1}) = R - \frac{1}{n}$

and $\lambda^{(n)} = \max \lambda_i \leq 4\epsilon$

///

(Converse) " if \exists a seq. of $(2^{nR}, n)$ codes with $d^{(n)} \rightarrow 0$, $R \leq C$ "

(special case)

Supp. $P_e^{(n)} = 0$, Ass: W : uniformly dist. over $\{1, 2, \dots, 2^{nR}\}$

$$nR \stackrel{?}{=} H(W)$$

$$\begin{aligned} &= H(W|Y^n) + I(W; Y^n) \\ &\quad \text{Fig 2.2} \end{aligned}$$

($\because P_e^{(n)} = 0$; $\hat{W} = g(Y^n)$)

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

: Data processing - $W \rightarrow X^n \rightarrow Y^n$

$$\leq \sum_{i=1}^n I(X_i; Y_i)$$

1.9.2

$$\leq nC$$

def. of C

$\therefore R \leq C$

7.9.1 Lemma (Fano's inequality)

For a DMC, W : uniformly distributed. $H(X^n|Y^n) \leq 1 + P_e^{(n)} \cdot nR$

7.9.2. For DMC, $I(X^n; Y^n) \leq nC$

pf. $I(X^n; Y^n) = H(Y^n) - H(Y^n|X^n)$

$$= H(Y^n) - \sum H(Y_i | Y_1, \dots, Y_{i-1}, X^n)$$

$$= H(Y^n) - \sum H(Y_i | X_i)$$

← DMC

$$\leq \sum H(Y_i) - \sum H(Y_i | X_i)$$

$$= \sum I(X_i; Y_i) \leq nC$$

* Converse of the Coding theorem.

" if \exists a seq. of $(2^{nR}, n)$ codes with $d^{(n)} \rightarrow 0$, $R \leq C$ "

Pf. Since $p_e^{(n)} \leq \lambda^{(n)}$, $\lambda^{(n)} \rightarrow 0$ implies $p_e^{(n)} \rightarrow 0$.

Ans. ω : uniform dist with $p_e^{(n)} = P(\hat{\omega} \neq \omega)$.

$$nR = H(\omega) = H(\omega | \gamma^n) + I(\omega; \gamma^n)$$

$$\leq H(\omega | \gamma^n) + I(x^n(\omega); \gamma^n)$$

$$\omega \rightarrow x^n \rightarrow \gamma^n$$

$$\leq (1 + p_n^{(e)}) nR + I(x^n; \gamma^n)$$

$$\leq (1 + p_n^{(e)}) nR + nC$$

Dividing by n , $R \leq \frac{1}{n} + p_n^{(e)} R + C$.

As $n \rightarrow \infty$, $R \leq C$.

((/)

Equation in the converse

We consider the special case $\frac{P_n^{(e)}}{n} = 0$: zero error code

$$nR \stackrel{?}{=} H(W)$$

$W \rightarrow \text{code}$

$$\stackrel{?}{=} H(W|W') + I(W; W')$$

$\stackrel{0}{=} \text{ (}\because P_n^{(e)} = 0\text{)}$

Fig 2.2. X^n must be distinct.

$$\stackrel{?}{\leq} I(X^n(W); Y^n)$$

data processing req. : $W \rightarrow X^n \rightarrow Y^n \rightarrow \hat{W}$.

$$= H(Y^n) - H(Y^n | X^n)$$

$$I(Y^n; X^n(W) | W) = 0 = I(X^n; Y^n | W)$$

$$= H(Y^n) - \sum_{i=1}^n H(Y_i | X_i) \leftarrow \text{DMC}$$

$$\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \leftarrow Y_i: \text{indep.}$$

$$= \sum_{i=1}^n I(X_i; Y_i) \leftarrow \text{by def.}$$

$$\leq nC$$

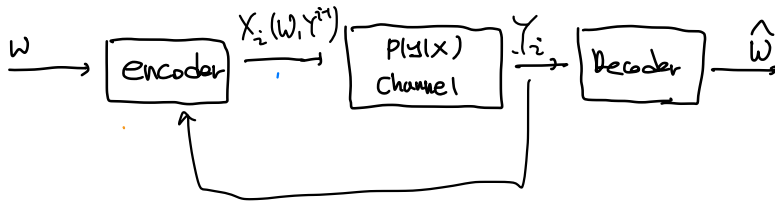
$\leftarrow \text{by def.}$

$$X_i \sim p^*(x_i)$$

$$\parallel$$

$$\text{argmax } I(X; Y)$$

§ 11.12 Feedback capacity



- $(2^{nR}, n)$ feedback code: a sequence of mappings $x_i(w, Y^{i-1})$,
 seq of decodg fns, $g: Y^n \rightarrow \{1, 2, \dots, 2^{nR}\}$

The capacity with feedback, C_{FB} , of a DMC

: a supremum of all rates achievable by feedback codes.

$$1.12.1. \quad C_{FB} = C = \max_{p(x)} I(X; Y)$$

pt. A non feedback code is a special code of a feedback code

$$C_{FB} \geq C.$$

$$\text{Aim: } C_{FB} \leq C.$$

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

$$(*) \quad \leq (1 + P_e^{(n)}) nR + I(W;\hat{W})$$

$$\leq 1 + P_e^{(n)} nR + \underline{I(W; Y^n)} \leftarrow W \rightarrow Y^n \rightarrow \hat{W}.$$

$$I(\omega; \gamma^n) = H(\gamma^n) - H(\gamma^n | \omega)$$

$$= H(\gamma^n) - \sum_{i=1}^n H(\gamma_i | \gamma_1, \dots, \gamma_{i-1}, \omega)$$

$$\approx H(\gamma^n) - \sum_{i=1}^n H(\gamma_i | \gamma_1, \gamma_2, \dots, \gamma_{i-1}, \omega, X_i) \quad \text{def } X_i = X_i(\omega, \gamma^{i-1})$$

$$= H(\gamma^n) - \sum_{i=1}^n H(\gamma_i | X_i)$$

DMC

$$\leq \sum_{i=1}^n H(\gamma_i) - \sum_{i=1}^n H(\gamma_i | X_i)$$

Chain rule + 2.6.1

$$= \sum_{i=1}^n I(X_i; \gamma_i)$$

def

$$H(X|Y) \leq H(X)$$

$$\leq nC.$$

def

from (*) , $nR \leq P_e^{(n)} nR + 1 + nC.$

Dividing n , $n \rightarrow \infty$, $R \leq C.$

$$C_{FB} \leq C.$$

//

$$3.1.2-2: P(A_{\frac{\epsilon}{2}}^{(n)}) > 1 - \epsilon.$$

From ch. coding thm, we can transmit the indices with the arbitrary small prob

$$\text{of error if } H(U) + \epsilon = R < C$$

Can we construct V^n to agree with the transmitted sequence with high prob.

(why, input size $= n$ = output size)

Precisely,

$$\begin{aligned} P(V^n \neq \hat{V}^n) &\leq P(V^n \notin A_{\frac{\epsilon}{2}}^{(n)}) + P(g(\hat{V}^n) \neq V^n | V^n \in A_{\frac{\epsilon}{2}}^{(n)}) \\ &\leq \epsilon + \epsilon = 2\epsilon \end{aligned}$$

$$\therefore H(U) < C.$$

Converse: " $P(V^n \neq \hat{V}^n) \rightarrow 0 \Rightarrow H(U) \leq C$ "

Given the following encoder and decoder

$$x^n(v^n): U^n \rightarrow \mathcal{X}^n$$

$$g_n(\hat{v}^n): \mathcal{Y}^n \rightarrow U^n$$

By Fano's ineq.

$$H(V^n | \hat{V}^n) \leq (1 + P(\hat{V}^n \neq V^n)) \log |U^n| = (1 + P(\hat{V}^n \neq V^n)) n \log |U|$$

||
||ⁿ (entropy seq.)

$$H(U) \leq \frac{H(X_1, \dots, X_n)}{n} \quad \text{Ch 4. } H(U) = \lim_{n \rightarrow \infty} \frac{H(U_1, \dots, U_n)}{n}, \quad \neq \text{dec. seq.}$$

$$= \frac{H(V^n)}{n}$$

$$= \frac{1}{n} H(V^n | \hat{V}^n) + \frac{1}{n} I(V^n; \hat{V}^n)$$

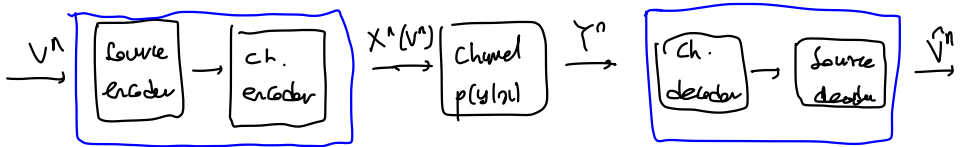
$$\leq \frac{1}{n} (1 + P(\hat{V}^n \neq V^n)) n \log |U| + \frac{1}{n} I(V^n; \hat{V}^n) \quad \text{Fano}$$

$$\leq \frac{1}{n} (1 + P(\hat{V}^n \neq V^n)) n \log |U| + \frac{1}{n} I(X^n; Y^n) \quad \text{Data processing}$$

$$\leq \frac{1}{n} (1 + P(\hat{V}^n \neq V^n)) n \log |U| + C \quad \text{Memoryless}$$

As $n \rightarrow \infty$, $\therefore H(U) \leq C$

(/)



Summary, Data Compression thm: a consequence of the AEP.

Data transmission thm: a consequence of the joint AEP.