

"Classical" Information Theory  
⇕  
Shannon Information Theory

Claude Shannon "A mathematical theory of communication" 1948

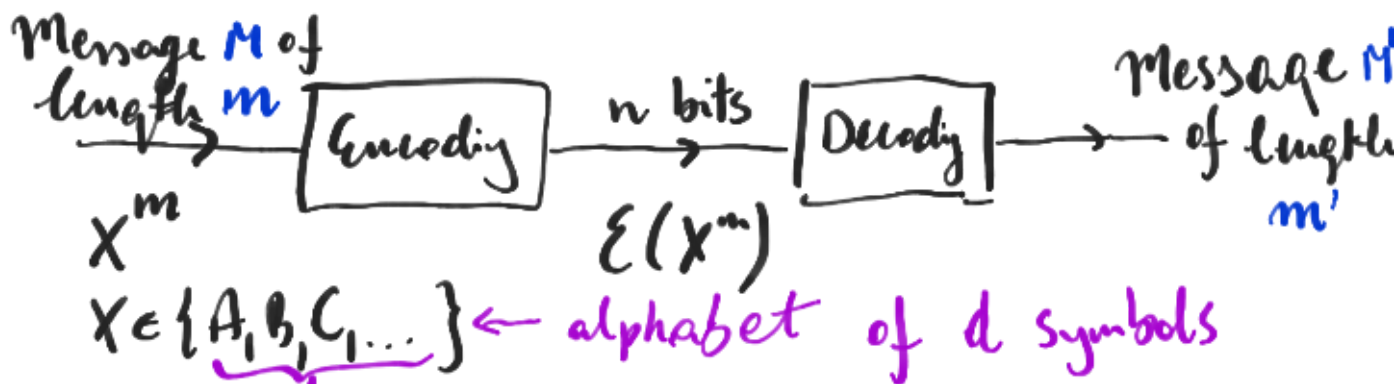
MIT → Princeton → Bell Labs → MIT

- (a) Data compression: (Shannon) Coding Theorem
- (b) Communication: (Shannon) Channel Capacity
- (c) Cryptography:
  - Security of one-time pad (OTP)
  - (Shannon) Secret Capacity

Book: Cover's Thomas: "Elements of Information Theory"

(a) Data Compression

Motivation: File compression {ZIP, ARJ, RAR, TAR, 7z, ...}



Q: How many bits (at least) per symbol?  
for lossless compression  $M=M'$

Assume:  $X$  are distributed according to

some PDF but i.i.d.!

[N.B. for "text" we can use also correlation e.g. "Ala ma kota" ... ]  
vowel after consonant

### Example

alphabet = {'a', 'b', 'c', 'd'} symbols (d=4)

|      |     |     |     |     |
|------|-----|-----|-----|-----|
|      | 'a' | 'b' | 'c' | 'd' |
| p(x) | 1/2 | 1/4 | 1/8 | 1/8 |

codewords

"Simplest" encoding: 'a' = 00, 'b' = 01, 'c' = 10, 'd' = 11  
with 2 bits

twice as many:  $\frac{\# \text{bits}}{\text{msg. length}} = \frac{n}{m} = \frac{m \cdot 2}{m} = 2$  ← average codeword length (in bits)

"Smarter" encoding: 'a' = 0, 'b' = 10, 'c' = 110, 'd' = 111

"frequencies" = "probabilities"

unambiguous decoding!

1000 b de, 11000 c ee, 111000 d haa

$$\frac{n(m)}{m} \underset{m \rightarrow \infty}{=} \frac{1}{m} \left( \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 \right) = 1.75$$

Q Shannon: What is the optimal  $\frac{n(m)}{m}$  as  $m \rightarrow \infty$ ?

Answer: (Shannon Coding Thm.)

$$\lim_{m \rightarrow \infty} \left\{ \frac{n(m)}{m} \right\} = H(x)$$

Shannon Entropy

$$H(x) = - \sum_{x=1}^d p(x) \log_2 p(x) \quad \left\{ \begin{array}{l} \log \equiv \log_{10} \\ \lg \equiv \log_2 \end{array} \right.$$

→ measure of randomness of the random variable  $X$  }  $\ln = \log_e$   
*d-outcomes (symbols in alphabet)*  
 $X \sim p(x_i) \quad p(x_i) = \delta_{ik} \quad 0 \leq H(X) \leq \lg d \quad p(x_i) = \frac{1}{d}$

Example:  $H(X) = -\frac{1}{2} \lg \frac{1}{2} - \frac{1}{4} \lg \frac{1}{4} - \frac{1}{8} \lg \frac{1}{8} - \frac{1}{8} \lg \frac{1}{8} =$   
 $= \frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{3}{8} = 1 \frac{3}{4} = 1.75$

### Intuition of the proof

$X^m$ : as sequence length  $m \rightarrow \infty$   
the symbol  $X=x$  expected  
to appear  $m p(x)$  times

⇒ "typical" sequences

[consequence of "Law of Large Numbers"]  
i.e.  $X_m = \frac{\sum_{i=1}^m X_i}{m} \rightarrow E[X]$

• Probability of obtaining a typical sequence:

$$P_{\text{typ}} = P(X_0 X_1 \dots X_m \in T) = \prod_{x=1}^d p(x)^{m p(x)}$$

$$= 2^{\lg \prod_x p(x)^{m p(x)}} = 2^{\sum_x \lg p(x)^{m p(x)}} = 2^{m \sum_x p(x) \lg p(x)}$$

• # all sequences (in bits)  $d^m = 2^{\lg d^m} = 2^{\underline{m \lg d}}$

• # typ =  $\frac{1}{n^m} = 2^{-m \sum_x p(x) \lg p(x)} = 2^{\underline{m H(X)}}$

" seq. type

$n$

$\Rightarrow$  rather than mlgd we need  $m H(x)$  bits to label all typical sequences

$$\lim_{m \rightarrow \infty} \frac{\binom{n}{m}}{m} = \frac{m H(x)}{m} = H(x)$$

$H(x)$  is the optimal compression rate  $\Leftrightarrow$  Shannon Coding Theorem

How to construct the code?  $\Rightarrow$  Huffman Coding

### Properties of entropies

Shannon Entropy:  $X \sim p(x)$   $H(x) = - \sum_x p(x) \lg p(x)$

Joint Entropy:  $X, Y \sim p(x, y)$   $H(X, Y) = - \sum_{x, y} p(x, y) \lg p(x, y)$

Conditional Entropy:  $X \sim p(x|y)$  given  $y$   
 $H(X|Y=y) = - \sum_x p(x|y) \lg p(x|y)$

$$\underline{H(X|Y)} = \sum_y p(y) H(X|Y=y) =$$

$$= - \sum_{x, y} p(x, y) \lg p(x|y) = - \sum_{x, y} p(x, y) \lg \frac{p(x, y)}{p(y)} =$$

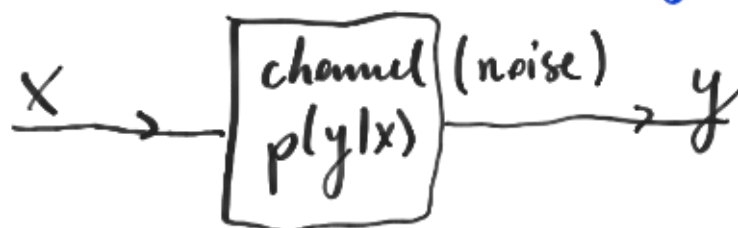
$$\sum_y p(y) \left( - \sum_x p(x|y) \lg p(x|y) \right) =$$

$$\begin{aligned}
&= -\sum_{x,y} p(x,y) \lg p(x,y) + \sum_{x,y} p(x,y) \lg p(y) - \\
&= -\sum_{x,y} p(x,y) \lg p(x,y) - \left( -\sum_y p(y) \lg p(y) \right) = \\
&= H(X,Y) - H(Y)
\end{aligned}$$



- (i)  $0 \leq H(X) \leq \lg d$  ← uniform PDF
- (ii)  $H(X,Y) \leq H(X) + H(Y)$  ↑ indep<sup>nt</sup> variables
- (iii) chain rule:
 
$$H(X,Y) = H(X|Y) + H(Y)$$
- (iv)  $H(X|Y) \leq H(X)$
- (v)  $H(X,Y|Z) = H(X|YZ) + H(Y|Z)$

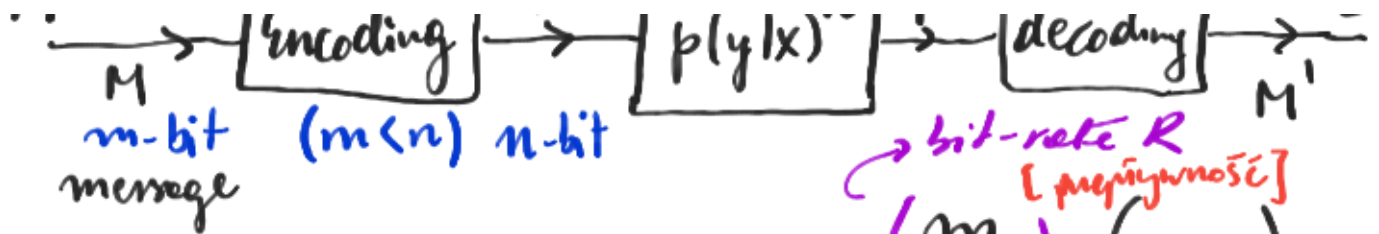
(b) Communication: Shannon Channel (Theorem)  
Capacity



How to send a message reliably over a channel?

Encode m-bit message into n-bit codeword  
 $n > m$  (use redundant data)





Channel Capacity:  $C = \lim_{m \rightarrow \infty} \left( \frac{m}{n} \right) \left( \leq 1 \right)$

bit-rate  $R$   
[propriety]

### Examples

1) error-correction codes eg. CD disc:  $142 \rightarrow 588$   
 $m$   $n$   
 logical bits physical bits  
 $C \approx \frac{1}{3}$   $\leftarrow \approx 1 : \approx 3$

2) Transmission via an optical link eg. optical fibre  
 size: 800MHz 2.46G  
 [bits/s]  $C = B \lg(1 + \text{SNR})$  bandwidth

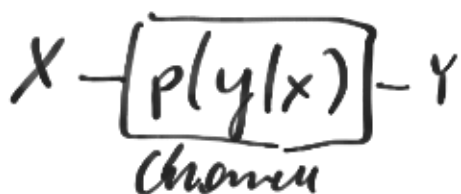
(SNR) Signal to noise ratio:  $\frac{\text{Power of signal}}{\text{Power of noise}}$  in dB  
 $10 \log_{10} \text{SNR}$   
 $\{1000 \rightarrow 30 \text{ dB}\}$

Answer:  
Shannon Channel Theorem

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

(Shannon Mutual Information)

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



How much about "X" we know from knowing "Y"?

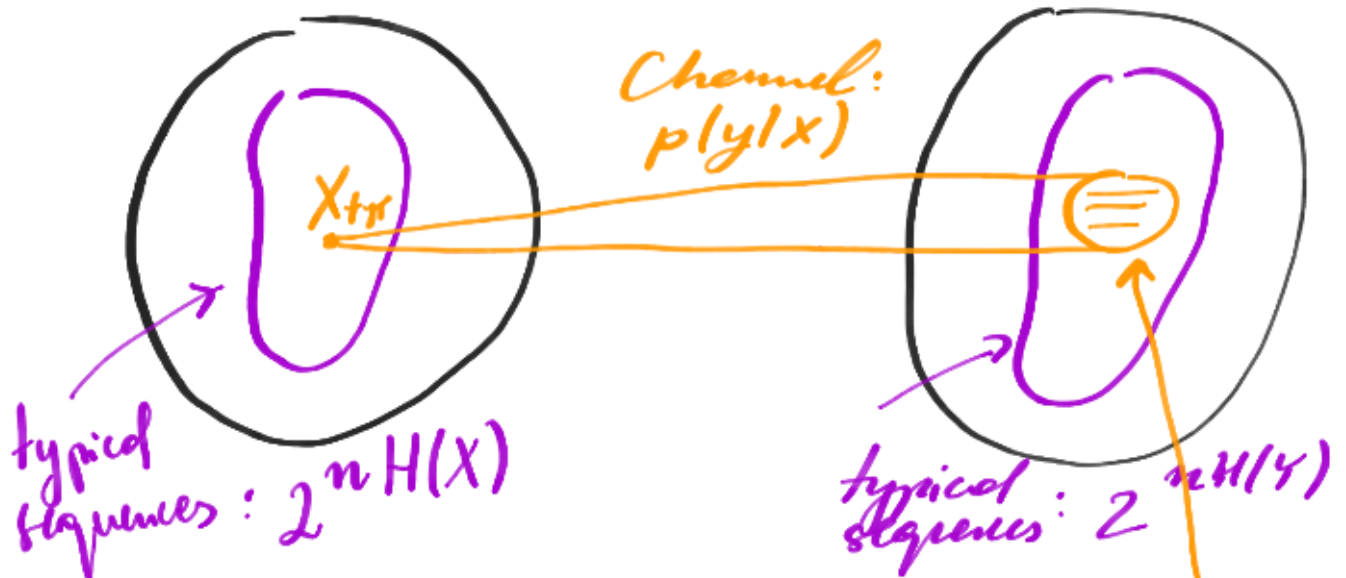
symmetric.

$$I(X:Y) = H(X) + H(Y) - H(X,Y)$$

## Intuition of the proof

all sequences  
of  $X^n$   
[n bits :  $|X^n| = 2^n$ ]

all sequences  
of  $Y^n$   
[ $|Y^n| = 2^n$ ]



• given  $X=x$  there will be  $2^{H(Y|X=x)}$   
(n=1) typical (fixed)  $Y$  it can lead to.

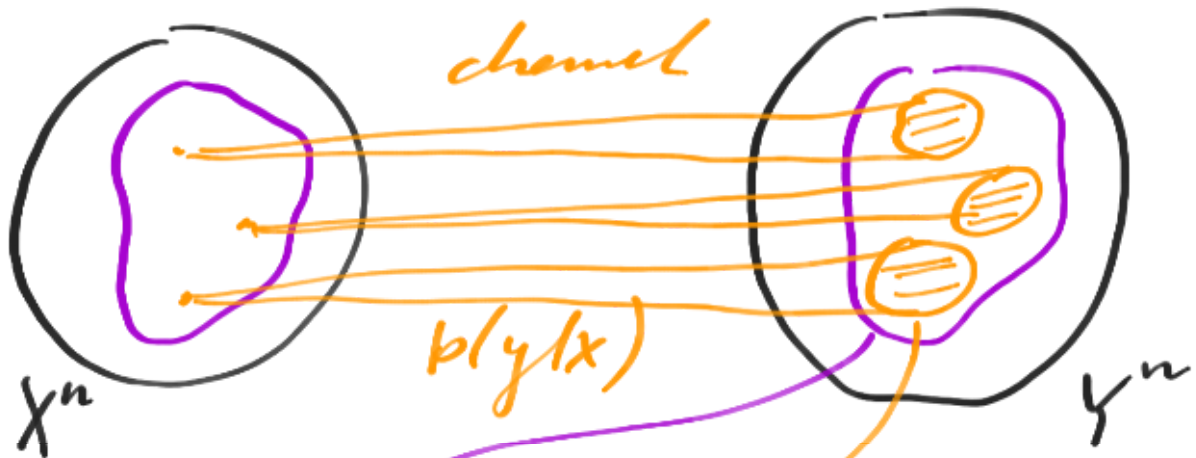
• ok. but there are  $\approx n p(x)$   
(n)  $X=x$  in a typical sequences

$X_{typ} \Rightarrow$  leads to

$$\prod_x (2^{H(Y|x)})^{n p(x)} = 2^{n \sum_x p(x) H(Y|x)} = 2^{n H(Y|X)}$$

How many different input (typical)

sequences can be distinguished from typical output sequences?

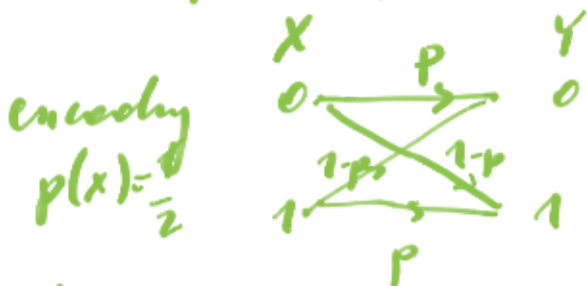


$$\frac{2^{n H(Y)}}{2^{n H(Y|X)}} = 2^{n \underbrace{I(X:Y)}_{\text{bit rate } R}}$$

Still, freedom to optimise over encodings

$$X \sim p(x) \Rightarrow C = \max_{p(x)} I(X:Y)$$

Example (A)



$$I(X:Y) = 1 - h(p)$$

$\{ h(p) = -p \log p - (1-p) \log (1-p) \}$  binary entropy



$$\Rightarrow \begin{aligned} p=1 \text{ (100\%)} & \quad I(X:Y)=1 \\ p=\frac{1}{2} \text{ (50\%)} & \quad I(X:Y)=0 \end{aligned}$$

(no noise) perfectly, completely



U.S. History / Chapter 10 / 10.1 / 10.1.1 / 10.1.1.1

Last modified: 16:36