

ENTROPY RATE of a sequence of random variables

$\{X_t\}_{t \in \mathbb{N}}$ joint distrib. $\underline{P}_N(x_1, \dots, x_N)$
for N variables

marginal distrib. over $A \subset \{1, \dots, N\}$; $x_{-A} = \{x_i, i \in A\}$
(obtained by summing over $x \in \bar{A} = \{1, \dots, N\} \setminus A$)

$$\underline{P}(x_{-A}) = \sum_{x_{\bar{A}}} \underline{P}_N(x_1, \dots, x_N)$$

Simple examples:

$$\underline{P}_N(x_1, \dots, x_N) = \prod_1^N p(x_i) \quad (\text{i.i.d.})$$

Markov chains $\underline{P}_N(x_1, \dots, x_N) = p_1(x_1) \prod_1^{N-1} w(x_t \rightarrow x_{t+1})$
initial state transition probabilities
 $\sum_{y \in X} w(x \rightarrow y) = 1$
 $\forall x \in X$

ENTROPY RATE

for $X_{-N} = \{X_t\}_{t=1, \dots, N}$

$$h_X = \lim_{N \rightarrow \infty} H_{X_{-N}} / N$$

(entropy grows linearly with the number of variables)

- i.i.d. with $p(x)$ $\rightarrow h_X = H(p) = - \sum_x p(x) \log_2 [p(x)]$

- Markov chain with $p^*(x)$ stationary distr.

$$p^*(x) = \lim_{t \rightarrow \infty} p(x_t) \quad \text{empirical of } X_t$$

$$h_x = - \sum_{x, y \in X} p^*(x) w(x \rightarrow y) \log_2 [w(x \rightarrow y)]$$

Correlated variables X, Y

(conditional entropy & mutual information
joint prob. distrib. $p(x, y)$

conditional prob. $p(y/x)$

Bayes theorem $\rightarrow p(y/x) = \frac{p(x, y)}{p(x)}$

$$p(x) = \sum_y p(x, y)$$

X, Y are independent $\rightarrow p(y/x) = p(y)$

X, Y are correlated \rightarrow measure of correlation?

Conditional entropy: Shannon entropy of the conditional prob. (over y) averaged over x

$$H_{Y/X} = - \sum_x p(x) \sum_y p(y/x) \log_2 p(y/x)$$

joint entropy $H_{X,Y} = - \sum_{x,y} p(x,y) \log_2 [p(x,y)]$

CHAIN RULE

$$H_{X,Y} = H_X + H_{Y/X}$$

$$= H_Y + H_{X/Y}$$

X, Y are independent $\rightarrow H_{Y/X} = H_Y$

X, Y are completely correlated

$$\rightarrow H_{X,Y} = H_X + H_Y$$

$Y = f(X) \rightarrow H_{Y/X} = 0 \quad (H_{X/Y} = 0)$

$$\Rightarrow H_{X,Y} = H_X = H_Y$$

MUTUAL INFORMATION

$$I_{X,Y} = \sum_{x,y} p(x,y) \log_2 \left[\frac{p(x,y)}{p(x)p(y)} \right]$$

$$= D(p_{X,Y} \parallel p_X p_Y) = -H_{X,Y} + H_X + H_Y$$

$$I_{X,Y} = H_Y - H_{Y/X} = H_X - H_{X/Y}$$

$I_{X,Y}$ is symmetric for $X \leftrightarrow Y$

$$I_{X,Y} \geq 0$$

- X, Y independent $(\Rightarrow) I_{X,Y} = 0$

- complete correlation
($X = f(Y)$) $(\Rightarrow) I_{X,Y} = H_Y = H_X$
 $= H_{X,Y}$

Mutual information is degraded
when data are transmitted

DATA-PROCESSING INEQUALITY:

Markov chain $X \rightarrow Y \rightarrow Z$

$$p(x, y, z) = p_a(x) w_2(x \rightarrow y) w_3(y \rightarrow z)$$

$$I_{X,Z} \leq I_{X,Y}$$

$$I_{X, f(Y)} \leq I_{X,Y}$$

This holds also if $Z = f(Y)$

DATA COMPRESSION

Source (input data) $\underline{X} = \{X_1, \dots, X_N\}$ X_i are

random variables

how to store the info in \underline{X} efficiently?

optimal compression \rightarrow entropy rate of \underline{X}

• Source code: mapping $w: X^N \rightarrow \{0, 1\}^*$

• codeword: any binary string
set of binary strings of arbitrary length

x_1, \dots, x_N part of a computer stream

Coding in 3 steps: (compression)

- 1) break stream into blocks of length N
- 2) each block encoded through mapping w
- 3) codewords are glued together to form the output stream \rightarrow compressed data

$\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(t)}$

$w(\underline{x}^{(1)}), w(\underline{x}^{(2)}), \dots, w(\underline{x}^{(t)})$

Careful: "padding" the output stream into codewords in a unique way

for any concatenation of codewords!

w uniquely decodable code

strange condition \rightarrow instantaneous code

$\forall \underline{x}, \underline{x}' \quad w(\underline{x})$ is not a prefix of $w(\underline{x}')$

(decoding can be done on the spot)

Quality of compression

length of the codewords:

length of the string $w(x)$

$$L(w) = \sum_{\underline{x} \in X^N} p(\underline{x}) \ell_w(\underline{x})$$

sum over all possible codewords

"Best code" = code with shortest $L(w)$

THEOREM

L_N^* shortest average length

(instantaneous codes) for $\underline{X} = \{x_1, \dots, x_N\}$

$H_{\underline{X}}$ = Shannon entropy of \underline{X}

$$(1) \forall N \geq 1 \quad H_{\underline{X}} \leq L_N^* \leq H_{\underline{X}} + 1$$

(2) if the source \underline{X} has a finite entropy rate

$$h = \lim_{N \rightarrow \infty} H_{\underline{X}}/N \Rightarrow h = \lim_{N \rightarrow \infty} L_N^*/N$$

Shannon Entropy of source sets the optimal compression that can be achieved

(1) is the difficult part in the proof:

instantaneous code \Rightarrow codewords not too short

instantaneous codes $(= \sum_{\underline{x}} 2^{-l_w(\underline{x})} \leq 1$

minimize $L(w)$ $\xrightarrow{\underline{x}}$ Kraft's inequality

subject to the constraint $\sum_{\underline{x}} 2^{-l_w(\underline{x})} = 1$

$$\left\{ \begin{array}{l} l_w(\underline{x}) = -\log_2 p(\underline{x}) \\ \left[\begin{array}{l} l_w^*(\underline{x}) = (-\log_2 p(\underline{x})) \\ \text{(integer part)} \end{array} \right] \end{array} \right.$$

optimal
solution

\rightarrow SHANNON'S CODE

for $N \gg 1$

NOT optimal for finite N ($p(\underline{x}) \ll 1$ $\rightarrow l_w(\underline{x}) \gg 1$)

Shannon code is impractical (for $N \gg 1$)
from a computational perspective

• build the code once and for all:

\rightarrow requires $\mathcal{O}(|X|^N)$ memory

• reconstruct the code every time

\rightarrow it takes $\mathcal{O}(|X|^N)$ CPU operations

DATA TRANSMISSION

may be degraded
by noise

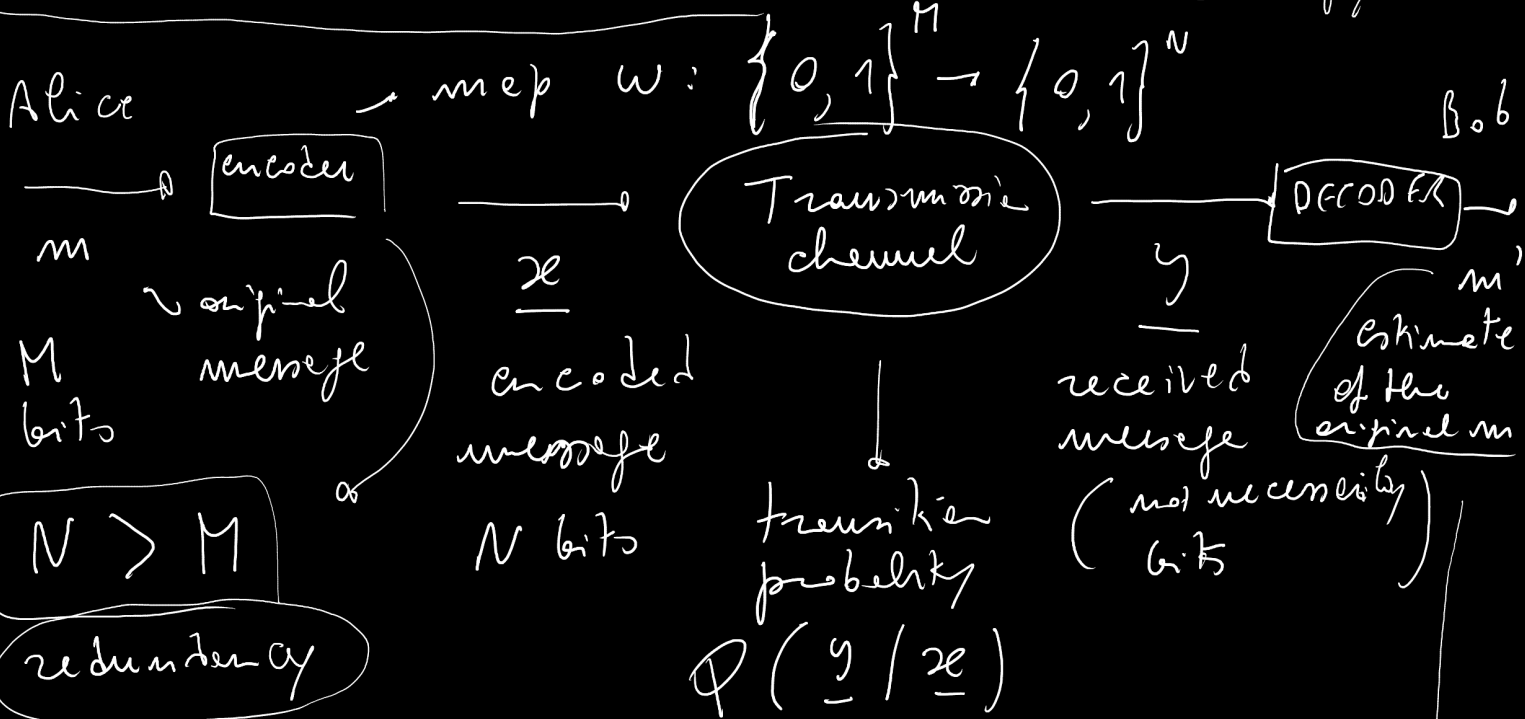
to counter transmission noise \rightarrow add redundancy
in source coding

(compression due to transmission)

Key result: CHANNEL CODING THEOREM (Shannon 1948)

Level of redundancy needed = maximal level of noise that can be tolerated

COMMUNICATION SYSTEM



$N =$ block length

$\underline{y} = \{y_1, \dots, y_N\}$ from our alphabet $\mathcal{Y}; \in \mathcal{Y}$

decoding map $d(\underline{y}) = m'$

noise $\rightarrow P(\underline{y} | \underline{x})$

Memoryless channels: noise on each bit

$$P(\underline{y}, \underline{x}) = \prod_{i=1}^N P(y_i | x_i)$$