

## 2 - Coding (compression)

4 października 2010  
13:04

### 2.1 Motivation



Imagine we use a 4-letter alphabet  $\{a, b, c, d\}$  and probabilities for respective symbols in the message are:

	a	b	c	d
$P(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$

How to encode symbols in order to obtain largest  $\frac{m}{n}$

(i) stupid  $a=00$   $b=01$   $c=10$   $d=11$   $n=2 \cdot m$   $\frac{m}{n} = \frac{1}{2}$

(ii) better  $a=0$   $b=10$   $c=110$   $d=111$   
 $n = m \left( \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 \right) = m \cdot 1,75$   $\frac{m}{n} = \frac{1}{1,75}$

### Intuitions

- for long sequences of symbols we will expect roughly  $p(x) \cdot m$  symbols of  $x$  (typical sequences) other sequences will almost never happen so we only need to encode typical sequences

- my kodowanie (ii) ma być wydłużenie a i 1  
 my kodowanie (i) jest 0 i 1 - nie optymalne  
 wyhaszanie kantu

### 2.2 Law of large numbers (empirical mean approach)

(The most important law of statistics) the mean for large number of repeats

$X$  - random variable  $x \in X$

$X_1, \dots, X_N$  -  $N$  independent realisation of  $X$

$\bar{X}_N = \frac{\sum X_i}{N}$   
 Estimator of the mean

we want to prove that  
 $\bar{X}_N \xrightarrow{N \rightarrow \infty} E(X) = \sum p_i x_i$

#### Theorem:

Assume  $X_i$  indep. random variables with  $E(X) < \infty$ ,  $E(X^2) < \infty$

Then for any  $\epsilon$   $p(|\bar{X}_N - E(X)| > \epsilon) \xrightarrow{N \rightarrow \infty} 0$

#### Proof:

$$P(x_1, \dots, x_N) = \prod_i p(x_i)$$

$$E(\bar{X}_N) = E(X) \quad E(|\bar{X}_N - E(X)|^2) = E(\bar{X}_N^2) - E(X)^2$$

$$E(X_N^2) = \frac{N E(X^2)}{N^2} + \frac{1}{N^2} \sum_{i \neq j} E(X)^2 = \frac{E(X^2)}{N} + \frac{N-1}{N} E(X)^2$$

$$E(|\bar{X}_N - E(X)|^2) = \frac{E(X^2)}{N} - E(X)^2 = \frac{E((X - E(X))^2)}{N} \xrightarrow{N \rightarrow \infty} 0$$

Ważne uwagi: Główny  $\int_{\epsilon > 0} p(|\bar{X}_N - E(X)| > \epsilon) \xrightarrow{N \rightarrow \infty} a$   $a > 0$

to oznacza, że  $E(|\bar{X}_N - E(X)|^2) \geq a \cdot \epsilon^2$  sprzeczność, ☹

jeśli  $a \epsilon^2$  maleje szybciej niż  $\frac{1}{N}$  to ok. czyli

$$\epsilon^2 \leq \frac{1}{N} \quad \epsilon \sqrt{N} \xrightarrow{N \rightarrow \infty} 0$$

Uwaga: Dla jakiegokolwiek  $N$  prawdopodobieństwa  $\epsilon$  istnieje  $N$  takie, że  $|x_N - E(X)| > \epsilon$  będzie więcej niż  $a$

$$\frac{E(|x - E(x)|^2)}{N} > a \epsilon^2 \quad N < \frac{\text{Var } X}{a \epsilon^2}$$

prawdopodobieństwa  $|x_m - E(X)| > \epsilon$  i.s.  $a < \frac{\text{Var } X}{N \epsilon^2}$

Fact With probability  $> 1 - \frac{\text{Var } X}{N \epsilon^2}$  a sequence is  $|x_m - E(X)| < \epsilon$

### 2.3 Typical sequences

$x \in X$   $\underbrace{x_1, \dots, x_m}_{x^m}$   $m$  independent realizations

$$P(x^m) = p(x_1) \cdot \dots \cdot p(x_m)$$

$$\log p(x^m) = \sum_{x \in X} n_x \log p(x)$$

↑ number of occurrences of  $x$

By law of large numbers:

$$\frac{1}{m} \log p(x^m) = \sum_{i=1}^m \frac{1}{m} \log p(x_i) \xrightarrow{m \rightarrow \infty} E(\log p(x)) = \sum_x p(x) \log p(x)$$

Definition:

Shannon entropy  $H(X) = - \sum_x p(x) \log p(x)$

$$\left( \begin{array}{l} x=0 \quad p \\ x=1 \quad 1-p \end{array} \quad p=\frac{1}{2} \quad H(x)=1, \quad \begin{array}{l} p=1 \\ p=0 \end{array} \quad H(x)=0 \right)$$

Definition (the most important concept of information theory)

A sequence  $x^m = \{x_1, \dots, x_m\}$  is called  $\epsilon$ -typical iff

$$\left| -\frac{1}{m} \log p(x^m) - H(x) \right| < \epsilon$$

$$T_\epsilon^m = \{x^m : x^m \text{ is } \epsilon\text{-typical}\}$$

↳ set of  $\epsilon$ -typical sequences  $x^m$

Facts  $\forall \delta > 0$ , for sufficiently large  $m$ :

(i)  $\text{pr}[x^m \in T_\epsilon^m] > 1 - \delta$  asymptotically all sequences are typical

(ii)  $(1 - \delta) 2^{m(H(x) - \epsilon)} \leq |T_\epsilon^m| \leq 2^{m(H(x) + \epsilon)}$

Proof:

(i) by law of large numbers

$$\text{pr}[x^m \in T_\epsilon^m] > 1 - \frac{\text{Var} \log p(x)}{m \epsilon^2}$$

(ii)  $x^m \in T_\epsilon^m \quad 2^{-m(H(x) + \epsilon)} \leq p(x^m) \leq 2^{-m(H(x) - \epsilon)}$

$$\sum_{x^m \in T_\epsilon^m} p(x^m) \leq 1$$

$$2^{-n(H(x) + \epsilon)} \sum_{x^m \in T_\epsilon^m} \leq 1 \Rightarrow |T_\epsilon^m| \leq 2^{n(H(x) + \epsilon)}$$

$$p(x^m) \leq 2^{-n(H(x) - \epsilon)}$$

$$|T_\epsilon^m| > (1 - \delta) 2^{n(H(x) - \epsilon)}$$

$$\sum_{x^m \in T_\epsilon^m} p(x^m) > 1 - \delta$$

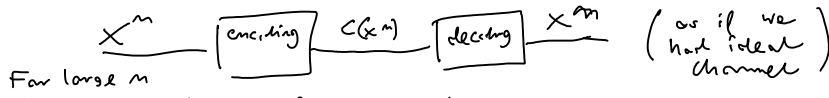


All typical sequences are almost equally probable (Asymptotic equipartition property),  $\approx 2^{-n(H(x) \pm \epsilon)}$

All proofs may assume that we have only typical sequences since atypical sequences are negligible for large  $n$

## 2.4 Shannon coding theorem (data compression)

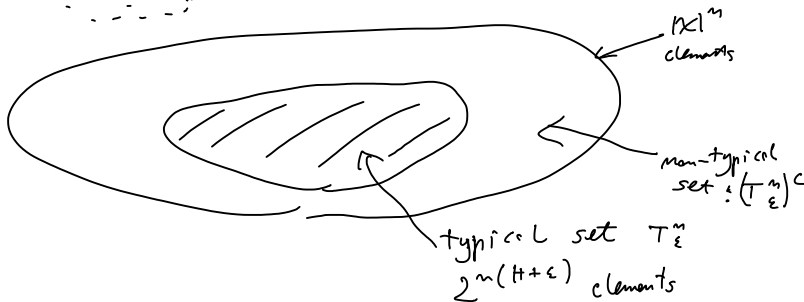
$X^m$  i.i.d distributed  $p(x)$



For large  $n$   
The average length of code words per symbol

$$H(x) - \epsilon \stackrel{(ii)}{\leq} \mathbb{E}[|C(X^m)|] \stackrel{(i)}{\leq} H(x) + \epsilon$$

Proof:  
(i)



We label typical sequences using  $n(H + \epsilon) + 1$  bit sequences + one bit 0 for a prefix

The atypical sequences can be encoded using  $\leq n \log |\mathcal{X}| + 1$  bits (actually less but it does not matter...)

+ one bit 1 for a prefix

$|C(X^m)|$  length of a codeword encoding  $x^m$

On average expected length of codeword will be:

$$\mathbb{E}[|C(X^m)|] = \sum_{x^m \in T_\epsilon^m} p(x^m) |C(x^m)| + \sum_{x^m \in (T_\epsilon^m)^c} p(x^m) |C(x^m)|$$

$$\leq \sum_{x^m \in T_\epsilon^m} p(x^m) [n(H + \epsilon) + 2] + \sum_{x^m \in (T_\epsilon^m)^c} p(x^m) \cdot (n \log |\mathcal{X}| + 2)$$

$$\leq n(H + \epsilon) + 2 + \delta (n \log |\mathcal{X}| + 2) = n(H + \epsilon')$$

$$\epsilon' = \epsilon + \frac{2}{n} + \delta (\log |\mathcal{X}| + \frac{2}{n}) \rightarrow \text{arbitrarily small}$$

— 1 k r . . . . .

$$\text{So } E\left(\frac{1}{n} \log(L(x^n))\right) \leq H(x) + \epsilon$$

(ii) (to prove strictly we need Kraft inequality...)

But intuitively it is clear we need at least  $2^{n(H(x)-\epsilon)}$  codewords or  $n(H(x)-\epsilon)$  bits since messages are uniformly distributed.

More precisely if someone took a smaller set  $|S| \leq 2^{n(H(x)-\epsilon-\epsilon')}$

$$\Pr(x^n \in S) \leq \Pr(x^n \in S \cap T_\epsilon^n) + \Pr(x^n \in S \cap T_\epsilon^c) \leq 2^{n(H(x)-\epsilon-\epsilon')} \cdot 2^{-n(H(x)-\epsilon)} + \delta = 2^{-n\epsilon'} + \delta$$

so it will go to 0 for  $n \rightarrow \infty$ .

## 2.5 Properties of Entropies

Shannon entropy  $H(x) = -\sum_x p(x) \log p(x)$  (measure of uncertainty)

Joint entropy:

$$H(X, Y) = -\sum_{x, y} p(x, y) \log p(x, y)$$

Conditional entropy

$$H(X|Y=y) = -\sum_x p(x|y) \log p(x|y)$$

$$H(X|Y) = -\sum_x p(y) p(x|y) \log p(x|y)$$

Lemma:  $p(x), q(x)$  - prob. distribution

$$D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \geq 0 \quad (\text{relative entropy})$$

$$\left\{ \begin{array}{l} D(p||q) = 0 \text{ iff } p=q \\ D(p||q) \neq D(q||p) \end{array} \right.$$

Proof:

$\log x$  - concave function  $\sum_i p_i \log x_i \leq \log \sum_i p_i x_i$

$$= \sum_x p(x) \log \frac{p(x)}{q(x)} = \sum_x p(x) \log \frac{p(x)}{p(x)} \leq 0 \quad \square$$

Properties:

(i)  $0 \leq H(x) \leq \log |X|$

(ii)  $H(X, Y) \leq H(X) + H(Y)$  { equality for independent variables

(iii)  $H(X, Y) = H(X|Y) + H(Y)$  { chain rule

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

(v)  $H(X, Y|Z) = H(X|Y, Z) + H(Y|Z)$  { chain rule for cond probs

Proofs:

(i) take  $q(x) = \frac{1}{|X|}$ ,  $p(x)$

$$D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = \sum_x p(x) \log p(x) + \sum_x p(x) \log |X| \geq 0$$

$$H(x) \leq \log |\mathcal{X}|$$

□

(ii)  $D(p(x,y) || p(x) \cdot p(y)) =$   $p(x), p(y)$  - niezależne zmienne losowe

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = -H(x,y) + H(x) + H(y) \geq 0$$

□

(iii)  $H(x,y) = -\sum_{x,y} p(x,y) \log p(x,y) = -\sum_{x,y} p(x,y) \log p(x)p(y) =$   
 $= H(x,y) + H(y)$

(iv) from (ii) and (iii) □

## 2.6 Huffman coding

$x, p(x)$

→ we take two  $x$  with smallest probabilities assign 0 and 1 and group them into one new symbol

Example:

codeword	$x$	$p(x)$
10	1	0.25
01	2	0.25
00	3	0.2
110	4	0.1
1111	5	0.1
1110	6	0.1

*(Note: The original image shows a tree diagram with probabilities being updated and crossed out as nodes are merged. For example, 0.25 and 0.25 merge to 0.5, then 0.5 and 0.3 merge to 0.8, etc.)*

on average:  $\frac{1}{2} \cdot 2 + \frac{1}{5} \cdot 2 + 0.1 \cdot 3 + 0.2 \cdot 4 = 1.4 + 0.3 + 0.8 = 2.5$  bits

$H(x) = 2.461$