

Boris Kudryashov

ITMO University

December 18, 2016

- ① Universal coding task
- ② Useful combinatorial formulas
- ③ Two pass encoding
- ④ Enumerative coding
- ⑤ Asymptotic bounds of redundancy
- ⑥ Adaptive coding
- ⑦ Algorithm comparison

- *Encoding redundancy* for a model class Ω is

$$r_n(\Omega) = \sup_{\omega \in \Omega} [\bar{R}_n(\omega) - H_\omega] . \quad (1)$$

- Coding is called *Universal* if for algorithm holds

$$\lim_{n \rightarrow \infty} r_n(\Omega) = 0,$$

Useful combinatorial formulas

- Consider sequences $\mathbf{x} = (x_1, \dots, x_n)$, where x_i has one of M_i values, $i = 1, \dots, n$. Number of different \mathbf{x} is

$$|\{\mathbf{x} = (x_1, \dots, x_n) : x_i \in \{0, \dots, M_i - 1\}, i = 1, \dots, n\}| =$$
$$(2)$$



$$A_M^n = M(M-1) \times \dots \times (M-n+1) = \frac{M!}{(M-n)!}.$$
$$(3)$$

Useful combinatorial formulas

- Number of combinations

$$\begin{aligned}C_M^n &= \binom{M}{n} = \frac{A_M^n}{P_n} = \\&= \frac{M(M-1) \times \dots \times (M-n+1)}{n!} = \\&= \frac{M!}{n!(M-n)!}.\end{aligned}\tag{4}$$

- Number of combinations

$$\binom{n}{k} = \begin{cases} \frac{n!}{k!(n-k)!}, & \text{если } n \geq k \geq 0 \\ 1, & \text{если } n \geq 0 \text{ и } k = 0 \text{ или } k = n \\ 0, & \text{если } k < 0 \text{ или } k > n \end{cases}\tag{5}$$

Useful combinatorial formulas

- binomial coefficient

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}.$$

- Number of binary sequences of length n , which contain τ_1 ones and $\tau_0 = n - \tau_1$ zeros.

$$N(\tau_0, \tau_1) = \binom{n}{\tau_0} = \frac{n}{\tau_0! \tau_1!}. \quad (6)$$

- Composition of sequence \mathbf{x} is vector $\boldsymbol{\tau}(\mathbf{x}) = (\tau_0(\mathbf{x}), \dots, \tau_{M-1}(\mathbf{x}))$, where $\tau_i(\mathbf{x})$ denotes number of elements $x_t = i$ in sequence $\mathbf{x} = (x_1, \dots, x_n)$.

Useful combinatorial formulas

- For $M = 3$

$$N(\tau) = \binom{n}{\tau_0} \binom{n - \tau_0}{\tau_1} = \frac{n!}{\tau_0!(n - \tau_0)!} \frac{(n - \tau_0)!}{\tau_1!(n - \tau_0 - \tau_1)!} = \frac{n!}{\tau_0! \tau_1! \tau_2!}$$

- For arbitrary M

$$N(\tau) = \frac{n!}{\tau_0! \dots \tau_{M-1}!}. \quad (7)$$

- Newton formula generalization

$$(a_0 + \dots + a_{M-1})^n = \sum_{\tau: \tau_0 + \dots + \tau_{M-1} = n} N(\tau) \prod_{i=0}^{M-1} a_i^{\tau_i}.$$

Useful combinatorial formulas

- Consider the following lemma:

Lemma

$n \in \mathbb{N}$ can be written as sum of M non-negative integer terms in $\binom{n+M-1}{M-1}$ ways.

- Number of different compositions of sequence of length n over M -size alphabet is

$$N_{\tau}(n, M) = \binom{n+M-1}{M-1} \quad (8)$$

- Stirling formula

$$\sqrt{2\pi n} n^n e^{-n} \exp \left\{ \frac{1}{12n+1} \right\} < n! < \sqrt{2\pi n} n^n e^{-n} \exp \left\{ \frac{1}{12n} \right\}. \quad (9)$$

Useful combinatorial formulas

- Consider

$$N(\boldsymbol{\tau}) < (2\pi n)^{-\frac{M-1}{2}} 2^{n \log n - \sum_i \tau_i \log \tau_i} \left(\prod_i \frac{n}{\tau_i} \right)^{1/2} \times \\ \times \exp \left\{ \frac{1}{12n} - \sum_i \frac{1}{12\tau_i + 1} \right\}. \quad (10)$$

- Logarithm of number of sequences with specified composition

$$\log N(\boldsymbol{\tau}) < nH(\hat{\boldsymbol{p}}) - \frac{M-1}{2} \log(2\pi n) - \frac{1}{2} \sum_i \log(\hat{p}_i), \quad (11)$$

Useful combinatorial formulas

- More compact estimation

$$\log N(\boldsymbol{\tau}) < nH(\hat{\boldsymbol{p}}) - \frac{M-1}{2} \log(2\pi n) + \frac{1}{2} \log \frac{n}{n-M+1}. \quad (12)$$

- Recurrent formula holds

$$\binom{n+1}{w} = \binom{n}{w} + \binom{n}{w-1}. \quad (13)$$

-

$$\binom{n+1}{w} = \binom{n}{w} + \binom{n-1}{w-1} + \dots + \binom{n-w+1}{1}. \quad (14)$$

Two pass encoding



IF_WE_CANNOT_DO_AS_WE_WOULD_WE_SHOULD_DO_AS_WE_CAN

(15)



$$l_2 = 6 + 6 + 12 \times 2 + 5 \times 3 + \dots + 6 = 178.$$

- $00010000010100110111101101111.$

Two pass encoding



$$l_1 = 29 + 8 \times 15 = 149 \text{ bit.}$$



$$l = l_1 + l_2 = 149 + 178 = 327 \text{ bit.} \quad (16)$$

Two pass encoding

Table: Huffman code for text (15)

Character	Number of iterations	Codeword length	Codeword
I	1	6	010000
F	1	6	010001
_	12	2	00
W	5	3	100
E	4	4	0101
C	2	5	01001
A	4	4	1010
N	3	4	1011
O	5	3	110
T	1	6	011110
D	4	4	0110
S	3	4	1110
U	2	4	1111
L	2	5	01110
H	1	6	011111

Two pass encoding

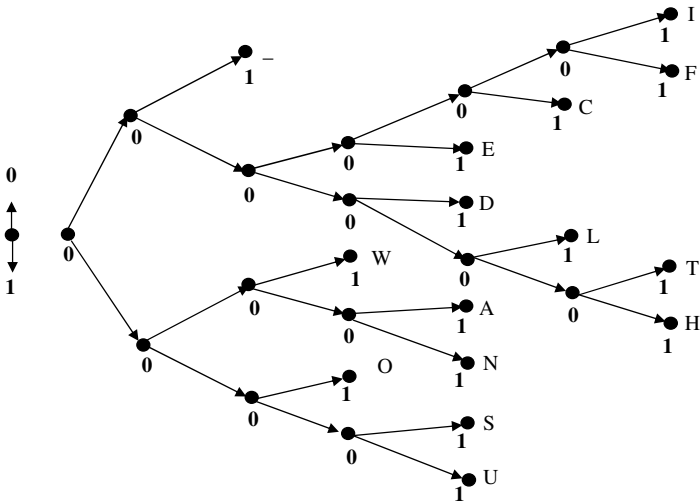


Figure: Huffman codetree for (15)

Two pass encoding

Table: Regular Huffman code

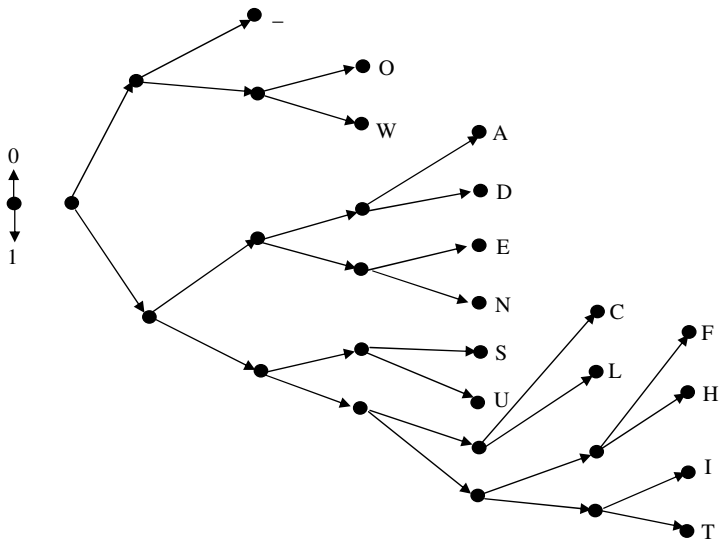
Character	Codeword length	Codeword
—	2	00
O	3	010
W	3	011
A	4	1000
D	4	1001
E	4	1010
N	4	1011
S	4	1100
U	4	1101
C	5	11100
L	5	11101
F	6	111100
H	6	111101
I	6	111110
T	6	111111

Two pass encoding

Table: Number of bits for regular code tree transmitting

Level	Number of nodes	Number of leaves n_i	Range of values n_i	Expenses in bits
0	1	0	0...1	1
1	2	0	0...2	2
2	4	1	0...4	3
3	6	2	0...6	3
4	8	6	0...8	4
5	4	2	0...4	3
6	4	4	0...4	3
Bcero				19

Two pass encoding



Two pass encoding

- Enough for transmitting information about letters that are associated with regular codetree nodes:

$$\left\lceil \log \binom{256}{1} \right\rceil + \left\lceil \log \binom{255}{2} \right\rceil + \left\lceil \log \binom{253}{6} \right\rceil + \\ + \left\lceil \log \binom{247}{2} \right\rceil + \left\lceil \log \binom{245}{4} \right\rceil = 105 \text{ бит}$$

- More precise

$$I = 178 + 19 + 105 = 302 \text{ бит.} \quad (17)$$

Theorem

For two pass coding with Huffman code of Discrete Memoryless Source with alphabet size M and entropy H , average code rate satisfies

$$\bar{R} \leq H + 1 + \frac{1}{n} (M \log M + 3M - 1). \quad (18)$$

Two pass encoding

Proof.

- $l_1(\mathbf{x}) \leq 2M - 1 + M \lceil \log M \rceil \leq M \log M + 3M - 1.$
-

$$\begin{aligned} l_2(\mathbf{x}) &\stackrel{(a)}{=} \sum_{i=1}^n l(x_i) = \\ &\stackrel{(b)}{=} \sum_{x \in X} \tau_n(x) l(x) = \\ &\stackrel{(c)}{=} n \sum_{x \in X} \frac{\tau_n(x)}{n} l(x) = \\ &\stackrel{(d)}{=} n \sum_{x \in X} \hat{p}_n(x) l(x) = \\ &\stackrel{(e)}{=} n \mathbf{M}_{\hat{p}_n} [l(x)] \leq \\ &\stackrel{(f)}{\leq} n(H(\hat{\mathbf{p}}_n) + 1). \end{aligned} \tag{19}$$

Proof.

•

$$\bar{R}(\mathbf{x}) = \frac{l(\mathbf{x})}{n} = \frac{l_1(\mathbf{x}) + l_2(\mathbf{x})}{n} \leq \quad (20)$$

$$\leq H(\hat{\mathbf{p}}_n) + 1 + \frac{1}{n} (M \log M + 3M - 1). \quad (21)$$

•

$$\mathbf{M}[H(\hat{\mathbf{p}}_n)] \stackrel{(a)}{\leq} H(\mathbf{M}[\hat{\mathbf{p}}_n]) \stackrel{(b)}{=} H(\mathbf{p}) = H. \quad (22)$$

•

$$\mathbf{M}[\hat{\mathbf{p}}_n] = \mathbf{p}, \quad (23)$$

•

$$\mathbf{M}\left[\frac{\tau_n(a)}{n}\right] = p(a), \quad a \in X.$$

Proof.

- $$\chi_a(x) = \begin{cases} 1, & \text{при } x = a, \\ 0, & \text{при } x \neq a. \end{cases}$$

- $$\mathbf{M}[\chi_a(x)] = 1 \times p(a) + 0 \times (1 - p(a)) = p(a).$$

- $$\begin{aligned} \mathbf{M}\left[\frac{\tau_n(a)}{n}\right] &= \frac{1}{n} \mathbf{M}\left[\sum_{i=1}^n \chi_a(x_i)\right] = \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{M}[\chi_a(x_i)] = \\ &= p(a), \quad a \in X. \end{aligned}$$

Two pass encoding

- Note, that coding redundancy satisfies

$$r = \bar{R} - H \leq 1 + \frac{K}{n}, \quad (24)$$

- When using arithmetic coding, the redundancy can be achieved:

$$r(n) = \frac{M-1}{n} \log n + \frac{K}{n}, \quad (25)$$

where M alphabet size, K is a constant.

Two pass encoding

Algorithm comparison

Table: Universal coding algorithm comparison

Algorithm	Number of traverses	Asymptotic redundancy	codeword length for text (15)
2-traverse coding, Huffman code	2	$1 + K_1/n$	302
Enumerative coding	2	$\frac{M \log n + K_3}{2n}$	283
Adaptive coding (A)	1	$\frac{M \log n + K_4}{2n}$	291
Adaptive coding (D)	1	$\frac{M \log n + K_5}{2n}$	283