

## Lecture 2: Data compression

Guest lecturer: Grigory Kabatiansky

`g.kabatyansky@skoltech.ru`

Course instructor: Alexey Frolov

`al.frolov@skoltech.ru`

Teaching Assistant: Stanislav Kruglik

`stanislav.kruglik@skolkovotech.ru`

February 2, 2017

- 1 Practical applications of entropy
- 2 Source coding
- 3 Maximal data compression
- 4 Universal coding

# Practical applications of entropy

- ① Book authorship (if two books compressed in the same manner then they have the same author)
- ② Randomness of sequences (if a sequence is “really random” it cannot be compressed)

Let us have a finite alphabet of size  $m$  and a source which emits letter  $i$  from it with probability  $p_i$ . We assume that the next symbol appears independently from the previous one.

We want to represent each letter of initial alphabet as a sequence from another alphabet (codeword) in such a manner that different letters of initial alphabet may have different lengths (variable-length code).

We want to minimize the average sequence length. It's obvious that the most probable letter should have the smallest length and we have to solve minimization task  $\min \sum p_i l_i$ . By  $l_i$  we denote the length of the sequence corresponding to the letter  $i$ .

Code is uniquely decodable if there is only one possible decision for each data sequence.

Code is prefix if no one codeword can be the origin (prefix) of any other codeword.

# Link between uniquely decodable and prefix codes

Prefix code  $\Rightarrow$  uniquely decodable code

Uniquely decodable code  $\nRightarrow$  prefix code

Prefix codes are a subset of uniquely decodable codes.

## Example

Code  $C = \{z_1 = 0, z_2 = 10, z_3 = 11, z_4 = 111\}$  is prefix and uniquely decodable.

Code  $C = \{z_1 = 0, z_2 = 01, z_3 = 011, z_4 = 111\}$  is not prefix, but uniquely decodable

## Theorem

*For any prefix code over the alphabet of size  $D$ , the codeword lengths  $l_1, \dots, l_m$  must satisfy the inequality*

$$\sum_{i=1}^m D^{-l_i} \leq 1.$$

*Conversely, given a set of codeword lengths that satisfy this inequality, there exists a prefix code with these word lengths.*



## Proof.

Consider a  $D$ -ary tree in which each node has  $D$  children. Let the branches of the tree represent the symbols of the codeword. Each codeword is represented by a leaf on the tree. The prefix condition on the codewords implies that no codeword is an ancestor of any other codeword on the tree. Hence, each codeword eliminates its descendants as possible codewords.

Let  $l_{\max} = \max\{l_1 \dots l_m\}$ , then due to the prefix condition we have

$$\sum_{i=1}^m D^{l_{\max} - l_i} \leq D^{l_{\max}}$$



## Theorem

*For any uniquely decodable code over the alphabet of size  $D$ , the codeword lengths  $l_1, \dots, l_m$  must satisfy the Kraft inequality*

$$\sum_{i=1}^m D^{-l_i} \leq 1.$$

*Conversely, given a set of codeword lengths that satisfy this inequality, there exists an uniquely decodable code with these word lengths.*

# Bounds on the optimal codelength

Optimization problem: minimize the average codelength ( $L = \sum p_i l_i$ ) under condition subject to constraints  $l_1, l_2, \dots, l_m$  are integers and  $\sum D^{-l_i} \leq 1$ .

## Theorem

*Let  $l_1^*, l_2^*, \dots, l_m^*$  be the optimal codeword lengths for a source distribution  $P_X$  and a  $D$ -ary alphabet, and let  $L^*$  be the associated expected length of the optimal code ( $L^* = \sum p_i l_i^*$ ). Then*

$$\frac{H(X)}{\log D} \leq L^* < \frac{H(X)}{\log D} + 1$$

# Bounds on the optimal codelength

Proof.

Optimal codeword lengths (method of Lagrange multipliers)

$$l_i = \left\lceil \log_D \frac{1}{p_i} \right\rceil$$



# Huffman code

Optimal code – code that minimizes average length of symbol representation.

Below we illustrate the Huffman algorithm. The main idea is to unite least probable letters in one virtual letter and construct multiple level tree that correspond to the code

Codeword length	Codeword	$X$	Probability			
2	01	1	0.25	0.3	0.45	0.55
2	10	2	0.25	0.25	0.3	0.45
2	11	3	0.2	0.25	0.25	
3	000	4	0.15	0.2		
3	001	5	0.15			

This code has average length 2.3 bits.

## Lemma

*For any distribution, there exists an optimal prefix code (with minimum expected length) that satisfies the following properties:*

- *If  $p_j > p_k$ , then  $l_j \leq l_k$*
- *The two longest codewords have the same length*
- *The two longest codewords differ only in the last bit and correspond to the two least likely symbols*

## Theorem

*Huffman coding is optimal, i.e., if  $C^*$  is the Huffman code and  $C$  is any other code, then  $L(C^*) \leq L(C)$ .*

Despite of optimality, Huffman code is not widely used in practice. The reason is as follows: one need to know the exact values of probabilities  $p_i$  to construct the code. If you slightly change the probability distribution you have to reconstruct corresponding code tree. In the majority of cases optimum does not mean good practical aspects.

As the probabilities of letters may change we introduce classes in which all vectors have the same probability and first encode a number of this class (which is proportional to  $n^m$ ) and then encode vectors of each class by some source code such as Huffman code. Because we additionally code class of vector we lose some useful information. It means, that (assume we encode  $n$  symbols)

$$\frac{\sum_{i=1}^n l_i}{n} \leq H(X) + \frac{\log_2 n}{n}$$

Note, that the length of the resulting sequence is bigger here in comparison to optimal source code ( $H(X) + \frac{\log_2 n}{n}$  versus  $H(X) + \frac{1}{n}$ ). But the difference is negligible in asymptotic regime.



Thank you for your attention!