

Seminar 3

Algorithm 1 (Variable byte code) A number n is encoded in variable byte code in the following procedure:

1. Take a binary representation of n with padding to the length of a multiple of 7.
2. Split into of 7 bit blocks right-to-left.
3. Add 1 to the beginning of the last block and 0 to the beginning of all previous blocks.

Example: $VB(824) = 0000011010111000$

Definition 1 (Unary code) Unary code, also referred to as α code, is a coding type where a number n is represented by a sequence of n 1s (or 0s) and terminated with one 0 (or 1). That is, 6 in unary code is 1111110 (or 0000001). The alternative representation in parentheses is equivalent but for this course we use the default representation.

Definition 2 (γ code) γ code is a coding type, that consists of an offset and its length: $\gamma(n) = \text{length of offset}(n) \text{ in } \alpha, \text{offset}(n)$. Offset is a binary representation of a number n without the highest bit (1). Length of this offset encoded in the unary (α) code. Then the number 60 is encoded in γ as 111110,11100.

Definition 3 (δ code) A number n is encoded in δ code in the following way. First calculate the offset of n and the length of n encode with γ code. Then add the offset of n . The final form is $\delta(n) = \text{length of offset}(n) \text{ in } \gamma, \text{offset}(n)$. Analogously, 600 is encoded in δ as 1110,001,001011000.

Definition 4 (Zipf's law) Zipf's law says that the i -th most frequent term has the frequency $\frac{1}{i}$. In this exercise we use the dependence of the Zipf's law $cf_i \propto \frac{1}{i} = ci^k$ where cf_i is the number of terms t_i in a given collection with $k = -1$.

Definition 5 (Heaps' law) Heaps' law expresses an empiric dependency of collection size (number of all words) T and vocabulary size (number of distinct words) M by $M = kT^b$ where $30 \leq k \leq 100$ and $b \approx \frac{1}{2}$.

Exercise 1

Count variable byte code for the postings list $\langle 777, 17\,743, 294\,068, 31\,251\,336 \rangle$. Bear in mind that the gaps are encoded. Write in 8-bit blocks.

Encode the list of gaps $\langle 777, 16\,966, 276\,325, 30\,957\,268 \rangle$. Variable byte code of the gaps:

- $VB(777) = 0000011010001001$
- $VB(16\,966) = 000000010000010011000110$
- $VB(276\,325) = 000100000110111011100101$
- $VB(30\,957\,268) = 00001110011000010011110111010100$

Result: $VB(\langle 777, 17\,743, 294\,068, 31\,251\,336 \rangle) = 00000110100010010000000100000100110001100001000011011011101110010100001110011011000010011110111010100$

Exercise 2

Count γ and δ codes for the numbers 63 and 1023.

According to the definition 2 it is necessary to count the offsets as binary representations without the highest bit $63_{10} = 11111_2$ and $\text{offset}(63) = 11111$. Offset length is encoded in α as $|11111| = 5 \rightsquigarrow \alpha(5) = 111110$. Finally, $\gamma(63) = 111110, 11111$. Analogically for 1023. $1023_{10} = 111111111_2$, offset is 111111111 , its length is $|111111111| = 9 \rightsquigarrow \alpha(9) = 1111111110$. Then $\gamma(1023) = 1111111110, 111111111$.

δ is a little more complicated. First we count the offset $63 = 11111$ and its length $|11111| = 5$. The value of 5 we encode in γ so $\gamma(5) = 110, 01$. By definition 3 we have $\delta(63) = 110, 01, 11111$. And finally, $\delta(1023) = 1110, 010, 111111111$.

Exercise 3

A collection of documents contains 4 words: *one*, *two*, *three*, *four* of decreasing word frequencies f_1 , f_2 , f_3 and f_4 . The total number of tokens in the collection is 5000. Assume that the Zipf's law holds for this collection perfectly. What are the word frequencies?

We use the Zipf's law in Definition 4. The least frequent term is *four*, then *three*, *two* and the most frequent is *one*. Applying the Zipf's law we get

$$\begin{aligned} cf_1 + cf_2 + cf_3 + cf_4 &= 5000 \\ c \cdot 1^{-1} + c \cdot 2^{-1} + c \cdot 3^{-1} + c \cdot 4^{-1} &= 5000 \\ c + \frac{1}{2}c + \frac{1}{3}c + \frac{1}{4}c &= 5000 \\ 12c + 6c + 4c + 3c &= 60000 \\ 25c &= 60000 \\ c &= 2400 \end{aligned}$$

and, plugging in to the formula $cf_i = ci^{-1}$, we obtain the term frequency values:

$$\begin{aligned} cf_1 &= 2400 \frac{1}{1} = 2400 \\ cf_2 &= 2400 \frac{1}{2} = 1200 \\ cf_3 &= 2400 \frac{1}{3} = 800 \\ cf_4 &= 2400 \frac{1}{4} = 600 \end{aligned}$$

Exercise 4

How many distinct terms are expected in a document of 1,000,000 tokens? Use the Heaps' law with parameters $k = 44$ and $b = 0.5$

By Definition 5,

$$44 \times 1,000,000^{0.5} = 44,000.$$