

1 Wzory

Informacja zdarzenia A :

$$I(A) = -\log_x P(A)$$

Entropia źródła X ze zdarzeniami A_1, \dots, A_n :

$$H(X) = \sum_{i=1}^n P(A_i) I(A_i)$$

Średnia długość kodu C :

$$I(C) = \sum_{i=1}^n P(C_i) \cdot l_i$$

Nierówność Krafta (warunek konieczny jednoznacznej dekodowalności):

$$K(C) = \sum_{i=1}^n 2^{-l_i} \leq 1$$

Współczynnik informacji kodu C :

$$\frac{1}{n} \log |\mathcal{C}|$$

2 Kod Huffmana

Znajdź dwa najrzadziej występujące elementy i połącz je w jeden element o prawdopodobieństwie $p_1 + p_2$. Rozróżnij je 0 lub 1. Powtórz ten krok na liście $n - 1$ długiej aż zostanie jeden element.

Jeśli nie znamy prawdopodobieństw, to możemy drzewo tworzyć dynamicznie, traktując ilość wystąpień jako wagę, które łączymy tworząc poddrzewa.

3 Kod Shannon-Fano

Dla symboli a_1, \dots, a_n o prawdopodobieństwach p_1, \dots, p_n , ustalmy kody długości $l_n = \lceil -\log p_i \rceil$. Następnie zdefiniujmy zmienne pomocnicze w_1, \dots, w_n jako:

$$w_1 = 0, w_j = \sum_{i=1}^{j-1} 2^{l_j - l_i}$$

Jeżeli $\lceil \log w_j \rceil = l_j$ to j -te słowo kodowe jest binarną reprezentacją w_j . Jeżeli $\lceil \log w_j \rceil < l_j$ to reprezentację uzupełniamy zerami z lewej strony.

4 Kod Tunstalla

Chcemy stworzyć kod na n bitach dla a_1, \dots, a_m symboli o prawdopodobieństwach p_1, \dots, p_m . Tworzenie kodu Tunstalla polega na iteracyjnym wyborze ze zbioru symbolu o największym prawdopodobieństwie S i łączenie go z wszystkimi innymi symbolami tworząc symbole Sa_m , nadając im prawdopodobieństwa $P \cdot p_m$. Proces ten powtarzamy aż do uzyskania kodu o długości n .

5 Kodowanie Eliasa

$$n = \lfloor \log_2(x) \rfloor + 1$$

5.1 γ

$$\gamma(x) = 0^{n-1}(x)_2$$

5.2 δ

$$\delta(x) = \gamma(n) + (x)_2$$

5.3 ω

Na koniec umieszczane jest 0, potem kodowana jest liczba $k = x$. Potem ten krok jest powtarzany dla $k = n - 1$ gdzie n to liczba bitów z poprzedniego kroku.

$$\omega(x) = \omega(n - 1) + (x)_2 + 0$$

6 Kodowanie Fibonacciego

$$\begin{aligned} f_0 &= f_1 = 1 \\ f_n &= f_{n-1} + f_{n-2} : n \geq 2 \\ x &= \sum_{i=0} a_i \cdot f_i, a_i \in \{0, 1\} \end{aligned}$$

7 Kodowanie arytmetyczne

- $d = p - l$
- $p = l + d \cdot F(j + 1)$
- $l = l + d \cdot F(j)$
- $d = p - l$
- Wybieramy takie j , że $F(j) \leq \frac{x-l}{d} < F(j+1)$
- $p = l + d \cdot F(j+1)$
- $l = l + d \cdot F(j)$

8 Kodowanie słownikowe

8.1 LZ77

$$(o, l, k) = C_{i-o} \cdots C_{i-o+l} k$$

8.2 LZ78

1. Szukaj w słowniku najdłuższy prefiks aktualnego okna, jeśli nie znajdziesz to użyj ϵ .
2. Dodaj prefiks + znak do słownika.
3. Zakoduj symbol jako (i, k) , gdzie i to numer prefiksu w słowniku, a k to symbol.

$$(i, k) = s(i) + k$$

8.3 LZW

Podobne do LZ78, tylko że zaczynamy ze słownikiem.

$$(i) = s(i)$$

Jeśli napotkasz symbol, którego nie ma w słowniku, ale masz jego prefiks, np.: $s(5) = ab?$, to ? to pierwsza litera $s(5)$.

9 bzip2/BWT

Układamy tabelę z dwoma kolumnami. Pierwsza kolumna to słowo posortowane leksykograficznie. Druga kolumna to poprzedni znak. Na podstawie tej tabeli zapisujemy ostatnią kolumnę, i numer wiersza w którym w pierwszej kolumnie znajduje się początek słowa, a w drugiej kolumnie jego koniec.

e	h	0	1	2	3	4
h	o	e	h	l	l	o
ll	e	2	0	3	4	1
lo	l					
o	l					

10 Move To Front

Jest to transformacja zmniejszająca entropię. Zaczynamy od tabeli liter ze słowa posortowanych alfabetycznie. Następnie dla każdej litery ze słowa kodujemy jej pozycję w tabeli, a następnie przesuwamy ją na początek tabeli. W ten sposób hello to 11203.

11 PPM

Dla każdego symbolu z tekstu, sprawdzamy jego kontekst, zakładając daną maksymalną długość kontekstu. Kolejno sprawdzamy drzewa kontekstowe dla długości $n, n - 1, \dots, 0$ gdzie interesuje nas tylko to ile razy dany symbol występuje w tekście, aż dojdziemy do kontekstu -1 gdzie zaznaczamy tylko czy symbol występuje czy nie. Następnie na podstawie takiej serii drzew, można zbudować solidny kod, chociażby Huffmana.

Symbol	Symbol	Licznik	Kontekst	Symbol	Licznik	Kontekst	Symbol	Licznik
t	ESC	1	t	ESC	1	th	ESC	1
h	t	1	h	h	1		l	1
l	h	1	l	ESC	1	hl	ESC	1
s	l	1	s	l	1	is	ESC	1
-	s	2	s	s	2	-	-	1
-	-	1	s	ESC	1	s-	ESC	1
			-	-	1	-i	ESC	1
				i	1		s	1

12 Kody Hamminga

Dla długości kodu $n = 2^m - 1$, mają liczbę bitów informacji $k = n - m$. Macierz parzystości, powstaje poprzez zapis binarnych reprezentacji liczb $1, \dots, n$ w postaci kolumn.

$$H_H(3) = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

Aby stworzyć macierz generującą, musimy potraktować H_H jako macierz równania, gdzie pierwsze k wierszy to macierz identyczności, a pozostałe to przekształcone równania z macierzy parzystości.

$$G_H(3) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \end{bmatrix}$$