

ZAKRES MATERIAŁU

I. Entropia
II. Kodowania Huffmana a entropia

I. Entropia

[DEF 01] (Entropia)

Niech X – zmienna losowa przyjmująca skończenie wiele wartości zgodnie z rozkładem prawdopodobieństwa $p(X)$. Entropię tego rozkładu definiujemy

$$H(x) = - \sum_{i=1}^n p_i \log_2 p_i \quad (\text{dla } p_i \neq 0)$$

[Przykład 01] Entropia jednego rzutu symetryczną kostką do gry jest równa 1, ponieważ informację, że wypadło 1/2/3/4/5/6 oczek można przekazać za pomocą jednej jednostki informacji.

[Przykład 02] Entropia n niezależnych rzutów symetryczną kostką do gry jest równa n , ponieważ informację o n rzutach kostką można przekazać za pomocą n jednostek informacji.

[Przykład 03] X jest zmienną losową przyjmującą wartości ze zbioru $A = \{a, b, c, d\}$ z prawdopodobieństwami $p(a)=1/4, p(b)=1/2, p(c)=1/8, p(d)=1/8$. Policz entropię związaną z wylosowaniem jednej wartości ze zbioru A .

$$H(x) = - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{8} \log_2 \frac{1}{8} = \frac{1}{4} \cdot 2 + \frac{1}{2} \cdot 1 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = \frac{7}{4}$$

Entropia wylosowania jednej wartości ze zbioru A wynosi $7/4$.

Entropia składników systemu kryptograficznego.

Możemy policzyć

- ♦ entropię $H(K)$, gdzie K jest zmienną losową (reprezentującą klucz) przyjmującą wartości zgodnie z rozkładem prawdopodobieństwa p_K ,
- ♦ entropię $H(P)$, gdzie P jest zmienną losową (reprezentującą przestrzeń tekstu jawnego),
- ♦ entropię $H(C)$, gdzie C jest zmienną losową (reprezentującą przestrzeń tekstu zaszyfowanego)

[Przykład 04] Niech $P = \{a, b, c\}$ będzie przestrzenią tekstów jawnych oraz niech prawdopodobieństwa a priori wystąpienia poszczególnych znaków tego zbioru będą równe $p_P(a) = 1/4$ i $p_P(b) = 1/2$, $p_P(c) = 1/4$. Entropię $H(P)$ liczymy

$$H(P) = - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} = \frac{1}{4} \log_2 4 + \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 4 = \frac{1}{4} \cdot 2 + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 = \frac{3}{2}$$

II. Kodowania Huffmana a entropia

[DEF 02] (kodowanie)

Niech X – zmienna losowa przyjmująca wartości ze skończonego zbioru zgodnie z rozkładem prawdopodobieństwa $p(X)$.

Kodowanie zmiennej losowej X – przekształcenie $f: X \rightarrow \{0,1\}^*$.

$\{0,1\}^*$ - zbiór wszystkich skończonych ciągów 0-1

Niech x_1, \dots, x_n będzie ciągiem zdarzeń.

Kodowaniem f nazywamy (też) funkcję $f(x_1, \dots, x_n) = f(x_1) | \dots | f(x_n)$, gdzie Znak $|$ oznacza konkatencję ciągów.

[Własności] (kodowania f)

- ♦ Funkcje f powinny być **przekształceniem różnowartościowym**, ponieważ tylko w takim przypadku zaszyfrowany (skompresowany / zakodowany) tekst da się odszyfrować (zdekompresować / zdekodować).
- ♦ Funkcja f musi być **funkcją prefiksową** (inaczej mówimy **funkcją wolną od przedrostków**) tzn. nie może zdarzyć się sytuacja, że dla dwóch elementów $a, b \in X$ i ciągu $x \in \{0,1\}^*$ zachodzi $f(a) = f(b) | x$ (takie sytuacje powodowałyby niejasności przy dekodowaniu).

[Przykład 05] Niech eksperyment polega na n -krotnym rzucie monetą

Niech $X = \{O, R\}$ (O – orzeł, R - reszka)

- ♦ $f(O) = 01$, $f(R) = 0110$ - f nie jest funkcją prefiksową
- ♦ W takiej sytuacji (np. podczas rozkodowania ciągu 01101...) nie ma jasności jaki jest pierwszy znak (O czy R).
- ♦ $f(O) = 01$, $g(R) = 11$ - g jest funkcją prefiksową
- ♦ W takiej sytuacji zawsze mamy jednoznaczną odpowiedź, jaki jest kolejny znak (O czy R).

[DEF 03] (średnia długość kodowania elementu X)

Miara skuteczności kodowania f definiowana następująco

$$\ell(f) = \sum_{x \in X} p(x) |f(x)|$$

$|y|$ -ozn. długość ciągu y .

- ♦ Dąży się do wyszukania takiego różnowartościowego kodowania f , które minimalizuje wartość $\ell(f)$.

[Algorytm] (kompresji i dekompresji Huffmana)

- ♦ Zachłanny algorytm Huffmana generuje kodowanie f , które
 - ♦ jest różnowartościowe
 - ♦ prefiksowe
 - ♦ spełnia nierówności $H(X) \leq \ell(f) < H(X) + 1$ (entropia przybliża średnią długość kodowania Huffmana)

[Idea] (algorytmu kompresji Huffmana)

KOMPRESJA

- (1) Wyznaczamy tablicę częstości wystąpień poszczególnych znaków w tekście i na jej podstawie określamy rozkład prawdopodobieństwa na zbiorze X możliwych znaków.
- (2) Generujemy uporządkowaną niemalejąco (ze względu na wartość prawdopodobieństwa) listę (np. ArrayList) jednowęzłowych drzew, których węzły zawierają następujące pola (klucz (prawdopodobieństwo wystąpienia), znak (lub ‘_’ w przypadku węzłów wewnętrznych), referencje do lewego i prawego poddrzewa)
- (3) Na podstawie listy tworzymy drzewo Huffmana. W tym celu, w każdej iteracji bierzemy z listy dwa drzewa $L1$ i $L2$ o najmniejszych prawdopodobieństwach w korzeniach i łączymy je w jedno drzewo. Kluczem korzenia powstałego drzewa jest suma prawdopodobieństw korzeni łączonych drzew. Jego lewym poddrzewem jest $L1$, a prawym – $L2$. Nowe drzewo umieszczamy na liście w taki sposób, by nadal była uporządkowana niemalejąco. Proces kończymy w momencie uzyskania jednego drzewa, które ma tę własność, że węzły wewnętrzne zawierają jedynie prawdopodobieństwa, a liście dodatkowo znaki występujące w kompresowanym tekście.
- (4) Generujemy kody Huffmana. Każdą lewą krawędź drzewa oznaczamy znakiem ‘0’, a prawą – ‘1’. Czytamy każdą ze ścieżek drzewa zaczynając od korzenia i kończąc na liściu i zapisujemy występujące na ścieżce zera i

- jedynki. Powstałe ciągi 0-1-kowe stanowią kody znaków występujących w liściach. Znaki i ich kody zapisujemy w strukturze HashMap (znak w roli klucza i kod jako obiekt).
- Kodujemy tekst korzystając z wygenerowanych kodów, a następnie każdy 0-1-kowy 8-elementowy podciąg zamieniamy na znak (zgodnie z tablicą kodów ASCII). Jeśli ostatni ciąg składa się z mniejszej niż osiem liczb znaków, dopełniamy go zerami.
 - Do pliku wyjściowego wklejamy tablicę częstości wystąpień poszczególnych znaków plus tekst w postaci zakodowanej.

DEKOMPRESJA

- Na podstawie informacji zawartej w skompresowanym pliku generujemy tablicę częstości wystąpień znaków.
- Tworzymy listę drzew jednowzłowych i na jej podstawie generujemy drzewo Huffmana.
- Zamieniamy zakodowany tekst na ciąg 0-1-kowy (każdy znak tekstu to wartość z zakresu 0-255, którą zapisujemy w postaci binarnej)
- Na podstawie tablicy częstości określamy ile znaków należy rozkodować. Czytamy kolejne zera i jedynki ciągu i na tej podstawie poruszamy się odpowiednio w lewo lub w prawo na drzewie Huffmana. Po rozkodowaniu fragmentu ciągu 0-1 (tj. po przypisaniu mu odpowiedniego znaku znajdującego się w liściu drzewa) rozpoczynamy dekodowanie kolejnego (startując na nowo od korzenia drzewa).

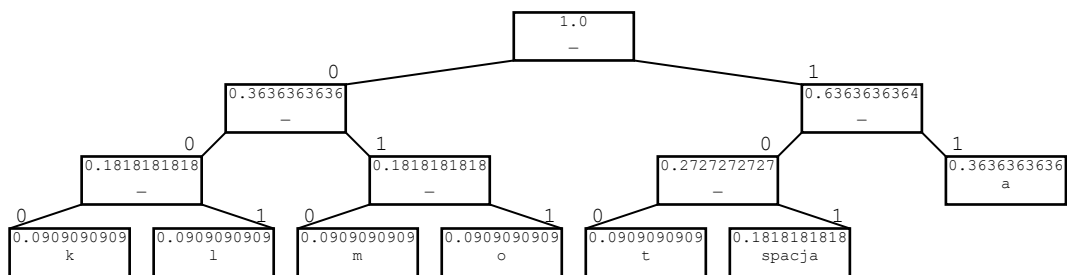
[Zadanie 01] Zastosuj algorytm kompresji i dekompresji Huffmana dla poniższego tekstu.
 ala ma kota

KOMPRESJA

- tablica częstości wystąpień poszczególnych znaków w tekście

znak	a	l	spacja	m	k	o	t
częstość	4	1	2	1	1	1	1
- uporządkowana niemalejąco lista drzew jednowzłowych

k	l	m	o	t	spacja	a
0.0909090909	0.0909090909	0.0909090909	0.0909090909	0.0909090909	0.1818181818	0.3636363636
- drzewo kodów Huffmana



(4) kody Huffmana

{spacja=101, o=011, k=000, m=010, a=11, t=100, l=001}

(5) zakodowany tekst z podziałem na podciągi 8-znakowe

a l a m a k o t a
 11001111 01010111 01000011 10011000 czerwone zera dopełniają ostatnią ósemkę

(6) zawartość skompresowanego pliku (tablica częstości i zakodowany tekst)

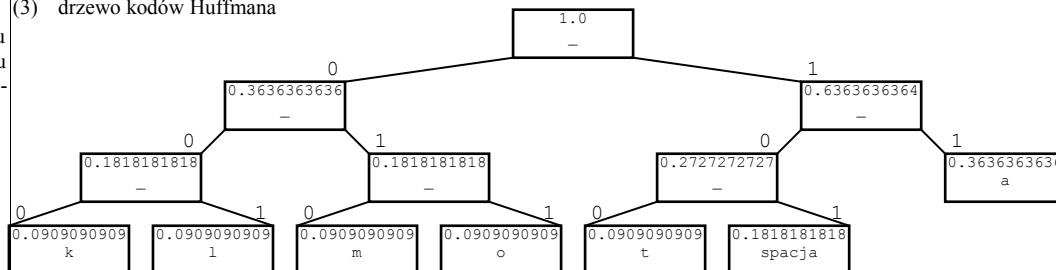
a-4,l-1,_-2,m-1,k-1,o-1,t-1
 gUh9@

DEKOMPRESJA

- tablica częstości wystąpień poszczególnych znaków w tekście

znak	a	l	spacja	m	k	o	t
częstość	4	1	2	1	1	1	1
- uporządkowana niemalejąco lista drzew jednowzłowych

k	l	m	o	t	spacja	a
0.0909090909	0.0909090909	0.0909090909	0.0909090909	0.0909090909	0.1818181818	0.3636363636
- drzewo kodów Huffmana



(4) zakodowany tekst z przypisanymi znakami

11|001|11|101|010|11|101|000|011|100|11|000 czerwone zera dopełniają ostatnią ósemkę
 a l a _ m a _ k o t a

[Zadanie 02] Zaimplementuj

- algorytm kompresji i dekompresji Huffmana. Dla ułatwienia, koduj 7-elementowe ciągi 0-1.
- algorytm wyznaczający entropię H(P), gdzie P jest zmienną losową (reprezentującą przestrzeń tekstu jawnego). Prawdopodobieństwa wystąpień poszczególnych znaków w tekstach polskich znajdują się w materiałach do jednych z wcześniejszych zajęć.
- algorytm wyznaczający średnią długość kodowania znaku dla konkretnego tekstu. Porównaj wartości z podpunktu (b) i (c) dla tekstu „ala ma kota”

[Przykład 06] Wyznacz średnią długość kodowania znaku algorytmem Huffmana dla tekstu „ala ma kota”

znak	a	l	spacja	m	k	o	t
częstość	4	1	2	1	1	1	1
p-stwo	0.3636363636	0.0909090909	0.1818181818	0.0909090909	0.0909090909	0.0909090909	0.0909090909
dł. kodu	2	3	3	3	3	3	3

$$\ell(f) = \sum_{x \in X} p(x) |f(x)|$$

$$\ell(f) =$$

$$= p('a') |f('a')| + p('l') |f('l')| + p('_') |f('_')| + p('m') |f('m')| + p('k') |f('k')| + p('o') |f('o')| + p('t') |f('t')| =$$

$$= 0.3636363636 * 2 + 5 * (0.0909090909 * 3) + 0.1818181818 * 3 =$$

$$= 0.7272727273 + 5 * 0.2727272727 + 0.5454545454 = 0.7272727273 + 1.3636363636 + 0.5454545454 = 2.636363632$$

