

Sprawozdanie z laboratorium Teorii informacji i kompresji danych

Lab nr 4 Kompresja bezstratna - Wstęp

Lab nr 5 Kompresja bezstratna – Kodowanie Huffmana

Wojciech Szczepaniak

Index: 136808

Data zajęć: 21.01.2020 (czwartek 13:30)

1. Kodowanie binarne

W celu realizacji zadania należało przygotować następujące metody:

a. Create

```
from bitarray import bitarray
def create(counted_signs):
    sign_len = len(counted_signs)
    code_len = len(bin(sign_len)[2:])
    code_tracker = {}
    counter = 0
    for sign_tuple in counted_signs:
        tmp = bitarray(bin(counter)[2:])
        for j in range(code_len-1):
            tmp.insert(0,False)
        code_tracker[sign_tuple[0]] = tmp[-code_len:]
        counter+=1
    return(code_tracker)
```

b. Encode

```
def encode(code,text):
    encoded_text = bitarray(0)
    for i in text:
        encoded_text= encoded_text + code[i]
    return(encoded_text)
```

c. Decode

```
def decode(code,coded_text):
    decoded_text = ''
    tmp = bytearray(0)
    for i in coded_text:
        tmp = tmp + bytearray(str(int(i)))
        try:
            for key,value in code.items():
                if tmp == value:
                    decoded_text = decoded_text + str(key)
                    tmp = bytearray(0)
        except:
            pass
    return(decoded_text)
```

d. Save

```
def save(encoded_text,filename,filename_code,code):
    f = open(filename, "w")
    f.write(str(len(code.values()))+"\n")
    for key, value in code.items():
        f.write(key+" "+str(value.to01())+"\n")
    f.close()
    with open(filename_code, 'wb') as fh:
        encoded_text.tofile(fh)
    return(1)
```

e. Load

```
def load(filename_text,filename_code):
    code = {}
    f = open(filename_code, "r")
    f_list = list(f)
    num_keys = f_list[0]
    del f_list[0]
    for i in range(int(num_keys)):
        code[f_list[i].rsplit(" ",1)[0]] = bytearray(f_list[i].rsplit(" ",1)[1].replace("\n",""))
    with open(filename_text,'rb') as file:
        text = bytearray(0)
        text.fromfile(file)
    return(code,text)
```

f. Zadanie 1

```
import string
def zadanie1(filename):
    filename = 'norm_wiki_sample_cut.txt'
    filename_save = 'norm_wiki_sample_save.txt'
    filename_code = 'norm_wiki_sample_code.txt'
    znaki = list(string.ascii_lowercase)
    liczby = list(string.digits)
    all_signs = znaki + liczby + list(" ")

    with open(filename, "r") as f1:
        data = f1.read()

    sign_tracker = {}
    for i in range(len(all_signs)):
        sign_tracker[all_signs[i]] = 0
    for i in data:
        try:
            sign_tracker[i] += 1
        except:
            #other signs
            pass

    sign_tracker = dict(sorted(sign_tracker.items(), key=lambda item: item[1], reverse = True))
    sign_list = [(k,v) for k,v in sign_tracker.items()]

    code = create(sign_list)
    encoded_text = encode(code,data[0:250])
    save(encoded_text,filename_save,filename_code,code)
    loaded_code,loaded_text = load(filename_code,filename_save)
    decoded_text = decode(loaded_code,loaded_text)
    if data[0:250] == decoded_text:
        msg = "poprawnie"
    else:
        msg = "niepoprawnie"
    return(code,msg,decoded_text,sign_list)

code,msg,decoded_text,alfabet = zadanie1('norm_wiki_sample.txt')
```

Wynik zadania pierwszego:

Wynikiem zadania pierwszego jest wygenerowany kod 6 bitowy, string "poprawnie" bądź "niepoprawnie" w zależności czy odkodowany tekst jest identyczny z tekstem początkowym oraz odkodowany tekst w tym wypadku identyczny tekst jak w pierwszych 250 znakach pliku 'norm_wiki_sample.txt'.

```
{' ': bytearray('000000'), 'e': bytearray('000001'), 'a': bytearray('000010'), 't': bytearray('000011'), 'i': bytearray('000100'), 'n': bytearray('000101'), 'o': bytearray('000110'), 'r': bytearray('000111'), 's': bytearray('001000'), 'h': bytearray('001001'), 'l': bytearray('001010'), 'd': bytearray('001011'), 'c': bytearray('001100'), 'm': bytearray('001101'), 'u': bytearray('001110'), 'f': bytearray('001111'), 'p': bytearray('010000'), 'g': bytearray('010001'), 'b': bytearray('010010'), 'w': bytearray('010011'), 'y': bytearray('010100'), 'v': bytearray('010101'), 'k': bytearray('010110'), '1': bytearray('010111'), '0': bytearray('011000'), '9': bytearray('011001'), '2': bytearray('011010'), 'j': bytearray('011011'), '8': bytearray('011100'), '3': bytearray('011101'), '5': bytearray('011110'), 'x': bytearray('011111'), '4': bytearray('100000'), '7': bytearray('100001'), '6': bytearray('100010'), 'z': bytearray('100011'), 'q': bytearray('100100')}
```

poprawnie

albert of prussia 17 may 1490 20 march 1568 was the last grand master of the teutonic knights who after converting to lutheranism became the first monarch of the duchy of prussia the secularized state that emerged from the former monastic state of t

Najkrótsza możliwa długość kodu dla korpusu to 6 bitów, należy zakodować 37 znaków $2^5 < 37 < 2^6$.

Stopień kompresji wyniesie 1.(3). 8 bitów / 6 bitów

Rozmiar pliku przed kompresją wynosi 2416 kb. Po kompresji pliki mają kod 1 kb a zakodowany tekst 1812 kb. Co daje współczynnik kompresji równy 75.04%.

2. Efektywność kodowania

```
import math
def zadanie2(alfabet):
    entropia = 0
    counter = 0
    for i in alfabet:
        counter+=i[1]
    for i in alfabet:
        entropia += (i[1]/counter * math.log(1/(i[1]/counter),2))
    srednia_dlugosc = 6
    efektywnosc = entropia/srednia_dlugosc
    return entropia,efektywnosc
```

```
zadanie2(alfabet)
```

```
(4.2803962467015655, 0.7133993744502609)
```

Dla kodowania z zadania nr 1 entropia wyniosła 4.28, średnia długość słów kodowych była stała i wynosiła 6 a efektywność wyniosła 71.34%.

3. Kodowanie Huffmana

Kodowanie huffmana zostało wykonane dzięki wykorzystaniu funkcji z zadania nr 1. Dodatkowo należało napisać funkcję tworzenia kodowania huffmana oraz funkcje odpowiedzialne za wykonanie kolejnych kroków.

a. Funkcja tworząca kodowanie Huffmana

```
def huffman_decode(node, array = bytearray(0)):
    if type(node) is str:
        return {node: array}
    l = node[0]
    r = node[1]
    d = dict()
    d.update(huffman_decode(l, array + '0'))
    d.update(huffman_decode(r, array + '1'))
    return d

def huffman_code(sign_list):
    nodes = []
    for i in range(len(sign_list)):
        g = sign_list[i]
        nodes.append(g)

    while len(nodes) > 1:
        (key1, value1) = nodes[-1]
        (key2, value2) = nodes[-2]
        nodes = nodes[:-2]
        node = (key1, key2)
        nodes.append((node, value1 + value2))
        nodes = sorted(nodes, key=lambda x: x[1], reverse=True)
    huffmanCode = huffman_decode(nodes[0][0])
    return huffmanCode
```

b. Funkcja wyliczająca średnią długość słowa w kodzie

```
def mean_huffman(code,p_list):
    mean_sum = 0
    for i in range(len(code)):
        mean_sum += len(code[i][1])*float(p_list[i][1])
    return(mean_sum)
```

c. Funkcja odpowiedzialna za wykonanie zadania nr 3

```
def zadanie3(code):
    filename = 'norm_wiki_sample.txt'
    filename_save = 'norm_wiki_sample_save.txt'
    filename_code = 'norm_wiki_sample_code.txt'
    znaki = list(string.ascii_lowercase)
    liczby = list(string.digits)
    all_signs = znaki + liczby + list(" ")

    with open(filename,"r") as f1:
        data = f1.read()

    sign_tracker = {}
    for i in range(len(all_signs)):
        sign_tracker[all_signs[i]] = 0
    counter = 0
    for i in data:
        try:
            sign_tracker[i] += 1
            counter+=1
        except:
            #other signs
            pass

    sign_tracker = dict(sorted(sign_tracker.items(), key=lambda item: item[1], reverse = True))
    sign_list = [(k,v/counter) for k,v in sign_tracker.items()]

    code = huffman_code(sign_list)
    code_list = [(k,v) for k,v in code.items()]
    encoded_text = encode(code,data[0:250])
    save(encoded_text,filename_save,filename_code,code)
    loaded_code,loaded_text = load(filename_code,filename_save)
    decoded_text = decode(loaded_code,loaded_text)
    if data[0:250] == decoded_text:
        msg = "poprawnie"
    else:
        msg = "niepoprawnie"

    p = mean_huffman(code_list,sign_list)
    entropia = zadanie2(alfabet,srednia_dlugosc = p)
    return(code,msg,decoded_text,entropia)

code,msg,decoded_text,entropia = zadanie3('norm_wiki_sample.txt')
print(code)
print(msg)
print(decoded_text)
print(entropia)
```

Wynik zadania trzeciego:

Wynikiem zadania pierwszego jest wygenerowany kod Huffmana, string “poprawnie” bądź “niepoprawnie” w zależności czy odkodowany tekst jest identyczny z tekstem początkowym oraz odkodowany tekst w tym wypadku identyczny tekst jak w pierwszych 250 znakach z pliku 'norm_wiki_sample.txt'.

```
{'e': bytearray('000'), 'm': bytearray('00100'), 'y': bytearray('001010'), 'k': bytearray('0010110'),  
'4': bytearray('001011100'), 'x': bytearray('001011101'), '5': bytearray('001011110'), '3': bytearray('00  
1011111'), 's': bytearray('0011'), 'w': bytearray('010000'), 'b': bytearray('010001'), 'c': bytearray('01  
001'), 'r': bytearray('0101'), 'o': bytearray('0110'), 'n': bytearray('0111'), 'i': bytearray('1000'),  
'd': bytearray('10010'), '2': bytearray('10011000'), '9': bytearray('10011001'), 'v': bytearray('100110  
1'), 'g': bytearray('100111'), 't': bytearray('1010'), 'p': bytearray('101100'), 'f': bytearray('10110  
1'), 'l': bytearray('10111'), 'a': bytearray('1100'), 'h': bytearray('11010'), '8': bytearray('11011000  
0'), 'j': bytearray('110110001'), '0': bytearray('11011001'), 'q': bytearray('1101101000'), 'z': bitarra  
y('1101101001'), '6': bytearray('1101101010'), '7': bytearray('1101101011'), '1': bytearray('11011011'),  
'u': bytearray('110111'), ' ': bytearray('111')}
```

poprawnie

```
albert of prussia 17 may 1490 20 march 1568 was the last grand master of the teutonic knights who af  
ter converting to lutheranism became the first monarch of the duchy of prussia the secularized state  
that emerged from the former monastic state of t
```

Dodatkowo funkcja ta zwraca entropie i efektywność Huffmana które w tym przypadku wyniosły:

Entropia: 4.28

Efektywność 72.41%

Średnia długość słowa kodu: 5.91

W przypadku kodowania Huffmana średnia długość słowa jest mniejsza niż w przypadku stałej długości, co oznacza, że skompresowany plik będzie mniejszych rozmiarów.