Hubert Makowski

Sprawozdanie z 1 zadania laboratoryjnego z przedmiotu Teoria informacji i kodowania.

Treść zadania:

Implementacja kodera/dekodera słownikowego. Metoda: LZ77

Analiza zadania

Algorytm LZ77 jest to metoda strumieniowej słownikowej kompresji danych, polega na skracaniu zapisów tekstowych poprzez zastępowanie słów wskazanymi na odpowiednie hasła słownika.

Metoda LZ77 opiera się na założeniu, że słowa w kodowanym tekście powtarzają się. Słownik stanowi skończony ciąg ostatnio kodowanych symboli. Ciąg interpretujemy jako zbiór haseł, czyli dowolnych podciągów symboli słownika.

Założenia

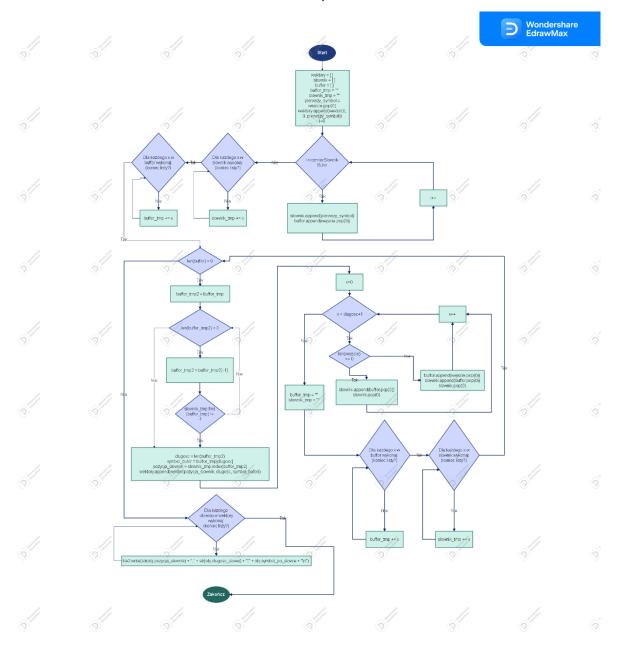
W zmiennej bufor zapisujemy aktualne kodowane symbole, zmienna wejście zawiera przeznaczona do kodowania sekwencje symboli wejściowych, słownik zawiera ostatnio kodowane symbole. Rozmiar bufora i słownika jest taki sam, a także jest on definiowany na ilość znaków w pliku wejściowym – 2.

Proces kodowania rozpoczyna się od inicjowania słownika. Następnie zmiennej słownik przypisuje się wartość początkowa sekwencje złożona z powtórzeń pierwszego kodowanego symbolu. Następnie symbol przesyłamy do dekodera, aby mógł on zainicjować słownik. W tym celu generujemy wektor kodowy w postaci <0, 0, s> gdzie s jest pierwszym elementem wejścia.

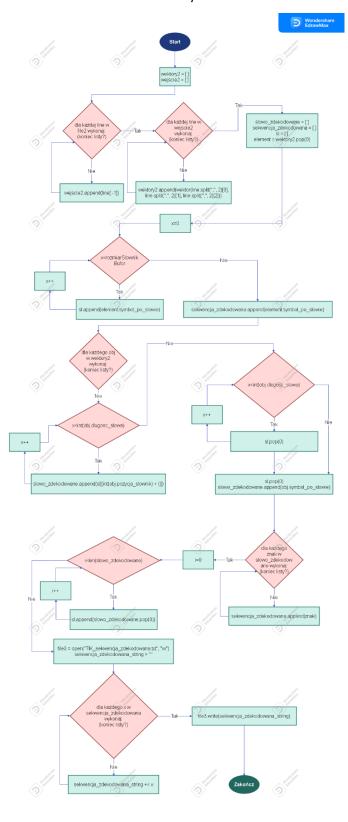
Algorytm

Schematy blokowe wykonałem w programie WondershareEdrawMax (darmowa licencja próbna, znak wodny).

Schemat blokowy kodera



Schemat blokowy dekodera



Implementacja

Początkowo algorytm chciałem implementować w języku c++, jednak ze względu na swoje ograniczenia związane z tym językiem, postanowiłem napisać kod w języku python. Korzystałem ze środowiska PyCharm community edition. Użytkownik podaje nazwę pliku wejściowego w terminalu (bez rozszerzenia .txt), nazwa pliku w którym zostaną zapisane wektory kodowe jest definiowana na TIK_output.txt, nazwa pliku do którego zostanie zapisana sekwencja zdekodowana jest zdefiniowana na TIK_sekwencja_zdekodowana.txt. Format zapisu wektorów kodowych ma postać -pozycja słowa w słowniku, długość słowa, symbol po słowie- (dane w pliku są oddzielone przecinkiem)

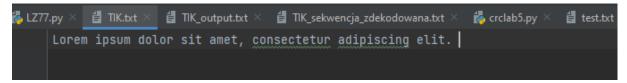
Testowanie poprawności działania

Podgląd terminala

```
C:\Users\Hubert\AppData\Local\Programs\Python\Python38-32\python.exe H:\pycharmproject\pythonProject\LZ77.py
Podaj nazwe pliku dla ktorego zostana zakodowane symbole: TIK

Process finished with exit code 0
```

Plik wejściowy TIK.txt



Plik z sekwencja zdekodowana TIK_sekwencja_zdekodowana.txt



Plik z zapisanymi wektorami kodowymi TIK_output.txt

```
🖧 LZ77.py ×

☐ TIK_output.txt

₫ TIK_sekwencja_zdekodowana.txt

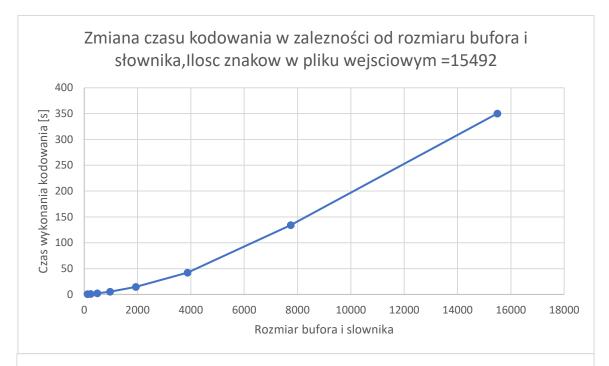
☐ TIK.txt

       0,0,L
       0,0,0
       0,0,r
       0,0,e
       0,0,m
       Θ,Θ,
       0,0,i
       Θ,Θ,ρ
       0,0,s
       Θ,Θ,υ
       48,2,d
       40,2,
       44,1,i
       0,0,t
       38,1,a
       35,1,e
       32,1,c
       31,1,e
       49,1,t
       43,2,U
       32,2,a
       25,1,i
       18,2,c
       14,1,n
       0,0,g
       10,1,e
       17,1,i
       21,1,.
```

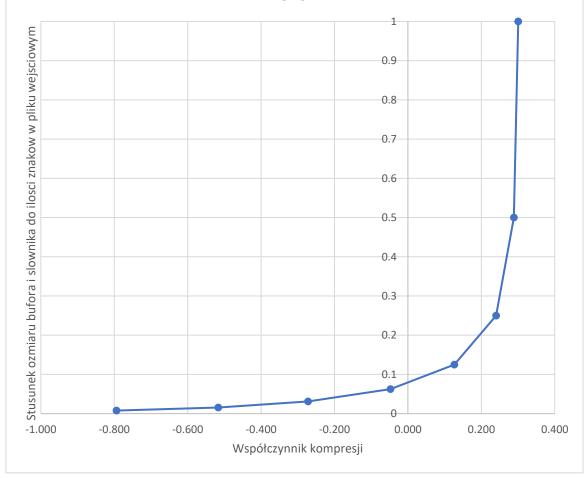
Testowanie stopnia kompresji

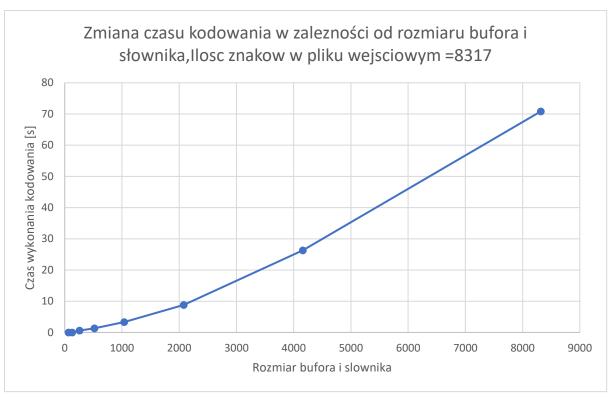
Wzór współczynnika kompresji-((Ilość znaków w pliku wejściowym- Ilość wektorów kodowych*5)/ Ilość znaków w pliku wejściowym), przy Ilości wektorów kodowych pojawia się współczynnik 5 ponieważ, format zapisu wektorów kodowych ma postać -pozycja słowa w słowniku, długość słowa, symbol po słowie- na co składa się 5 znaków.

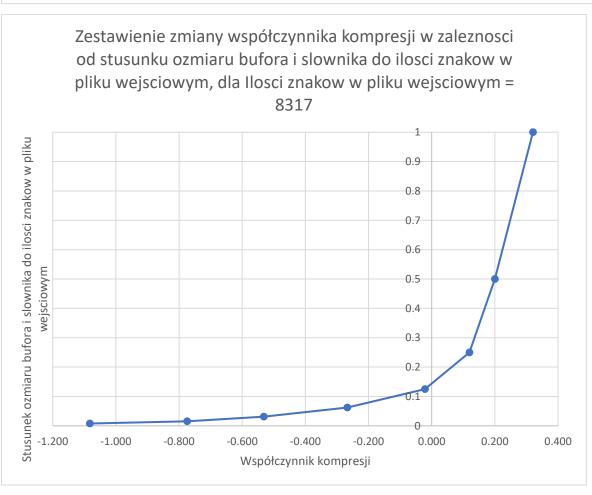
Tabela 1. Zestawienie wyników otrzymanych dla różnej ilości znaków w pliku wejściowym							
Rozmiar bufora i słownika	Ilość znaków w pliku wejściowym	llość wektorów kodowych	Współczynnik kompresji	Czas wykonania kodowania [s]	Czas wykonania dekodowania [s]		
15491	15492	2167	0.301	350.2	0.03		
7745	15492	2204	0.289	134.1	0.03		
3872	15492	2354	0.240	42.1	0.03		
1936	15492	2707	0.126	14.5	0.02		
968	15492	3246	-0.048	5.3	0.02		
484	15492	3941	-0.272	2.1	0.2		
242	15492	4699	-0.517	0.9	0.2		
121	15492	5559	-0.794	0.48	0.2		
8316	8317	1130	0.321	70.8	0.02		
4158	8317	1331	0.200	26.3	0.02		
2079	8317	1465	0.119	8.8	0.02		
1039	8317	1699	-0.021	3.3	0.02		
519	8317	2108	-0.267	1.3	0.02		
259	8317	2548	-0.532	0.59	0.02		
129	8317	2952	-0.775	0	0		
64	8317	3464	-1.082	0	0		
5144	5145	922	0.104	17	0.01		
2572	5145	950	0.077	6.6	0.01		
1286	5145	1089	-0.058	2.4	0.01		
643	5145	1273	-0.237	0.9	0.01		
300	5145	1514	-0.471	0.4	0.01		
150	5145	1787	-0.737	0.2	0.01		
100	5145	1948	-0.893	0	0		
50	5145	2270	-1.206	0	0		
20	5145	2790	-1.711	0	0		
10	5145	3330	-2.236	0	0		
1251	1252	335	-0.338	0	0		
625	1252	347	-0.386	0	0		
312	1252	379	-0.514	0	0		
156	1252	446	-0.781	0	0		
78	1252	499	-0.993	0	0		
39	1252	590	-1.356	0	0		
760	761	231	-0.518	0	0		
380	761	238	-0.564	0	0		
190	761	262	-0.721	0	0		
95	761	297	-0.951	0	0		
47	761	342	-1.247	0	0		
23	761	401	-1.635	0	0		
335	336	128	-0.905	0	0		
168	336	131	-0.949	0	0		
83	336	139	-1.068	0	0		
41	336	164	-1.440	0	0		
21	336	180	-1.679	0	0		



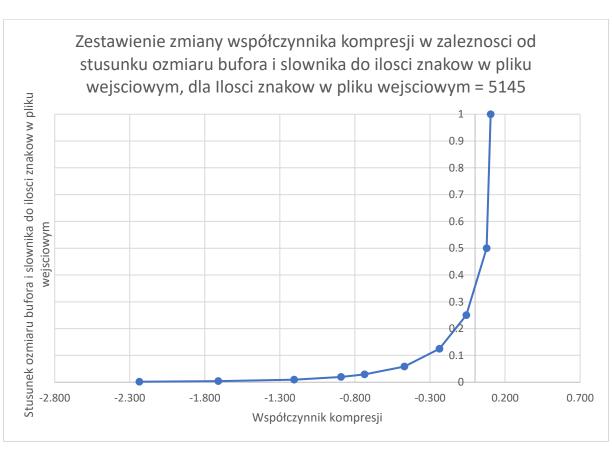
Zestawienie zmiany współczynnika kompresji w zaleznosci od stusunku ozmiaru bufora i slownika do ilosci znakow w pliku wejsciowym, dla Ilosci znakow w pliku wejsciowym = 15492











Wnioski

Z tabeli a także wykresów wynika, że największy współczynnik kompresji występuję dla rozmiaru słownika i bufora bliskiego ilości znaków w pliku wejściowym. Co ważne współczynnik kompresji rośnie wraz z wzrostem wielkości rozmiaru słownika i bufora (jest proporcjonalny do wielkości rozmiaru słownika i bufora).

Z danych w tabeli, a także wykresów wynika również, że empiryczna maksymalna wartość współczynnika kompresji dla danego algorytmu wynosi około 0.3, współczynnik kompresji asymptotycznie dąży do wartości 0.3.

Warto również zwrócić uwagę na czas kodowania, który rośnie proporcjonalnie z rozmiarem słownika i bufora. Jest to spowodowane tym, iż algorytm poświęca więcej czasu na przeszukiwanie danej sekwencji w słowniku, ze względu na większą ilość znaków w ciągu, zarówno w słowniku jak i buforze (funkcja find w python).

Program nie radzi sobie z znakiem nowej linii w pliku wejściowym, ze względu na zapis wektorów kodowych, zatem postanowiłem usunąć znaki nowej linii z listy *wejście* i nie kodować znaków nowej linii.

Złożoność czasowa dekodowania szacuje na T(n)=O(n²)

Złożoność czasowa kodowania szacuje na $T(n)=O(n^4)$ (funkcja find w python ma złożoność czasową $T(n)=O(n^2)$)

Wniosek końcowy – program nie sprawdzi się dla dużej ilości znaków w pliku wejściowym (czas wykonania programu jest nieakceptowalny).

Kod źródłowy programu

```
sl.append(element.symbol po slowie)
sekwencja zdekodowana.append(element.symbol po slowie)
        slowo zdekodowane.append(sl[(int(obj.pozycja slownik) + i)])
    slowo zdekodowane.append(obj.symbol po slowie) #dodaj do sekwencji
        sl.append(slowo zdekodowane.pop(0)) # dodaj do slownika slowo
    file3 = open("TIK sekwencja zdekodowana.txt", "w") # zapisz
file3.close()
wektory = []
buffor = []
for i in range (rozmiarSlownikBufor): #wypelnij slownik i bufor
    slownik.append(pierwszy symbol)
    buffor.append(wejscie.pop(0))
    buffor tmp2 = buffor tmp
```

```
symbol bufor = buffor tmp[dlugosc] #symbol po slowie
                slownik.append(buffor.pop(0)) #usun element z bufora i
                slownik.pop(0) #usun element ze slownika
                buffor.append(wejscie.pop(0)) #dodaj element na koniec
                slownik.append(buffor.pop(0)) #dodaj element na koniec
file name = input("Podaj nazwe pliku dla ktorego zostana zakodowane
wejscie = [] # lista w ktorej zostaną zapisane wektory kodowe
   char = file.read(1)
   wejscie.append(char) #dodaj znak do lisy
file.close()
ile znakow = wejscie.count('\n')
dekoduj("TIK output.txt", rozmiarSlownikBufor)
```