

Dokumentacja Egzaminu Praktycznego Algorytmy i Struktury Danych

Konrad Lubera, gr IB

21 maja 2020

Spis treści

1	Polecenie	3
2	Algorytm	3
2.1	Opis	3
2.2	Przykład działania	3
2.3	Schemat blokowy	4
3	Opis rozwiązania	5
3.1	Napotkane problemy	5
3.2	Pseudokod	5
4	Prezentacja wyników	7
5	Kod programu	9

1 Polecenie

Temat 9: Napisz program/aplikację który/która: Wyznaczy listę słów, które pojawiają się w tekście z jednego pliku, a nie pojawiają się w tekście z pliku drugiego przy pomocy algorytmu Boyer-Moore.

2 Algorytm

2.1 Opis

Do wykonania zadania użyto, zgodnie z poleceniem algorytmu Boyera-Moorea. Jest to wersja która korzysta z funkcji bad-character shift, bardzo podobna do tej która zastała zaprezentowana na wykładzie. Polega na porównywaniu z prawej strony wzorca, jeżeli wzorzec nie będzie pasował, to algorytm przeskakuje o długość wzorca. Jest to bardzo szybki i efektywny algorytm ponieważ ilość porównań jest relatywnie mała.

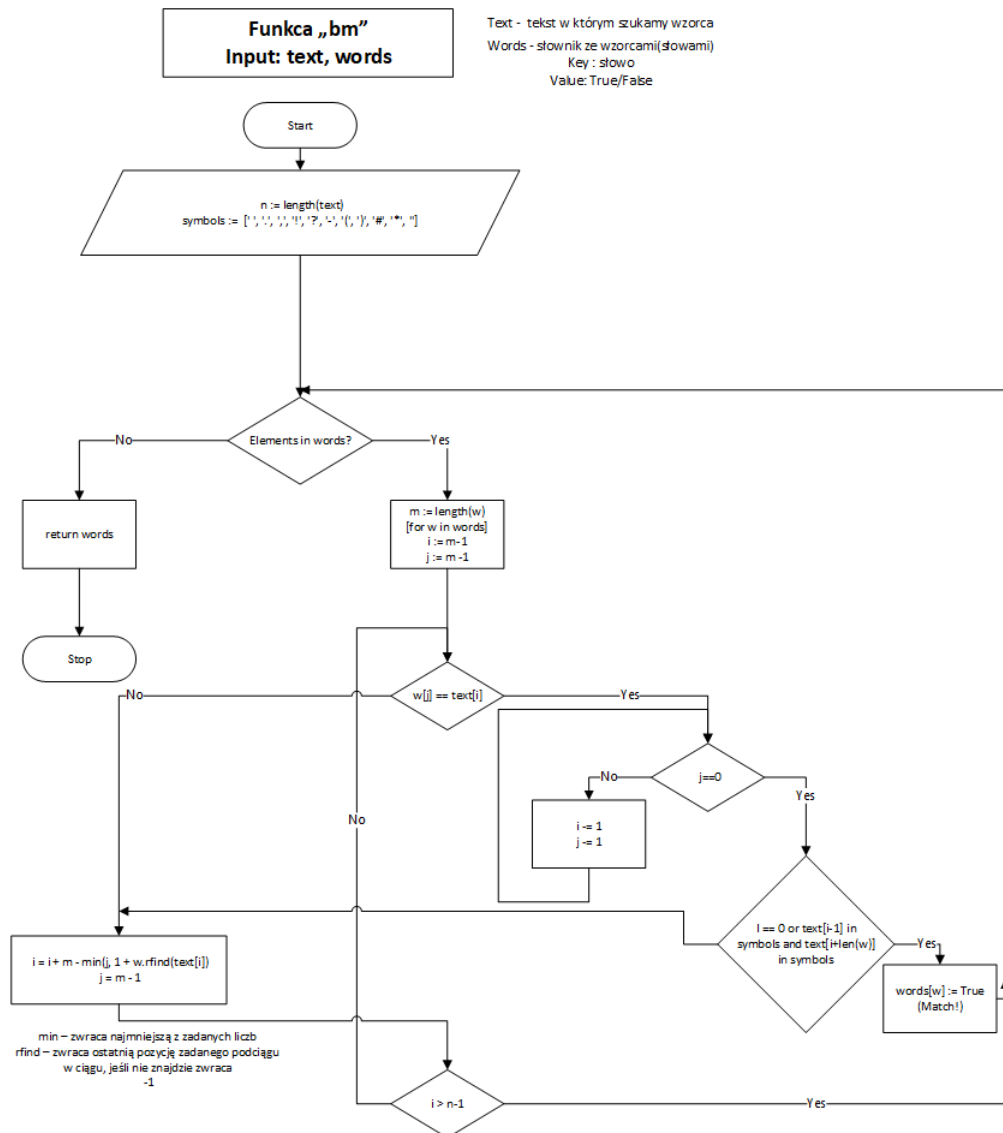
Implementację wykonano w języku Python w wersji 3.8.

2.2 Przykład działania

```
Tekst: Zdam szybko i efektywnie AiSD
      | | | | ||| |
efe | | | ||| |
      | | | ||| |
efe | | | ||| |
      | | | ||| |
efe | | | ||| |
      | | | ||| |
efe ||| |
      ||| |
efe|| |
      || |
efe| |
      | |
efe | i tak dalej...
      |
efe
```

Rysunek 1: Przykład działania algorytmu z notatek z wykładów

2.3 Schemat blokowy



Rysunek 2: Zastosowany algorytm Boyera Moorea

3 Opis rozwiązania

3.1 Napotkane problemy

Podczas implementacji napotkano problem związany z tym, że algorytm służy, w jego oryginalnym wydaniu, do wyszukiwania fraz-wzorca co generuje problemy w wyszukiwaniu słów. Dla przykładu gdy naszym wzorcem jest słowo „dam” a tekście mamy słowo „zdam” algorytm uzna, że wzorzec się pojawił. Problem ten rozwiązano dojadając warunek który sprawdza czy znak przed i za dopasowaniem znajduje się na liście która zawiera znaki interpunkcyjne, często używane symbole oraz spację. Dodatkowo aby wyeliminować problem z dopasowaniem słowa napisanego wielką literą i tego napisanego małą, wszystkie litery we wzorcu jak i tekście zostały zmienione na małe przy pomocy funkcji `lower()`.

3.2 Pseudokod

```
input : path
output: string text

open file on path as file:
    text = file.read
text = lower(text)
file.close
return text
```

Algorithm 1: Funkcja read to string

```
input : path
output: dictionary words

words:=
symbols := ‘.,!-?():’
open file on path as file:
    foreach line in file do
        foreach word in line do
            tmp:=word.strip(symbols)
            tmp = lower(tmp)
            if tmp not in words and tmp!=” then
                | words[tmp] = False ; /* przypisanie do słownika */
            end
        end
    end
return words
```

Algorithm 2: Funkcja read to dictionary

input : string *text*, dictionary *words*

output: updated dictionary *words*

```
n := lenght(text)
symbols := [ ' ', '.', ',', '!', '?', '-', '(', ')', '*', '']
foreach w in words do
    m := lenght(w)
    j := m-1
    i := m-1
    repeat
        if w[j]==text[i] then
            if j==0 then
                if i==0 or text[i - 1]insymbols and
                    text[i + len(w)]insymbols then
                    | words[w] = True
                    | break
                else
                    | i = i + m - min(j, 1+ w.find(text[i]))
                    | j = m-1
                end
            else
                | i -= 1
                | j -= 1
            end
        else
            | i = i + m - min(j, 1+ w.find(text[i]))
            | j = m-1
        end
    until i>n-1;
end
return words
```

Algorithm 3: Funkcja BM

Pseudokod został napisany tylko dla funkcji które są ważne z punktu widzenia działania samego algorytmu.

4 Prezentacja wyników

Teksty które zostały użyte do przetestowania programu zostały tak dobrane żeby sprawdzić czy problemy o których napisano w sekcji Algorytm - Opis zostały rozwiązane w rezultacie teksty mają dużo znaków interpunkcyjnych oraz słów które zawierają w sobie inne słowa. Teksty nie jest też zbyt długie żeby nie generować trudności podczas sprawdzania wyników, jednak przetestowano program używając dłuższych tekstów i też działa :D.

Oto teksty:

- file1.txt

Byliśmy bardzo smutni, że nie mogliśmy wrócić na uczelnię, jednak najważniejsze było żeby przetrwać ten trudny czas pandemii i co najważniejsze dobrze napisać egzamin z algorytmów i struktur danych. Mam nadzieję że to co zrobiłem wystarczy żebym zdał. Jeżeli nie będę bardzo zawiedziony! Kompletnie, nie, mam, pojęcia, jak, używać, znaków, interpunkcyjnych, więc, w, tym, zdaniu, przecinek, będzie, po, każdym, słowie. Umiejętności miękkie to nie jest moja mocna strona. Już nie wiem co tu napisać więc wkleję pierwszy akapit z Wikipedii o algorytmie Boyera Moore'a. Algorytm Boyera i Moore'a algorytm poszukiwania wzorca w tekście. Polega na porównywaniu, zaczynając od ostatniego elementu wzorca.

- file2.txt

Byliśmy bardzo zawiedzeni, że w tym semestrze nie wrócimy już na uczelnię, jednak pomimo tego ucze (specjalnie błąd) się pilnie i dam z siebie wszystko. Mam nadzieję, że taki stan rzeczy nie będzie trwać w nieskończoność. Są ważniejsze sprawy od spotkania z kolegami np. egzamin z aisd który mimo trudności się obył w bardzo przystępnej formie. Nie będę tu już więcej pisać jakichś głupot tylko wkleję znów akapit z Wikipedii i kilka słów z generatora losowych słów ów. Algorytm Boyera i Moore'a algorytm poszukiwania wzorca w tekście. Polega na porównywaniu, zaczynając od ostatniego elementu wzorca. Rzeczownik, przestrzeń, biblioteka, rosa parks, szubienica, ośmiornica, wzmacniacz, atletyczny, farmaceuta, jeżozwierz, frustracja, bezradność, anne frank, niedźwiedź, rękawiczki, kierownica, studzienka, dominikana, ciężarówka, sprzedawca.

To lista słów które pojawiają się w tekście z jednego pliku, a nie pojawiają się w tekście z pliku drugiego:

1. Zawiedzeni	23. Kolegami	45. Biblioteka
2. Semestrze	24. Np	46. Rosa
3. Wrócimy	25. Aisd	47. Parks
4. Pomimo	26. Który	48. Szubienica
5. Tego	27. Mimo	49. Ośmiornica
6. Ucze	28. Trudności	50. Wzmacniacz
7. Specjalnie	29. Obył	51. Atletyczny
8. Błąd	30. Przystępnej	52. Farmaceuta
9. Się	31. Formie	53. Jeżozwierz
10. Pilnie	32. Więcej	54. Frustracja
11. Dam	33. Pisać	55. Bezradność
12. Siebie	34. Jakichś	56. Anne
13. Wszystko	35. Głupot	57. Frank
14. Taki	36. Tylko	58. Niedźwiedź
15. Stan	37. Znów	59. Rękawiczki
16. Rzeczy	38. Kilka	60. Kierownica
17. Trwać	39. Słów	61. Studzienka
18. Nieskończoność	40. Generatora	62. Dominikana
19. Są	41. Losowych	63. Ciężarówka
20. Ważniejsze	42. Ów	64. Sprzedawca
21. Sprawy	43. Rzeczownik	65. Przeżyć
22. Spotkania	44. Przestrzeń	66. Samochód

Po wykonaniu programu otworzy się plik txt który zawiera tą listę oraz dodatkową statystykę. Jeżeli to nie nastąpi np gdy używamy systemu Linux wystarczy otworzyć plik „results.txt”. Istnieje możliwość zmiany tekstów na których operuje program wystarczy wyedytować zawartość plików „file1.txt” i „file2.txt”.

5 Kod programu

```
# -*- coding: utf-8 -*-
# Wyznaczy ęlist łósw, óktrę ąpojawią ęsi w ńtekcie z
# jednego pliku, a nie ąpojawią ęsi w ńtekcie z pliku
# drugiego przy pomocy algorytmu Boyer-Moore
import os

def read_to_string(path):
    with open(path, 'r', encoding="utf-8") as file:
        text = file.read()
        text = text.lower()
        text += "_" # dodatkowe zabezpieczenie
        #print(text)
        file.close()
    return text

def read_to_dict(path):
    words = {}
    symbols = '_.!-?() "\n' # warunki dla - w
    # ńrodku ! - strip łzaatwia sprwawe
    with open(path, 'r', encoding="utf-8") as file:
        for line in file:
            for word in line.split():
                tmp = word.strip(symbols)
                tmp = tmp.lower()
                # TYMCZASOWE ĄROZWIZANIE!
                #t = ' ' + tmp + ' '
                if tmp not in words and tmp != "":
                    words[tmp] = False
    # print("ńć Ilo łósw bez óńpowtrze: " + str(len(words)))
    ))
```

```

    #print(words.keys())
    file.close()
    return words

def bm(text, words):
    n = len(text)
    symbols = ['_', '.', ',', '!', '?', '-', '(', ')',
               '#', '*', '']
    for w in words:
        m = len(w)
        i = m - 1
        j = m - 1
        while True:
            if w[j] == text[i]:
                if j == 0:
                    if i == 0 or text[i-1] in symbols
                    and text[i+len(w)] in symbols:
                        words[w] = True
                        # print("Matched word: " + w +
                        "' on position: " + str(i))
                        # if i == 0: print("Before:
                        None")
                        # else: print("Before: " + text
                        [i-1])
                        # print("After: " + text[i +
                        len(w)])
                        break
                    else:
                        i = i + m - min(j, 1 + w.rfind(
                            text[i]))
                        #i = i + m - min(j, 1 + w.rfind
                        (text[i+len(w)])
                        j = m - 1
                else:
                    i -= 1
                    j -= 1
            else:
                i = i + m - min(j, 1 + w.rfind(text[i]))
                j = m - 1

```

```

        if i > n - 1:
            break
    return words

def count_words(path):
    words = []
    symbols = '_.!-?()"\n_\'
    with open(path, 'r', encoding="utf-8") as file:
        for line in file:
            for word in line.split():
                tmp = word.strip(symbols)
                tmp = tmp.lower()
                if tmp not in words and tmp != "":
                    words.append(tmp)
    file.close()
    return len(words)

def results(words_1, num_of_word_2):
    i = 0
    j = 0
    for w in words_1:
        if words_1[w]:
            j += 1
        elif not words_1[w]:
            # print(w)
            i += 1
    c = int(num_of_word_2) - j

    # print("ŚcIlo łósw w pliku 1. (bez óńpowtrze): " +
    #       str(len(words_1)))
    # print("ŚcIlo łósw w pliku 2. (bez óńpowtrze): " +
    #       str(num_of_word_2))
    # print("ŚcIlo łósw unikalne dla pliku 1: " + str(i)
    #       )
    # print("ŚcIlo łósw unikalne dla pliku 2: " + str(c)
    #       )
    # print("ŚcIlo łósw z pliku 1 znalezionych w pliku
    #       2: " + str(j))

```

```
save_results(words_1, num_of_word_2, i, c, j)
```

```
def save_results(words, n2, unique_1, unique_2, matches
):
    file = open("results.txt", "w")

    file.write("Egzamin_Praktyczny_AiSD_Konrad_Lubera_\n")
    file.write("Temat: Wyznaczenie list słów, które\n\
    pojawiają się w słankie z jednego pliku, a nie\n\
    pojawiają się w słankie z pliku drugiego przy\n\
    pomocy algorytmu Boyer-Moore\n")
    file.write("_____\n")

    file.write("Wyniki:\n")
    file.write("_____\n")
    file.write("Ilo słów w pliku 1. (bez óńpowtrze): \n"
    + str(len(words)) + "\n")
    file.write("Ilo słów w pliku 2. (bez óńpowtrze): \n"
    + str(n2) + "\n")
    file.write("_____\n")
    file.write("Słowa unikalne dla pliku 1: " + str(
    unique_1) + "\n")
    file.write("Słowa unikalne dla pliku 2: " + str(
    unique_2) + "\n")
    file.write("_____\n")
    file.write("Ilo słów z pliku 1 znalezionych w\n\
    pliku 2: " + str(matches) + "\n")
    file.write("_____\n")

    file.write(
    "Słowa które pojawiają się w słankie z jednego\n\
    pliku, a nie pojawiają się w słankie z pliku\n\
    drugiego:\n")
    i = 1
    for w in words:
        if not words[w]:
            file.write(str(i) + ". " + str(w.strip())
```

```
        capitalize()) + "\n")
    i += 1

file.close()

results(bm(read_to_string("file1.txt"), read_to_dict("
    file2.txt")), count_words("file1.txt"))
os.system('start_notepad_results.txt')
```