

**Michał Szczurek - lab 2**

## 1. Przygotowanie danych

Do ostatniego tekstu dodano na koniec \$ jako marker. W pozostałych rolę markera spełnia znak d

```
In [5]: f = open("1997_714_head.txt", "r", encoding="UTF-8")
data = ["bbbd","aababbd","ababcd","abcbcd","abccdd","f.read()"]
data[4] = data[4] + "S"
```

```
In [6]: print(data)

['bbbd', 'aababbd', 'ababcd', 'abcbcd', '\n\n\n\n\nDz.U. z 1998 r. Nr 144, poz. 930'\n
\n\n\n\n\nUSTAWA'\n
\n\n\n\n\nz dnia 20 listopada 1998 r.\n
\n\n\n\n\nosiąganych przez oso
by fizyczne\n
Rozdział 3\n
Przepisy ogólne\n
Art. 1.\nustawa reguluj
e opodatkowanie zryczałtowanym podatkiem dochodowym niektórych\nprzychodów (dochodów) osiągniętych przez osoby fizyczne
prowadzące pozarolniczą działalność gospodarczą oraz przez osoby duchowne.\n
Art. 2.\nOsoby fizyczne osiągające przychody z pozarolniczej działalności
gospodarczej opłacają zryczałtowany podatek dochodowy w formie:\n
1) ryczałtu od przychodów ewidencjonowanych,\n
2) karty podatkowej.\nOsoby duchowne, prawnie uznanych wyznawców, opłacają zryczałtowany podatek dochodowy od przy
chodów osób duchownych.\nOsoby wpływy z podatku dochodowego opłacanego w formie ryczałtu od przychodów ewidencjonowanych
oraz zryczałtowanego podatku dochodowego od przychodów\nosób duchownych stanowią dochód budżetu państwa.\nOsoby wpływy z karty podatkowej stanowią dochody gmin.\n
Art. 3.\nPrzychody (dochody) opodatковanych w formach zryczałtowych nie łączą się z\nprzychodami (dochodami) z innych źródeł podlegających opodatkowaniu na podstawie ustawy z dnia 28 lipca 1991 r. o podatku dochodowym od osób fi
zycznych (Dz. U. z 1993 r. Nr 98, poz. 416 i Nr 134, poz. 646, z 1994 r. Nr 43, poz. 163, Nr 90, poz. 411, Nr 113, p
oz. 547, Nr 123, poz. 682 i Nr 126, npoz. 626, z 1995 r. Nr 5, poz. 25 i Nr 133, poz. 654, z 1996 r. Nr 25, poz. 4nn11
3, Nr 87, poz. 395, Nr 137, poz. 638, Nr 147, poz. 686 i Nr 156, poz. 776, z 1997 r. Nr 28, poz. 153, Nr 30, poz. 16
4, Nr 71, poz. 449, Nr 85, poz. 538, Nr 96, poz. 592, Nr 121, poz. 770, Nr 123, poz. 776, Nr 137, poz. 926, Nr 13
7, poz. 932-934 i Nr 141, poz. 943 i 945 oraz z 1998 r. Nr 66, poz. 470, Nr 74, npoz. 471, Nr 108, poz. 685 i Nr 11
7, rozp., zmiana)\nustawa o podatkach dochodowych".\n']
```

## 2. Struktura Trie

```
In [7]: class Trie:
def __init__(self, char, depth):
    self.children = {}
    self.char = char
    self.depth = depth # used for giving node unique label, in order to draw tree plot (which did not work in th
e end ;))
    # self.end = False if we want to have only suffixes in trie

def add_child(self, char):
    d[char] = Node(char)

def has_child(self, char):
    return char in self.children

In [8]: class Trie:

def __init__(self, text):
    self.root = Node("$", 0)
    for i in range(len(text)-1, -1, -1):
        self.add(text[i:])
        #print(text[i:])

def add(self, text):
    i = 0
    node = self.root
    while i < len(text) and not node.has_child(text[i]):
        node = node.children[text[i]]
        i += 1
    while i < len(text):
        node.children[text[i]] = Node(text[i], node.depth + 1)
        node = node.children[text[i]]
        i += 1
    node.end = True

def contains(self, text):
    i = 0
    node = self.root
    while i < len(text) and node.has_child(text[i]):
        prev = node
        node = node.children[text[i]]
        i += 1
    return i == len(text) # and node.end == True
```

### 3. Drzewo suffiksowe

```

In [9]: class Suffix_node:
def __init__(self, start, end):
    self.start = start
    self.end = end
    self.children = {}

def __repr__(self):
    return f"({self.start}, {self.end})"
def __str__(self):
    return f"({self.start}, {self.end})"

def print(self):
    for key, child in self.children.items():
        print(child)
    for key, child in self.children.items():
        print("=====")
        print("Children of: " + str(child))
        child.print()

In [10]: class Suffix_tree:
def __init__(self, text):
    self.root = Suffix_node(-1, -2) # special values for root it is important that start > end for "contains" to
    work

    self.text = text
    for i in range(len(text)-1, -1, -1):
        self.add(i, self.root)
        # self.root.print() debug
        # print("=====")

def add(self, curr_id, curr_node):
    next_node = curr_node.children.get(self.text[curr_id])

    if next_node is None: # create new node representing text[i:]
        new_node = Suffix_node(curr_id, len(self.text)-1)
        curr_node.children[self.text[curr_id]] = new_node
        return

    i = next_node.start
    while i <= next_node.end and curr_id < len(self.text) and self.text[curr_id] == self.text[i]:
        i += 1
        curr_id += 1

    if curr_id == len(self.text): # whole new suffix is already represented
        return # new suffix already in tree, nothing to do

    elif i == next_node.end + 1: # index out of nodes range
        self.add(curr_id, next_node) # add remaining part of suffix to child
        return

    else: # text[curr_id] != text[i], we need to split branch
        old_branch = Suffix_node(i, next_node.end)
        old_branch.children = next_node.children
        new_branch = Suffix_node(curr_id, len(self.text)-1)
        next_node.end = i-1
        next_node.children = {}
        next_node.children[self.text[curr_id]] = new_branch
        next_node.children[self.text[i]] = old_branch
        return

def contains(self, pattern, curr_node = None):
    if curr_node == None:
        curr_node = self.root

    i = 0
    while i + curr_node.start <= curr_node.end and i < len(pattern):
        if self.text[curr_node.start + i] != pattern[i]:
            return False
        else:
            i += 1

    if i == len(pattern):
        return True

    else: # search in child
        next_node = curr_node.children.get(pattern[i])
        if next_node == None:
            return False
        return self.contains(pattern[i:], next_node)

```

#### 4. Testowanie poprawności algorytmów

#### 4.0 Uniwersalna funkcja testująca

```

in [11]: def test(trie, suffix, expected_true, expected_false):
        for pattern in expected_true:
            if trie.contains(pattern) is False:
                print("Error FALSE in trie for " + pattern)
                return False
            if suffix.contains(pattern) is False:
                print("Error FALSE in suffix tree for " + pattern)
                return False

        for pattern in expected_false:
            if trie.contains(pattern) is True:
                print("Error TRUE in trie for " + pattern)
                return False
            if suffix.contains(pattern) is True:
                print("Error TRUE in suffix tree for " + pattern)
                return False

        print("OK")
        return True

```

#### 4.1 Testy dla 1. tekstu:

```
In [12]: a_trie = Trie(data[0])
a_suffix = Suffix_tree(data[0])
expected_true = ["bb", "bd", "bbb", "bbbd", "b"]
expected_false = ["a", "db", "bbbbbb"]
test(a_trie, a_suffix, expected_true, expected_false)

OK

Out[12]: True
```

#### 4.2 Testy dla 2. tekstu:

```
In [13]: b_trie = Trie(data[1])
         b_suffix = Suffix_tree(data[1])
         expected_true = [data[1][i:j] for i in range(len(data[1])) for j in range(len(data[1]))]
         expected_false = ["aaa", "db", "bbbbbbb", "aba"]
         test(b_trie, b_suffix, expected_true, expected_false)

OK
```

```
4.3 Testy dla 3. tekstu:
In [14]: c_true = Trie(data[2])
          c_suffix = Suffix_tree(data[2])
          expected_true = [data[2][i:j] for i in range(len(data[2])) for j in range(len(data[2]))]
          expected_false = ["aaa", "db", "bbb", "abaaa", "dc", "baba"]
          test(c_true, c_suffix, expected_true, expected_false)

OK

Out[14]: True
```

#### 4.4 Testy dla 4. tekstu:

```
In [15]: d_trie = Trie(data[3])
         d_suffix = Suffix_tree(data[3])
         expected_true = [data[3][i:j] for i in range(len(data[3])) for j in range(len(data[3]))]
         expected_false = ["aaa", "db", "bbb", "aba", "dc", "baba", "cccc"]
         test(d_trie, d_suffix, expected_true, expected_false)

OK

Out[15]: True
```

#### 4.5 Testy dla 5. tekstu:

```
In [16]: e_trie = Trie(data[4])
e_suffix = Suffix_tree(data[4])
expected_true = ["anych oraz zryczaltowanego podatku dochodowego od przychodów",
                 "Nr 133, poz. 654, z 1996 r. Nr 25,",
                 "20 listopada 1998 r.",
                 "a"]
expected_false = ["aaa", "db", "bbb", "aba", "dc", "baba", "cccc"]
test(e_trie, e_suffix, expected_true, expected_false)

OK

Out[16]: True
```

#### 4. Badanie czasów budowy struktur

```
In [17]: print("TESTY DLA DRZEWY TRIE\n")
        for i in range(len(data1)):
            print(f"Test prędkości dla {i+1}. tekstu")
            %timeit Trie(data1[i])

TESTY DLA DRZEWY TRIE

Test prędkości dla 1. tekstu
9.63 µs ± 599 ns per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 2. tekstu
21.8 µs ± 1.6 µs per loop (mean ± std. dev. of 7 runs, 10000 loops each)
Test prędkości dla 3. tekstu
19.3 µs ± 1.83 µs per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 4. tekstu
24.9 µs ± 1.72 µs per loop (mean ± std. dev. of 7 runs, 10000 loops each)
Test prędkości dla 5. tekstu
3.07 s ± 139 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [18]: print("TESTY DLA DRZEWY SUFIKSWOWEGO\n")
        for i in range(len(data1)):
            print(f"Test prędkości dla {i+1}. tekstu")
            %timeit Suffix_tree(data1[i])

TESTY DLA DRZEWY SUFIKSWOWEGO

Test prędkości dla 1. tekstu
7.38 µs ± 272 ns per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 2. tekstu
13.1 µs ± 1.3 µs per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 3. tekstu
9.6 µs ± 384 ns per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 4. tekstu
11.1 µs ± 539 ns per loop (mean ± std. dev. of 7 runs, 100000 loops each)
Test prędkości dla 5. tekstu
25.7 ms ± 936 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)
```

Z powyższych testów wynika, że czas potrzebny na stworzenie drzewa sufikсового jest istotnie mniejszy od czasu budowy drzewa Trie pomimo takiej samej złożoności teoretycznej. Ponadto warto zwrócić uwagę, na fakt, że drzewa sufikсовые zajmują istotnie mniej pamięci.