

Raport z laboratorium 1

Filip Nikolow

23 marca 2021

1 Cel laboratorium

Celem laboratorium było zapoznanie się z drzewami suffiksowymi - trie oraz skompresowanymi, a także z metodami ich konstrukcji.

2 Realizacja

Zaimplementowałem budowę drzew suffiksowych w trzech wersjach:

- Drzewo trie (pamięć: $O(n^2)$, obliczenia: $O(n^2)$)
- Skompresowane drzewo suffiksowe z wykorzystaniem tylko procedury `słow_find` do znajdowania głów (pamięć: $O(n)$, obliczenia: $O(n^2)$)
- Po długiej walce: algorytm McCreighta (pamięć: $O(n)$, obliczenia: $O(n)$)

Przy implementacji algorytmu McCreighta tegoroczny wykład był znacznie bardziej pomocny niż zeszłoroczny (być może jest to też kwestia tego, że już przed wtorkowym wykładem sporo nad tym algorytmem myślałem). Kluczowy natomiast do doszlifowania algorytmu i usunięcia bugów był niezawodny MIMUW :) (smurf.mimuw.edu.pl).

3 Realizacja poszczególnych poleceń

3.1 Polecenia 1-2

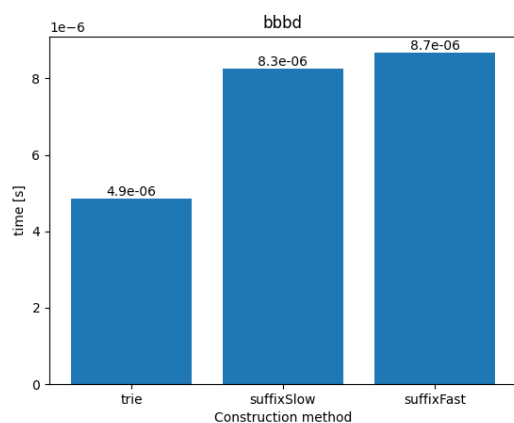
Łańcuchy 1-4 posiadają unikalny marker, natomiast do fragmentu ustawy na koniec dodałem znak Unicode "F000".

3.2 Polecenia 3-5

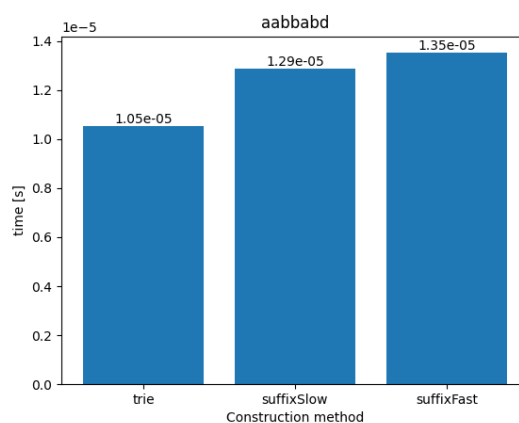
Zaimplementowałem algorytmy wyszczególnione w sekcji 2. Poprawność utworzonych struktur testowałem głównie wizualnie, poprzez porównanie drzewa trie z oczekiwanym, później porównanie drzewa trie z drzewem skompresowanym policzonym `słow_findem`, a następnie porównałem wyniki algorytmu ze `słow_findem` z algorytmem McCreighta. Nie wykryłem różnic - obrazki przedstawiające wyliczone drzewa dla przykładów z podzadania 1 znajdują się w sekcji "Drzewa". Dodatkowo załączam cały użyty kod na końcu raportu.

3.3 Testy porównujące szybkość działania algorytmów - polecenie 6

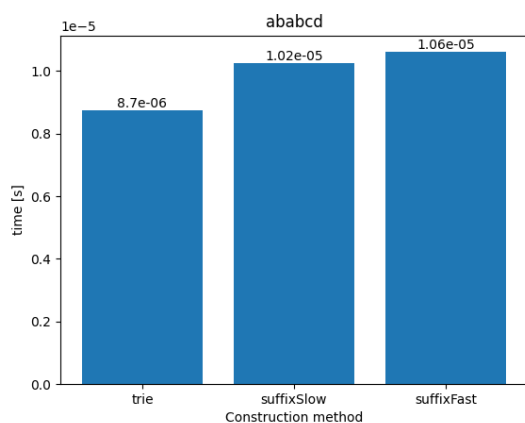
Poniżej załączam pomiary czasów działania każdego z trzech algorytmów: trie, algorytm ze `słow_findem` (suffixSlow), algorytm McCreighta (suffixFast). Testy przeprowadziłem na każdym z pięciu wymienionych w podzadaniu 1 tekstów. Każdy pomiar powtórzyłem 100-krotnie, aby zmniejszyć losowość przy krótkich tekstach.



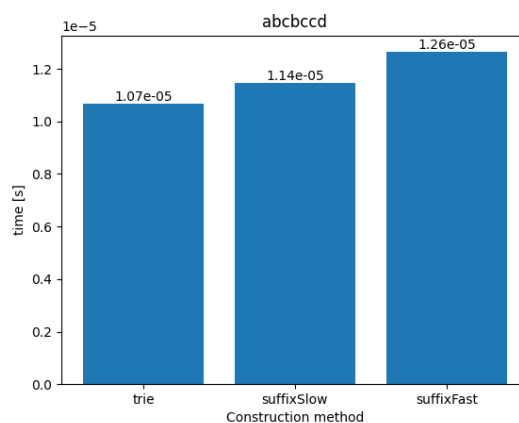
Rys. 1



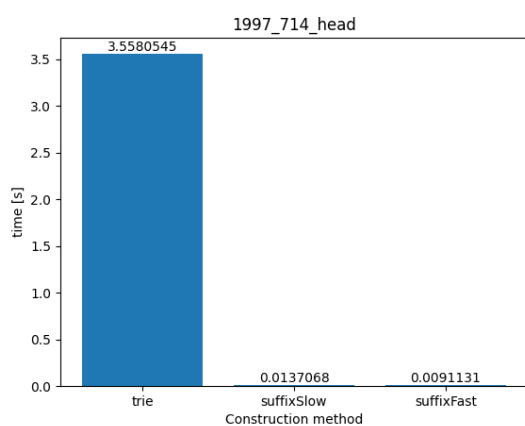
Rys. 2



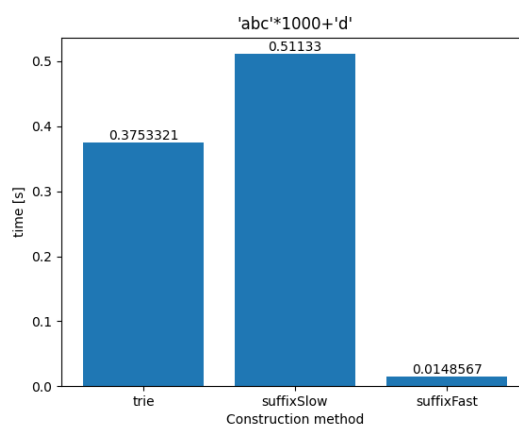
Rys. 3



Rys. 4



Rys. 5



Rys. 6

Rysunki od 1 do 4 pokazują, że na tak krótkich testach mała stała w algorytmie budującym trie jest bardziej znacząca niż dobra złożoność obliczeniowa i algorytm budujący trie wygrywa, natomiast algorytm McCreighta, właśnie przez większą stałą, jest ostatni. Na rysunku 5 kwadratowa złożoność algorytmu budującego trie jest już bardzo widoczna, natomiast ciekawe jest to, że suffixSlow tylko nieznacznie przegrywa z McCreightem - trie oprócz wysokiej złożoności ma także duży narzut związany z tworzeniem bardzo wielu obiektów (Node'ów), ale mimo to

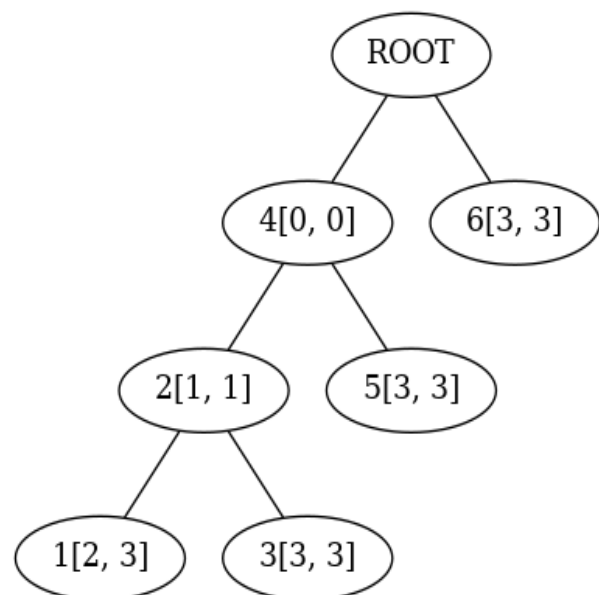
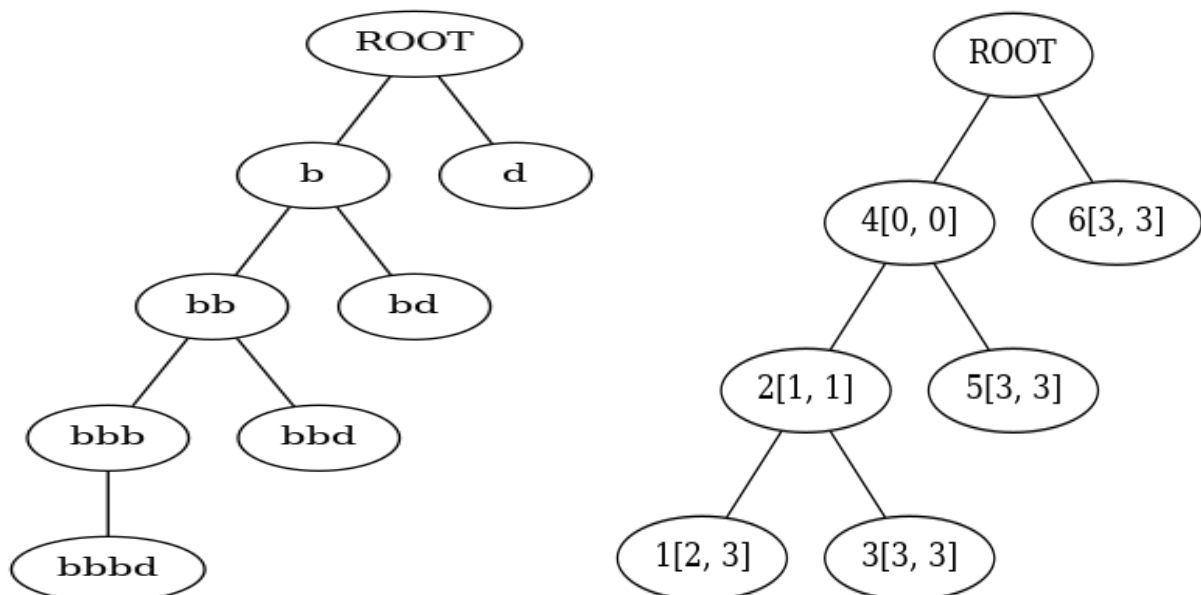
ewidentnie przykład ustawy jest korzystny dla funkcji `slov_find` (dodatkowo pewnie optymalizacje operacji na stringach w Pythonie pomagają) i linki w McCreighcie nie dają wiele - pomimo to McCreight jest szybszy mimo większej stałej. Aby sprawdzić czy faktycznie `suffixSlow` ma kwadratową złożoność (a zaimplementowany McCreight liniową) przeprowadziłem dodatkowy test z rysunku 6 - tam już wyniki są takie jakich można oczekiwać od strony teoretycznej :).

4 Szczegóły techniczne platformy testowej

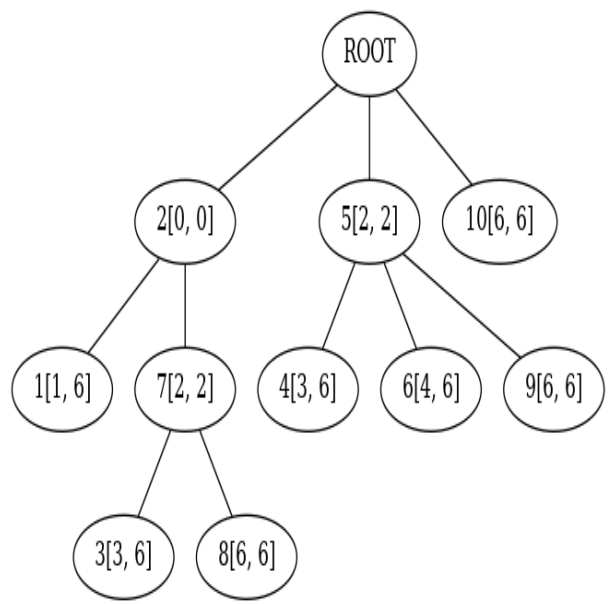
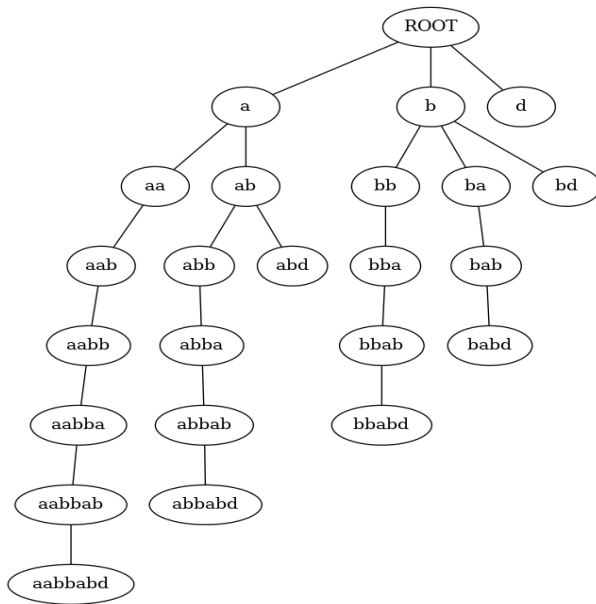
Wszystkie pomiary czasów zostały przeprowadzone na maszynie wirtualnej VirtualBox z systemem Ubuntu 20. Używana wersja języka Python to 3.8.5. Użyty procesor to AMD Ryzen 3900X.

5 Drzewa

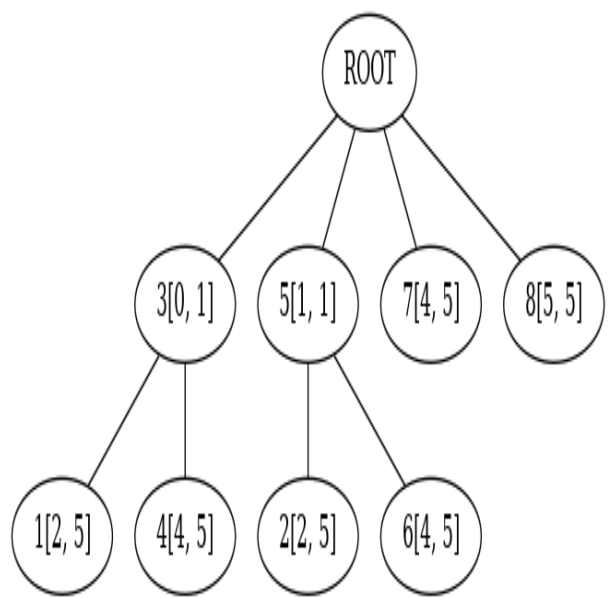
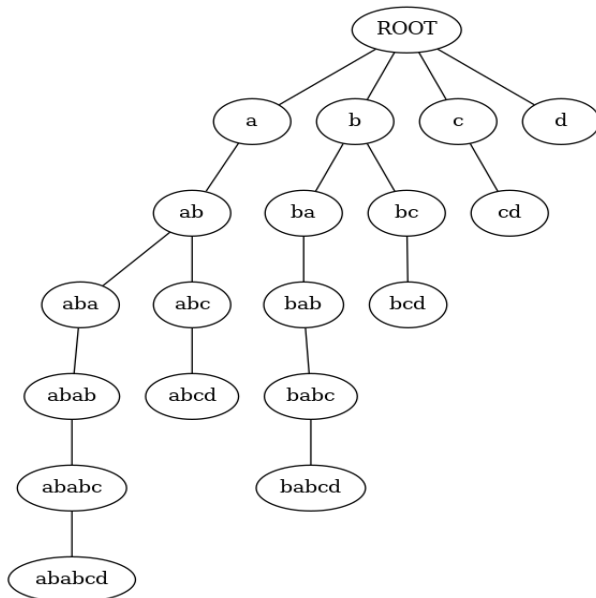
5.1 bbbd



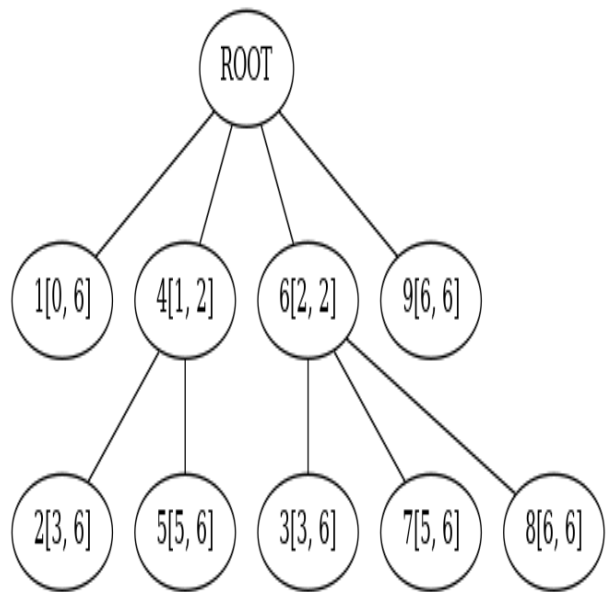
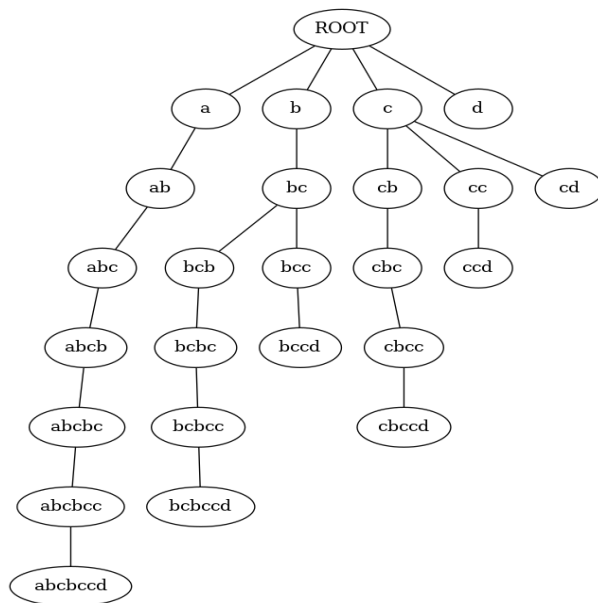
5.2 aabbabd



5.3 ababcd



5.4 abcdccd



6 Kod programów

6.1 Trie

```
1 import pydot
2 import tempfile
3 from PIL import Image
4
5
6 class Node:
7
8     def __init__(self):
9         self.children = dict()
10
11
12 def create_trie(text):
13     trie = Trie()
14     for i in range(len(text)):
15         trie.insert(text[i:])
16     return trie
17
18
19 class Trie:
20
21     def __init__(self):
22         self.root = Node()
23
24     def pretty_print(self, name=None, display=True):
25
26         def dfs_helper(node, label, graph):
27             for c in node.children.keys():
28                 graph.add_edge(pydot.Edge(label, label + c))
```

```

29         dfs_helper(node.children[c], label + c, graph)
30
31     graph = pydot.Dot(graph_type='graph')
32     for c in self.root.children.keys():
33         graph.add_edge(pydot.Edge("ROOT", c))
34         dfs_helper(self.root.children[c], c, graph)
35
36     if name is None:
37         fout = tempfile.NamedTemporaryFile(suffix=".png")
38         name = fout.name
39     else:
40         name += "_trie.png"
41     graph.write(name, format="png")
42     if display:
43         Image.open(name).show()
44
45     def insert(self, string):
46         p = self.root
47         for i, c in enumerate(string):
48             if c in p.children:
49                 p = p.children[c]
50             else:
51                 break
52         for j in range(i, len(string)):
53             c = string[j]
54             p.children[c] = Node()
55             p = p.children[c]
56
57     def search(self, string):
58         p = self.root
59         for c in string:
60             if c in p.children:
61                 p = p.children[c]
62             else:
63                 return False
64         return True
65
66
67     def create_treeplots(textlist):
68         for text in textlist:
69             T = create_trie(text)
70             T.pretty_print(text, False)
71
72
73     if __name__ == '__main__':
74         create_treeplots(["bbbd", "aabbabd", "ababcd", "abcbccd"])
75         # T = create_trie("aabbabbc")
76         # print(T.search("tki"))
77         # print(T.search("tkitek"))
78         # print(T.search("krotkitekst"))
79         # print(T.search("otko"))
80         # print(T.search("otkt"))
81         # T.pretty_print()

```

6.2 suffixSlow oraz McCreight

```
1  import pydot
2  import tempfile
3  from PIL import Image
4
5  ind = 0
6
7
8  class Node:
9
10     def __init__(self, label):
11         global ind
12         self.ind = ind
13         # print(ind)
14         ind += 1
15         self.label = label
16         self.letter = None
17         self.link = None
18         self.parent = None
19         self.children = dict()
20
21     def connect(self, child, key):
22         self.children[key] = child
23         child.parent = self
24         child.letter = key
25         return child
26
27     def length(self):
28         return self.label[1] - self.label[0] + 1
29
30
31 class suffixTree:
32
33     def __init__(self, text):
34         self.text = text
35         self.root = Node(None)
36         self.root.parent = self.root
37         self.root.link = self.root
38
39     def pretty_print(self, name=None, display=True):
40
41         def dfs_helper(node, graph):
42             for c in node.children.values():
43                 graph.add_edge(
44                     pydot.Edge(str(node.ind) + str(node.label),
45                               str(c.ind) + str(c.label)))
46                 dfs_helper(c, graph)
47
48         graph = pydot.Dot(graph_type='graph')
```

```

49     for c in self.root.children.values():
50         graph.add_edge(pydot.Edge("ROOT", str(c.ind) + str(c.label)))
51         dfs_helper(c, graph)
52     if name is None:
53         fout = tempfile.NamedTemporaryFile(suffix=".png")
54         name = fout.name
55     else:
56         name += "_tree.png"
57     graph.write(name, format="png")
58     if display:
59         Image.open(name).show()
60
61     def split(self, node, key, index):
62         child = node.children[key]
63         new_node = Node([child.label[0], index])
64         child.label[0] = index + 1
65         new_node.connect(child, self.text[index + 1])
66         node.connect(new_node, key)
67         return new_node
68
69     def fast_find(self, i, j, node):
70         # print("FAST")
71         if i > j:
72             return node
73         child = node.children[self.text[i]]
74         if (j - i + 1 > child.length()):
75             return self.fast_find(i + child.length(), j, child)
76         elif (j - i + 1 == child.length()):
77             return child
78         else:
79             return self.split(node, self.text[i], child.label[0] + j - i)
80
81     def slow_find(self, i, node):
82         # print("SLOW")
83         key = self.text[i]
84         if key not in node.children:
85             return (node, i)
86         child = node.children[key]
87         start, end = child.label
88         while start <= end:
89             if self.text[i] != self.text[start]:
90                 return (self.split(node, key, start - 1), i)
91             i += 1
92             start += 1
93         return self.slow_find(i, child)
94
95     def find(self, pattern, i, node):
96         key = pattern[i]
97         if key not in node.children:
98             return False
99         child = node.children[key]
100         start, end = child.label
101         while start <= end and i < len(pattern):

```



```

102         if pattern[i] != self.text[start]:
103             return False
104         i += 1
105         start += 1
106     if i == len(pattern):
107         return True
108     return self.find(pattern, i, child)
109
110 def mc_creight(self):
111     m = len(self.text)
112     last_head = head = self.root
113     leaf = self.root.connect(Node([0, m - 1]), self.text[0])
114     for i in range(1, m):
115         p = last_head.label
116         q = leaf.label
117         if last_head == self.root:
118             head, j = self.slow_find(q[0] + 1, self.root)
119         else:
120             u = last_head.parent
121             if u == self.root:
122                 v = self.fast_find(p[0] + 1, p[1], self.root)
123                 j = q[0]
124             else:
125                 # print("LINK")
126                 v = self.fast_find(p[0], p[1], u.link)
127                 j = q[0]
128             if len(v.children) == 1:
129                 head = v
130             else:
131                 head, j = self.slow_find(q[0], v)
132             last_head.link = v
133             # print(j, m - 1)
134             leaf = head.connect(Node([j, m - 1]), self.text[j])
135             last_head = head
136
137 def mc_creight_slow(self):
138     m = len(self.text)
139     for i in range(m):
140         head, j = self.slow_find(i, self.root)
141         head.connect(Node([j, m - 1]), self.text[j])
142
143
144 def suffixTreeFastWrapper(text):
145     T = suffixTree(text)
146     T.mc_creight()
147
148
149 def suffixTreeSlowWrapper(text):
150     T = suffixTree(text)
151     T.mc_creight_slow()
152
153
154 def create_treepLOTS(textlist):

```

```

155     for text in textlist:
156         global ind
157         ind = 0
158         T = suffixTree(text)
159         T.mc_creight()
160         T.pretty_print(text, False)
161
162
163 if __name__ == '__main__':
164     create_treepLOTS(["bbbd", "aabbabd", "ababcd", "abcbccd"])
165     # with open("ustawa.txt", "r") as f:
166     #     T = suffixTree(f.read() + "\0")
167     #     T.mc_creight()
168     #     T.pretty_print()
169     #     print(T.find("dochodów2137", 0, T.root))

```

6.3 Testowanie, generowanie wykresów itp.

```

1  from mc_creight_final_attempt import suffixTreeFastWrapper, suffixTreeSlowWrapper
2  from trie import create_trie
3  from time import time
4  import matplotlib.pyplot as plt
5
6
7  def benchmark(funcList,
8               datList,
9               labels,
10              titles,
11              axis_labels=None,
12              repeats=None,
13              save_to_file=False):
14      if repeats is None:
15          repeats = 1
16      for dat, title in zip(datList, titles):
17          times = dict()
18          for func, label in zip(funcList, labels):
19              for _ in range(repeats):
20                  t1 = time()
21                  func(dat)
22                  t2 = time()
23                  times[label] = times.get(label, 0) + (t2 - t1) / repeats
24      plt.bar(times.keys(), times.values())
25      plt.title(title)
26      for i, v in enumerate(times.values()):
27          plt.annotate(str(round(v, 7)), xy=(i, v), ha='center', va='bottom')
28      if axis_labels is not None:
29          plt.xlabel(axis_labels['x'])
30          plt.ylabel(axis_labels['y'])
31
32      if save_to_file:
33          plt.savefig(title + ".png")

```

```

34         plt.show()
35
36
37 if __name__ == '__main__':
38     with open("1997_714_head.txt") as ustawa:
39         funclist = [create_trie, suffixTreeSlowWrapper, suffixTreeFastWrapper]
40         datlist = [
41             "bbbd", "aabbabd", "ababcd", "abcbccd",
42             ustawa.read()[2000] + "\uF000", "abc" * 1000 + "d"
43         ]
44         labels = ["trie", "suffixSlow", "suffixFast"]
45         titles = [s for s in datlist]
46         titles[-2] = "1997_714_head"
47         titles[-1] = "'abc'*1000+'d'"
48         axis_labels = {'x': 'Construction method', 'y': 'time [s]'}
49         benchmark(funclist, datlist, labels, titles, axis_labels, 100, True)

```
