

# Algorytmy tekstowe – Drzewa sufiksów

Tomasz Gargula

## 1 Implementacja algorytmów

W celu wizualizacji drzew użyłem pakietu pydot, stąd w kodzie pojawiają się odniesienia do klas z tego pakietu.

### 1.1 Drzewo trie

Zdecydowałem się na implementację struktur, korzystając głównie z paradygmatu obiektowego. Poniżej kod odpowiedzialny za stworzenie węzła w drzewie Trie.

```
1 class TrieNode:
2     def __init__(self, letter=None, parent=None):
3         self.letter = letter
4         self.parent = parent
5         self.children = {}
6
7     def _save(self, graph, parent=None, prefix=""):
8         word = prefix + self.letter
9         node = pydot.Node(word, label=self.letter, shape="point")
10        graph.add_node(node)
11
12        if parent is not None:
13            graph.add_edge(pydot.Edge(parent, node, label=self.letter))
14
15        for child in self.children.values():
16            child._save(graph, node, word)
```

Oraz kod odpowiedzialny za stworzenie struktury:

```
1 class Trie:
2     def __init__(self, text):
3         self.root = TrieNode()
4         self.text = text
5         self._add(text)
6
7     def _add(self, text):
8         n = len(text)
9         for i in range(n):
10            pointer = self.root
11            for j in range(i, n):
12                letter = text[j]
13                if letter not in pointer.children:
14                    pointer.children[letter] = TrieNode(letter, pointer)
15                pointer = pointer.children[letter]
16
17     def save(self, filename):
18         graph = pydot.Dot(graph_type="digraph", label=f"Trie: {self.text}", labelloc='top')
19
20         root = pydot.Node("root", shape="point")
21         graph.add_node(root)
```

```

22         for child in self.root.children.values():
23             child._save(graph, root)
24
25     graph.write_png(filename)

```

## 1.2 Drzewo sufiksów

Analogicznie do drzew trie rozdzieliłem klasy opisujące całą strukturę drzewa sufiksów oraz jej węzły. W implementacji struktury korzystam z funkcji wyznaczającej długość wspólnego prefiksa dwóch podanych słów równej długości. Kod zamieszczam poniżej.

```

1 def common_suffix_size(A, B):
2     ctr = 0
3     for i, letter in enumerate(A):
4         if letter == B[i]:
5             ctr += 1
6         else:
7             return ctr
8     return ctr

```

Poniżej kod odpowiedzialny za stworzenie węzła w drzewie sufiksów.

```

1 class SuffixTreeNode:
2     def __init__(self, start=None, end=None, parent=None):
3         self.start = start
4         self.end = end
5         self.parent = parent
6         self.children = []
7
8     def length(self):
9         return self.end - self.start + 1
10
11     def slice(self):
12         return (self.start, self.end)
13
14     def subtext(self, text):
15         return text[self.start:self.end+1]
16
17     def _save(self, G, parent, text):
18         node = pydot.Node(str(self), shape="point")
19         G.add_node(node)
20         label = str(self.subtext(text) if self.children or
21                     self.start == self.end else self.slice())
22         edge = pydot.Edge(parent, node, label=label)
23         G.add_edge(edge)
24
25         for child in self.children:
26             child._save(G, node, text)
27
28     def _add(self, text, i):
29         for child in self.children:
30             if text[child.start] == text[i]:
31                 A = text[i:i + child.length()]
32                 B = text[child.start:child.start + len(A)]
33                 if A == B:
34                     child._add(text, i + child.length())
35                 return
36         size = common_suffix_size(A, B)

```

```

37         subnode = SuffixTreeNode(child.start + size, child.end, parent=child)
38         child.end = child.start + size - 1
39         subnode.children = child.children
40         leaf = SuffixTreeNode(i + size, len(text) - 1, parent=child)
41         child.children = [subnode, leaf]
42         return
43     leaf = SuffixTreeNode(start=i, end=len(text) - 1, parent=self)
44     self.children.append(leaf)

```

Oraz kod odpowiedzialny za stworzenie struktury:

```

1  class SuffixTree:
2      def __init__(self, text):
3          self.root = SuffixTreeNode()
4          self.text = text
5          for i in range(len(text)):
6              self.root._add(text, i)
7
8      def save(self, filename):
9          G = pydot.Dot(graph_type="graph", label=f"SuffixTree: {self.text}", labelloc='top')
10
11         root = pydot.Node(str(None), shape="point")
12         G.add_node(root)
13
14         for child in self.root.children:
15             child._save(G, root, self.text)
16
17         G.write_png(filename)

```

## 2 Sterowniki

### 2.1 Benchmarki

Poniżej znajduje się kod, który jest odpowiedzialny za przesetowanie algorytmów pod względem czasu i porównanie ich ze sobą dla zaproponowanych danych do ćwiczenia.

```

1  from pathlib import Path
2  from time import time
3
4  from trees import Trie, SuffixTree
5
6
7  def benchmark(Class, *args):
8      start = time()
9      Class(*args)
10     return time() - start
11
12
13  if __name__ == '__main__':
14      with open(Path('data', 'act.txt'), 'r') as f:
15          x = f.read()
16          x = ''.join(x) + '$'
17
18      data = ('bbbd', 'aabbabd', 'ababcd', 'abcbccd')
19
20      for text in data:
21          print(f'Trie ({text}):\t\t{benchmark(Trie, text)}')
22          print(f'SuffixTree ({text}):\t\t{benchmark(SuffixTree, text)}')

```

```

23
24     print(f'Trie (act.txt):\t\t{benchmark(Trie, x)}')
25     print(f'SuffixTree (act.txt):\t\t{benchmark(SuffixTree, x)}')

```

Poniżej zamieszczam wyniki:

```

1 Trie (bbbd):          1.9073486328125e-05
2 SuffixTree (bbbd):    0.0003688335418701172
3 Trie (aabbabd):       1.9550323486328125e-05
4 SuffixTree (aabbabd): 7.319450378417969e-05
5 Trie (ababcd):        6.723403930664062e-05
6 SuffixTree (ababcd):  1.6450881958007812e-05
7 Trie (abcbccd):       1.7404556274414062e-05
8 SuffixTree (abcbccd): 2.09808349609375e-05
9 Trie (act.txt):       8.21181058883667
10 SuffixTree (act.txt): 0.01808333396911621

```

Można zauważyć, że mimo faktu, że oba algorytmy mają złożoność  $O(n^2)$ , to stworzenie drzewa sufiksów trwa krócej.

## 2.2 Wizualizacja

Napisałem również kod, który tworzy pliki png z zaimplementowanych struktur. Poniżej znajduje się kod.

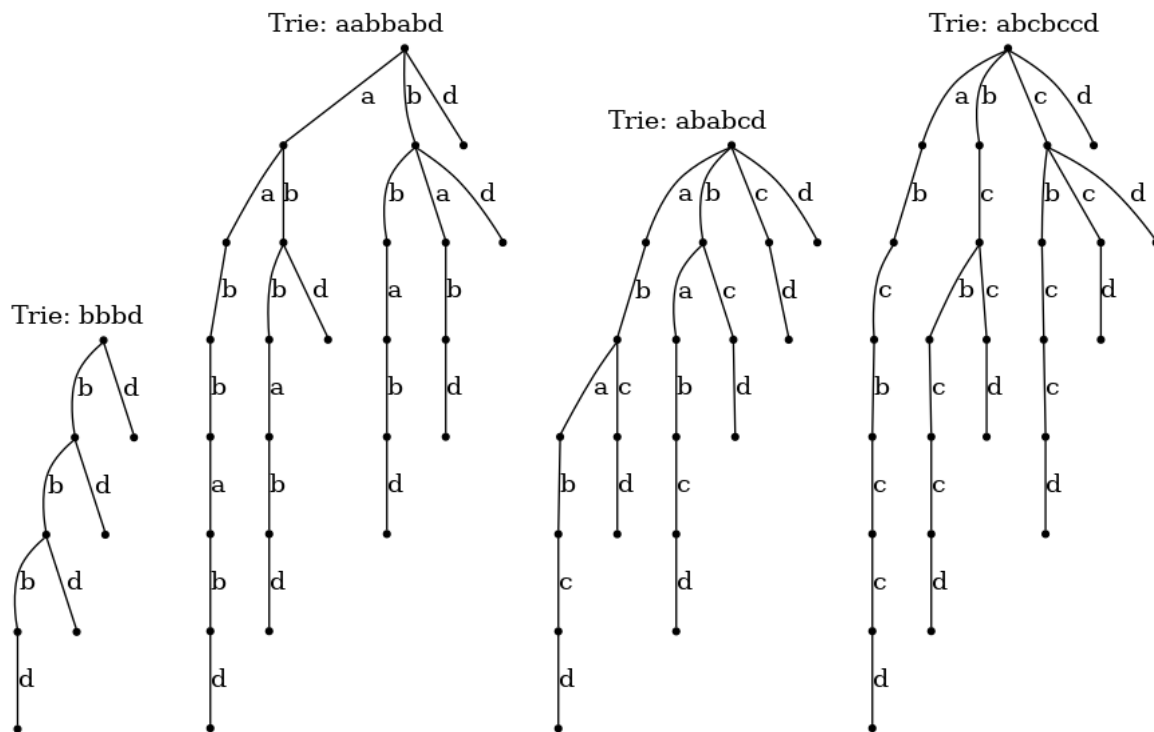
```

1 from pathlib import Path
2
3 from trees import Trie, SuffixTree
4
5 if __name__ == '__main__':
6
7     dir = {
8         Trie: 'trie',
9         SuffixTree: 'suffix_tree'
10    }
11
12    data = ('bbbd', 'aabbabd', 'ababcd', 'abcbccd')
13
14    for text in data:
15        filename = f'{text}.png'
16        T = Trie(text)
17        T.save(Path('out', dir[Trie], filename))
18        ST = SuffixTree(text)
19        ST.save(Path('out', dir[SuffixTree], filename))

```

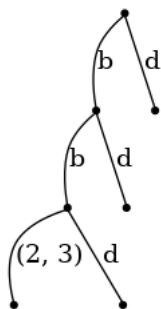
### 3 Wizualizacja

#### 3.1 Drzewa trie

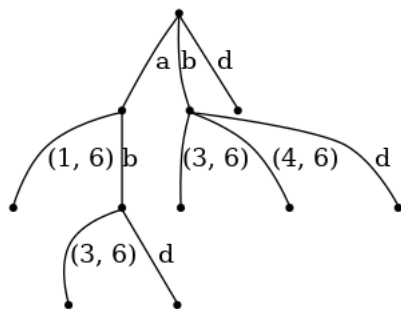


## 3.2 Drzewa sufiksów

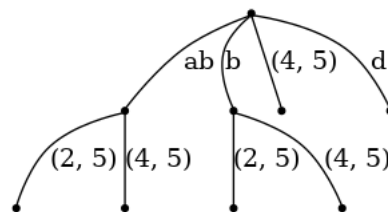
SuffixTree: bbbd



SuffixTree: aabbabd



SuffixTree: ababcd



SuffixTree: abcbccd

