

Algorytmy tekstowe – Algorytmy kompresji

Tomasz Gargula

1 Wstęp

W raporcie opiszę i porównam implementacje statycznego i dynamicznego drzewa Huffmana. Przy implementacji skorzystałem między innymi z bibliotek: `bitarray` oraz `collections.deque`. Ćwiczenie wykonałem, korzystając z paradygmatu obiektowego.

2 Implementacja algorytmów

W następnych podsekcjach opiszę implementację poszczególnych części programu.

2.1 Generator

Aby uzyskać pliki z jednostajnym rozkładem liter, skorzystałem z poniższego programu.

1

2.2 Funkcje pomocnicze

Poniżej zamieszczone są funkcje pomocnicze, które przydały się przy implementacji algorytmów Huffmana. Dodatkowo definiuję zmienne globalne `BITS_PER_LETTER` oraz `BITS_PER_SIZE`, które określają liczbę bitów przeznaczonych na przechowywanie znaków (16 bitów pozwala na kodowanie polskich znaków diakrytycznych) oraz na przechowywanie rozmiaru drzewa. Określenie rozmiaru struktury jest konieczne, ponieważ translacja do kodu binarnego powoduje dopisanie bitów do pełnych bajtów. Dopisywane są wtedy zera do zakodowanej struktury. Może więc zdarzyć się sytuacja, że po odkodowaniu otrzymamy jedną lub dwie litery więcej.

```
1 BITS_PER_LETTER = 16
2 BITS_PER_SIZE = 64
3
4 def extract_min(leafs, nodes):
5     if not leafs:
6         return nodes.pop()
7     if not nodes:
8         return leafs.pop()
9     return nodes.pop() if nodes[-1].weight < leafs[-1].weight else leafs.pop()
10
11
12 def count_letters(text):
13     dictionary = {}
14     for letter in text:
15         dictionary[letter] = dictionary.get(letter, 0) + 1
16     return dictionary
17
18
19 def utf8(letter):
20     return bin(ord(letter))[2:].zfill(BITS_PER_LETTER)
```

2.3 Węzły

Poniżej zamieszczam implementację węzłów wykorzystywanych w strukturach drzew Huffmana. Warto zauważyć, że określenie dzieci przy tworzeniu węzła, rejestruje węzeł automatycznie jako rodzica. Metoda `_save`, podobnie jak w innych strukturach, służy do wizualizacji drzewa przy pomocy biblioteki `pydot`.

```
1 class Node:
2     def __init__(self, weight, left=None, right=None):
3         self.left = left
4         self.right = right
5         self.weight = weight
6         self.parent = None
7         self.visited = False
8         if left: left.parent = self
9         if right: right.parent = self
10
11     def postorder(self):
12         return self.left.postorder() + self.right.postorder() + [self]
13
14     def _save(self, G, parent=None, i=None):
15         node = pydot.Node(str(self), label=str(self.weight), shape="circle")
16
17         G.add_node(node)
18         if parent:
19             edge = pydot.Edge(parent, node, label=str(i))
20             G.add_edge(edge)
21         for i, child in enumerate((self.left, self.right)):
22             child._save(G, node, i)
```

2.4 Liście

Zdecydowałem się wyodrębnić osobną klasę reprezentującą liście, która dodatkowo posiada metodę `code()`, która zwraca aktualną reprezentację bitową litery znajdującej się w liście.

```
1 class Leaf(Node):
2     def __init__(self, weight, letter):
3         super().__init__(weight)
4         self.letter = letter
5
6     def postorder(self):
7         return [self]
8
9     def code(self):
10        pointer = self
11        code = []
12        while pointer.parent is not None:
13            parent = pointer.parent
14            if parent.left == pointer:
15                code.append("0")
16            if parent.right == pointer:
17                code.append("1")
18            pointer = parent
19        return ''.join(reversed(code))
20
21     def _save(self, G, parent, i):
22        node = pydot.Node(
23            str(self), label=f'{self.letter}: {self.weight}', shape="square")
24        G.add_node(node)
```

```

25     edge = pydot.Edge(parent, node, label=str(i))
26     G.add_edge(edge)

```

2.5 Lista

Dynamiczne drzewo Huffmana wykorzystuje zależność, że jeśli T jest kompletnym ważonym drzewem binarnym z p liśćmi z dodatnimi wagami, a waga rodzica jest równa sumie wag dzieci, to węzły mogą być ułożone w ciąg niemalejący $(w(x_1), w(x_2), w(x_3), \dots, w(x_n))$ taki, że dla każdego i x_{2i} oraz x_{2i+1} są bliźniakami tj. mają tego samego rodzica. Gdy przetwarzamy nową literę, to zwiększamy jej wagę. Aby utrzymać poprawność struktury potrzebny jest szybki dostęp do następnego węzła. Aby to zapewnić, skorzystałem ze wzorca **Decorator**, obudowując zwykłą pythonową listę i dodając do każdego węzła informację o tym, na którym miejscu w liście się znajduje.

```

1  '''Decorator class for List that saves info about the index'''
2  class List:
3      def __init__(self, *nodes):
4          self.structure = []
5          self.extend(*nodes)
6
7      def __getitem__(self, index):
8          return self.structure[index]
9
10     def __setitem__(self, index, node):
11         self.structure[index] = node
12         node.index = index
13
14     def append(self, node):
15         self.structure.append(node)
16         node.index = len(self.structure) - 1
17
18     def extend(self, *nodes):
19         for node in nodes:
20             self.append(node)
21
22     def pop(self):
23         node = self.structure.pop()
24         node.index = None

```

2.6 Prototyp drzewa Huffmana

Statyczne i dynamiczne drzewo Huffmana mają kilka metod wspólnych lub działających na tej samej zasadzie. Korzystając więc z mechanizmu dziedziczenia, zdefiniowałem klasę `HuffmanTree`.

```

1  class HuffmanTree:
2      def __init__(self, root):
3          self.root = root
4          self.codes, self.letters = self._get_encodings()
5
6      def save(self, filename='out.png'):
7          G = pydot.Dot(graph_type="graph")
8          self.root._save(G)
9          G.write_png(filename)
10
11     def _get_encodings(self):
12         codes, letters = {}, {}
13         pointer = self.root
14         code = []
15         while pointer:
16             if pointer.left is None:

```

```

17         encoded = ''.join(code)
18         codes[pointer.letter] = encoded
19         letters[encoded] = pointer.letter
20         pointer.visited = True
21         code.pop()
22         pointer = pointer.parent
23     else:
24         if not pointer.left.visited:
25             code.append('0')
26             pointer = pointer.left
27         elif not pointer.right.visited:
28             code.append('1')
29             pointer = pointer.right
30         else:
31             pointer.visited = True
32             code = code[:-1]
33             pointer = pointer.parent
34     return codes, letters
35
36 def decode_size(self, binary):
37     return int(binary[:BITS_PER_SIZE].to01(), 2)

```

2.7 Statyczne drzewo Huffmana

Poniżej zamieściłem implementację statycznego drzewa Huffmana.

```

1 class StaticHuffmanTree(HuffmanTree):
2     def __init__(self, text=None):
3         if text is not None:
4             super().__init__(self._build(count_letters(text)))
5             self.encoded = self.encode(text).tobytes()
6         else:
7             self.root = None
8
9     def get_encoded(self):
10        return self.encoded
11
12    def _build(self, dictionary):
13        leaves = [Leaf(weight, letter) for letter, weight in
14                  sorted([item for item in dictionary.items()], key=lambda item: -item[1])]
15        nodes = deque()
16        if len(leaves) == 1:
17            nodes.append(leaves.pop())
18        while len(leaves) + len(nodes) > 1:
19            node1 = extract_min(leaves, nodes)
20            node2 = extract_min(leaves, nodes)
21            nodes.appendleft(Node(node1.weight + node2.weight, node1, node2))
22        return nodes.pop()
23
24    def _encode_tree(self):
25        return bytearray(''.join(
26            ['1' + utf8(node.letter) if isinstance(node, Leaf) else '0' for node in
27             ↪ self.root.postorder()]
28            ) + '0')
29
30    def _encode_size(self, text):
31        return bytearray(bin(len(text))[2:].zfill(BITS_PER_SIZE))

```

```

32     def _encode_text(self, text):
33         return bytearray(''.join([self.codes[letter] for letter in text]))
34
35     def encode(self, text):
36         return self._encode_size(text) + self._encode_tree() + self._encode_text(text)
37
38     def decode_tree(self, binary):
39         stack = []
40         j = 0
41         while j < len(binary):
42             j += 1
43             if binary[j-1]:
44                 stack.append(Leaf(None, chr(int(binary[j:j+BITS_PER_LETTER].to01(), 2))))
45                 j += BITS_PER_LETTER
46             else:
47                 if len(stack) == 1:
48                     self.root = stack.pop()
49                     self.codes, self.letters = self._get_encodings()
50                     return j
51                 node1, node2 = stack.pop(), stack.pop()
52                 stack.append(Node(None, node2, node1))
53
54     def decode_text(self, binary, size):
55         i, j = 0, 1
56         result = []
57         while i < len(binary) and (size is None or len(result) < size):
58             while binary[i:j].to01() not in self.letters:
59                 j += 1
60             result += self.letters[binary[i:j].to01()]
61             i = j
62         return ''.join(result)
63
64     def decode(self, binary):
65         size = self.decode_size(binary)
66         text_begin = self.decode_tree(binary[BITS_PER_SIZE:])
67         return self.decode_text(binary[BITS_PER_SIZE + text_begin:], size)

```

2.8 Dynamiczne drzewo Huffmana

Poniżej znajduje się implementacja dynamicznego drzewa Huffmana.

```

1  class DynamicHuffmanTree(HuffmanTree):
2      def __init__(self, stream=[]):
3          self.root = Leaf(0, None)
4          self.NYT = self.root
5          self.pointers = {}
6          self.encoded = bytearray()
7          self.nodes = List(self.NYT)
8          for letter in stream:
9              self.add(letter)
10
11     def get_encoded(self):
12         return bytearray(bin(self.root.weight)[2:].zfill(BITS_PER_SIZE)) + self.encoded
13
14     def decode(self, binary):
15         pointer = self.root
16         decoded = []
17         skipping = 0

```

```

18
19     size = decode_size(binary)
20     print(size)
21
22     binary = binary[BITS_PER_SIZE:]
23
24     for i, bit in enumerate(binary):
25         skipping = skipping - 1 if skipping > 0 else 0
26         if not size:
27             break
28         if not skipping:
29             if isinstance(pointer, Leaf):
30                 size -= 1
31                 if pointer is self.NYT:
32                     letter = chr(int(binary[i:i+BITS_PER_LETTER].to01(), 2))
33                     print(binary[i:i+BITS_PER_LETTER])
34                     skipping = BITS_PER_LETTER
35                     self.add(letter)
36                     decoded.append(letter)
37                     pointer = self.root
38                     continue
39                 else:
40                     decoded.append(pointer.letter)
41                     print(pointer.letter)
42                     self.add(pointer.letter)
43                     pointer = self.root
44
45             pointer = pointer.right if bit else pointer.left
46
47     skipping = skipping - 1 if skipping > 0 else 0
48     if not skipping:
49         if isinstance(pointer, Leaf):
50             if pointer is self.NYT:
51                 letter = chr(int(binary[i:i+BITS_PER_LETTER].to01(), 2))
52                 skipping = BITS_PER_LETTER
53                 self.add(letter)
54                 decoded.append(letter)
55                 pointer = self.root
56             else:
57                 decoded.append(pointer.letter)
58                 self.add(pointer.letter)
59                 pointer = self.root
60
61     return ''.join(decoded)
62
63 def next(self, node):
64     return self.nodes[node.index - 1] if node.index > 0 else None
65
66 def add(self, letter):
67     pointer = self.root
68     if letter in self.pointers: # has been transferred
69         leaf = self.pointers[letter]
70         self.encoded.extend(leaf.code())
71         p = leaf
72         next = self.next(p)
73         while next is not None and not (isinstance(p, Leaf) ^ isinstance(next, Leaf)) and
74             p.weight == next.weight:

```

```

74         self._swap(p, next)
75         next = self.next(p)
76     if p.parent.left is self.NYT:
77         leaf = p
78         p = p.parent
79     else:
80         leaf = None
81
82     else: # not yet transferred (NYT)
83         self.encoded.extend(self.NYT.code() + utf8(letter))
84         parent = self.NYT.parent
85         self.nodes.pop()
86         leaf = Leaf(0, letter)
87         node = Node(0, self.NYT, leaf)
88         self.nodes.extend(node, leaf, self.NYT)
89         if self.NYT is self.root:
90             self.root = node
91         if parent is not None:
92             parent.left = node
93         node.parent = parent
94         self.pointers[letter] = leaf
95         p = node
96
97     self._update(p, leaf)
98
99     def _swap(self, node1, node2):
100         index1 = node1.index
101         index2 = node2.index
102         self.nodes[index1], self.nodes[index2] = self.nodes[index2], self.nodes[index1]
103         parent1 = node1.parent
104         parent2 = node2.parent
105         if parent1.left is node1 and parent2.left is node2:
106             parent1.left, parent2.left = parent2.left, parent1.left
107         elif parent1.left is node1:
108             parent1.left, parent2.right = parent2.right, parent1.left
109         elif parent2.left is node2:
110             parent1.right, parent2.left = parent2.left, parent1.right
111         else:
112             parent1.right, parent2.right = parent2.right, parent1.right
113         node1.parent = parent2
114         node2.parent = parent1
115
116     def _update(self, node, leaf):
117         parent = node.parent
118         while node is not None:
119
120             next = self.next(node)
121             while next is not None and isinstance(next, Leaf) and next.weight == node.weight +
122                 ↪ 1:
123                 self._swap(node, next)
124                 next = self.next(node)
125             node.weight += 1
126             node = parent
127             if node is not None:
128                 parent = node.parent
129
130         if leaf is not None:

```

```

130         next = self.next(leaf)
131         while next is not None and not isinstance(next, Leaf) and next.weight ==
            ↳ leaf.weight:
132             self._swap(leaf, next)
133             next = self.next(leaf)
134         leaf.weight += 1

```

2.9 Testy

W celu sprawdzenia, czy algorytmy są w stanie odkodować zakodowany tekst napisałem test, który zamieszczam poniżej. W tym celu skorzystałem z możliwości biblioteki `pytest`

```

1  import pytest
2  from os import listdir
3  from bitarray import bitarray
4
5  from huffman.huffman import StaticHuffmanTree, DynamicHuffmanTree
6
7
8  def test_decompress():
9      for filename in listdir("data"):
10         encodedfile = ''.join(filename.split('.')[:-1] + [".bin"])
11         with open(f"data/{filename}", "r") as f:
12             text = f.read()
13
14         for HuffmanTree in (StaticHuffmanTree, DynamicHuffmanTree,):
15
16             with open(f"encoded/{HuffmanTree.__name__}/{encodedfile}", "rb") as f:
17                 array = bitarray()
18                 array.fromfile(f)
19                 decoded = HuffmanTree().decode(array)
20
21             assert text == decoded

```

2.10 Benchmarki

Poniżej zamieszczam kod odpowiedzialny za wykonanie pomiarów i sporządzenie wykresów.

```

1  from os import listdir
2  from timeit import default_timer as timer
3  import matplotlib
4  from matplotlib import pyplot as plt
5  import numpy as np
6
7  from huffman.huffman import DynamicHuffmanTree, StaticHuffmanTree
8
9
10 def benchmark(f, *args):
11     start = timer()
12     result = f(*args)
13     return timer() - start, result
14
15
16 if __name__ == '__main__':
17
18     data = ([], [], [])
19
20     for filename in sorted(listdir("data")):

```



```

21     data[0].append(filename)
22
23
24     for i, HuffmanTree in enumerate((StaticHuffmanTree, DynamicHuffmanTree)):
25         file = filename
26
27         with open(f"data/{file}", "r") as f:
28             elapsed, HT = benchmark(HuffmanTree, f.read())
29
30
31     data[i+1].append(elapsed)
32
33     file = ''.join(file.split('.')[:-1] + [".bin"])
34
35     # since we have built the tree, we can save it to file
36     with open(f"encoded/{HuffmanTree.__name__}/{file}", "wb") as f:
37         f.write(HT.get_encoded())
38
39
40 width = 0.35
41
42 independent = data[0][:3], data[1][:3], data[2][:3]
43 dependent = data[0][3:], data[1][3:], data[2][3:]
44
45 x = np.arange(len(independent[0]))
46
47 fig, ax = plt.subplots()
48 rects1 = ax.bar(x - width / 2, independent[1], width, label='static')
49 rects2 = ax.bar(x + width / 2, independent[2], width, label='dynamic')
50 ax.set_title('Static and dynamic huffman encoding comparison')
51 ax.set_ylabel('Time [s]')
52 ax.set_xticks(x)
53 ax.set_yscale('log')
54 ax.set_xticklabels(independent[0], rotation="vertical", fontsize="x-small")
55 ax.legend()
56
57 for rect, label in zip(rects1, independent[1]):
58     height = rect.get_height()
59     ax.text(rect.get_x() + rect.get_width() / 2,
60             height, round(label,3), ha='center', va='bottom', fontsize="x-small")
61
62 for rect, label in zip(rects2, independent[2]):
63     height = rect.get_height()
64     ax.text(rect.get_x() + rect.get_width() / 2,
65             height, round(label,3), ha='center', va='bottom', fontsize="x-small")
66
67 fig.tight_layout()
68
69 plt.show()
70
71 x = np.arange(len(dependent[0]))
72
73 fig, ax = plt.subplots()
74 rects1 = ax.bar(x - width / 2, dependent[1], width, label='static')
75 rects2 = ax.bar(x + width / 2, dependent[2], width, label='dynamic')
76 ax.set_title('Static and dynamic huffman encoding comparison')
77 ax.set_ylabel('Time [s]')

```

```

78 ax.set_xticks(x)
79 ax.set_yscale('log')
80 ax.set_xticklabels(dependent[0], rotation="vertical", fontsize="x-small")
81 ax.legend()
82
83 for rect, label in zip(rects1, dependent[1]):
84     height = rect.get_height()
85     ax.text(rect.get_x() + rect.get_width() / 2,
86             height, round(label,3), ha='center', va='bottom', fontsize="x-small")
87
88 for rect, label in zip(rects2, dependent[2]):
89     height = rect.get_height()
90     ax.text(rect.get_x() + rect.get_width() / 2,
91             height, round(label,3), ha='center', va='bottom', fontsize="x-small")
92
93 fig.tight_layout()
94
95 plt.show()

```

2.11 Wylczenie współczynników kompresji

Poniżej zamieszczam kod odpowiedzialny za wylczenie współczynników kompresji i sporządzenie wykresów.

```

1  import os
2  import matplotlib
3  from matplotlib import pyplot as plt
4  import numpy as np
5
6  if __name__ == '__main__':
7
8      data = ([], [], [])
9
10     for filename in sorted(os.listdir("data")):
11
12         original = os.path.getsize(f"data/{filename}")
13
14         encoded = ''.join(filename.split('.')[:-1] + [".bin"])
15
16         static = os.path.getsize(f"encoded/StaticHuffmanTree/{encoded}")
17
18         dynamic = os.path.getsize(f"encoded/DynamicHuffmanTree/{encoded}")
19
20         data[0].append(filename)
21         data[1].append(static / original * 100)
22         data[2].append(dynamic / original * 100)
23
24     x = np.arange(len(data[0])) # the label locations
25     width = 0.35 # the width of the bars
26
27     fig, ax = plt.subplots()
28     rects1 = ax.bar(x - width / 2, data[1], width, label='static')
29     rects2 = ax.bar(x + width / 2, data[2], width, label='dynamic')
30     ax.set_title('Static and dynamic huffman encoding comparison')
31     ax.set_ylabel('Compression ratio')
32     ax.set_xticks(x)
33     ax.set_xticklabels(data[0], rotation="vertical", fontsize="x-small")
34     ax.legend()
35

```

```

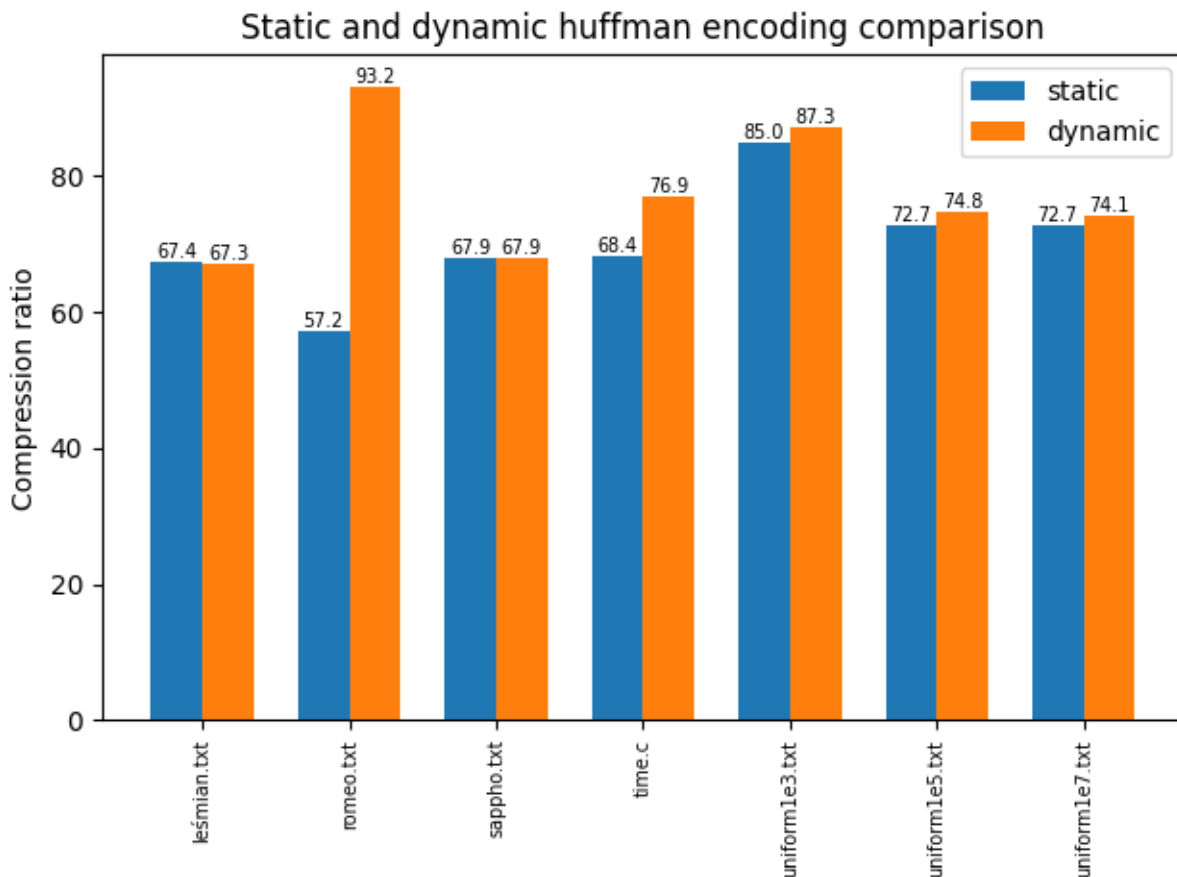
36 for rect, label in zip(rects1, data[1]):
37     height = rect.get_height()
38     ax.text(rect.get_x() + rect.get_width() / 2,
39             height, round(label,1), ha='center', va='bottom', fontsize="x-small")
40
41 for rect, label in zip(rects2, data[2]):
42     height = rect.get_height()
43     ax.text(rect.get_x() + rect.get_width() / 2,
44             height, round(label,1), ha='center', va='bottom', fontsize="x-small")
45
46 fig.tight_layout()
47
48 plt.show()

```

3 Porównanie współczynnika kompresji

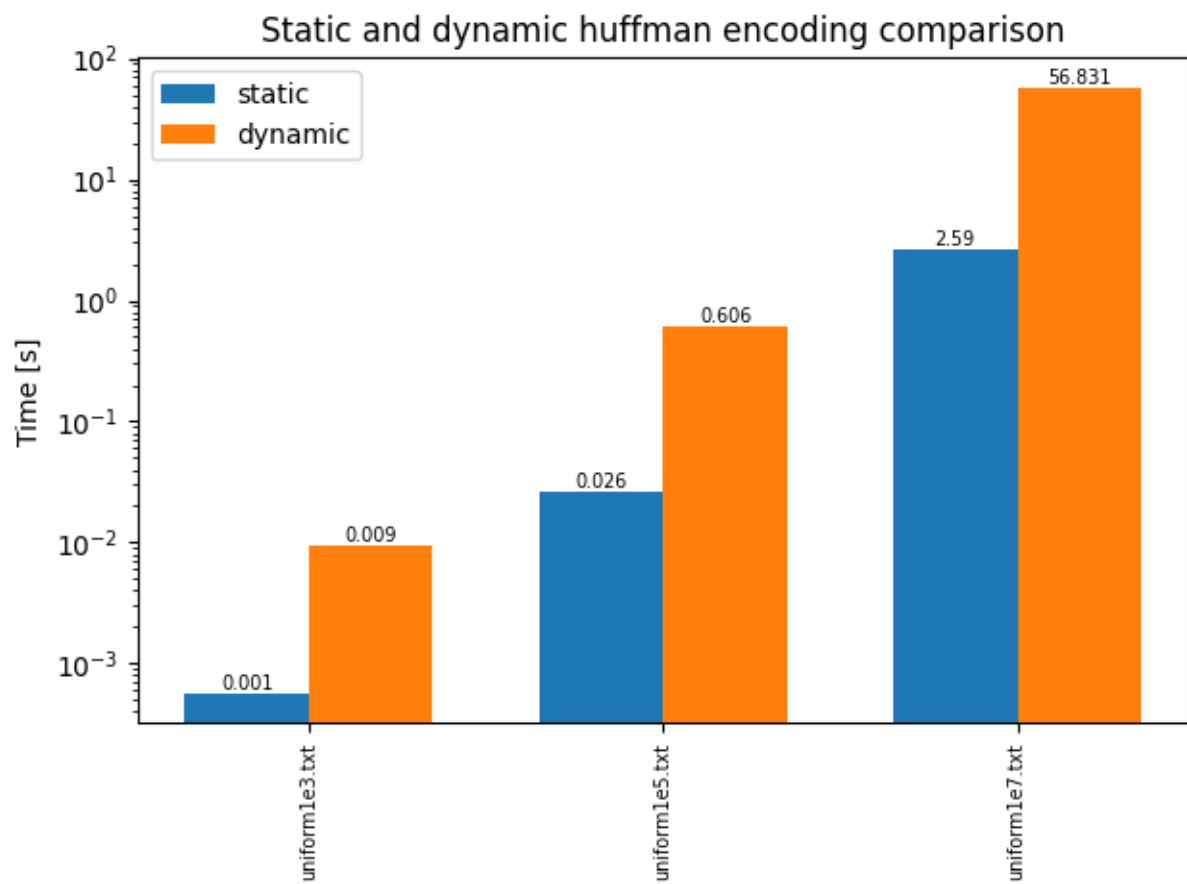
Poniżej zamieszczam porównanie współczynników kkompresji dla obydwu algorytmów wyliczanych jako

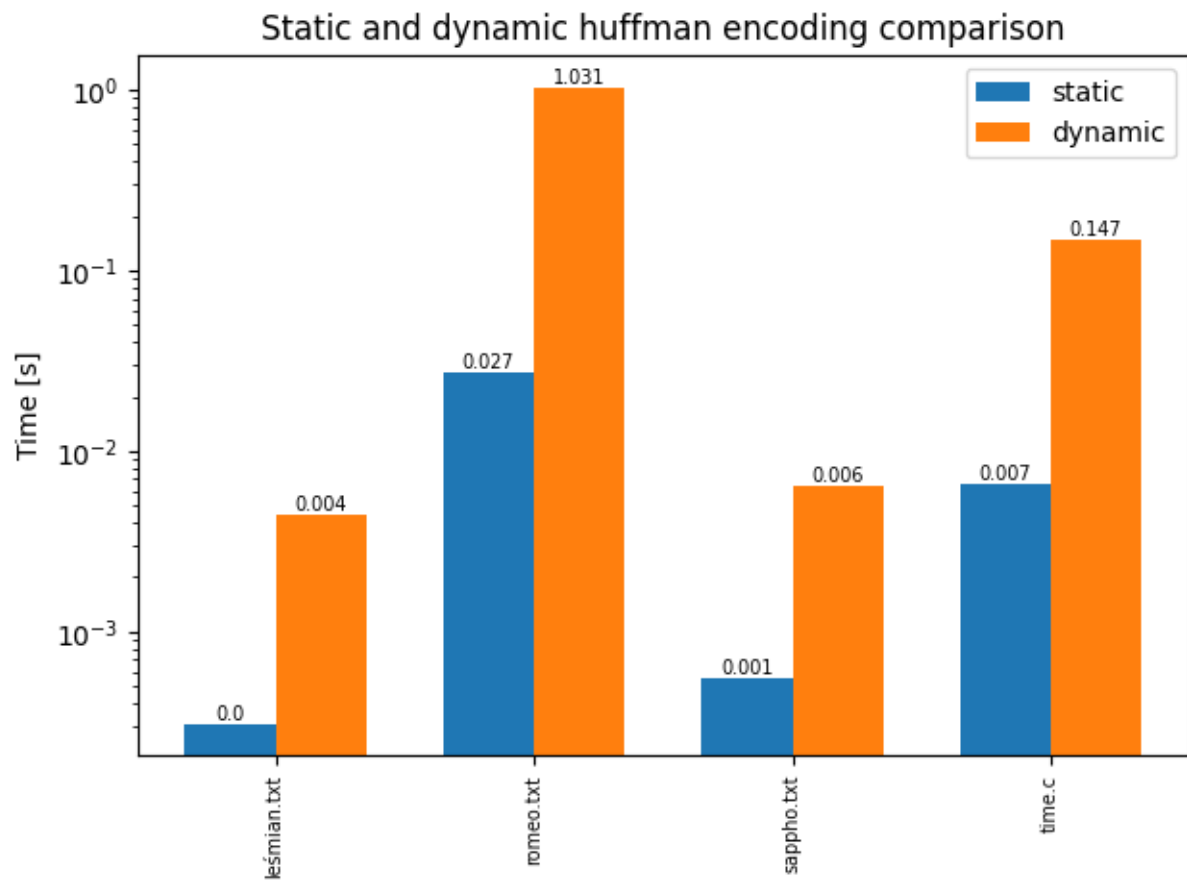
$$\frac{\text{rozmiar pliku skompresowanego}}{\text{rozmiar pliku oryginalnego}} \cdot 100\%$$



4 Porównanie czasów budowy struktur (kompresji)

Poniżej zamieszczam wykresy przedstawiające czasy wykonania algorytmów. Na drugim wykresie dane zostały wygenerowane przez zamieszczony wyżej generator dla kolejno: 1e3, 1e5 oraz 1e7 znaków.

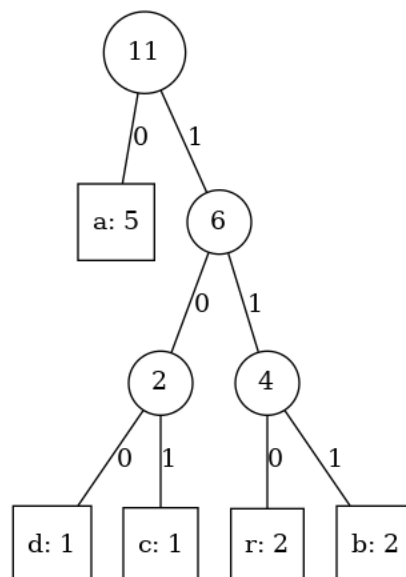




5 Wizualizacje

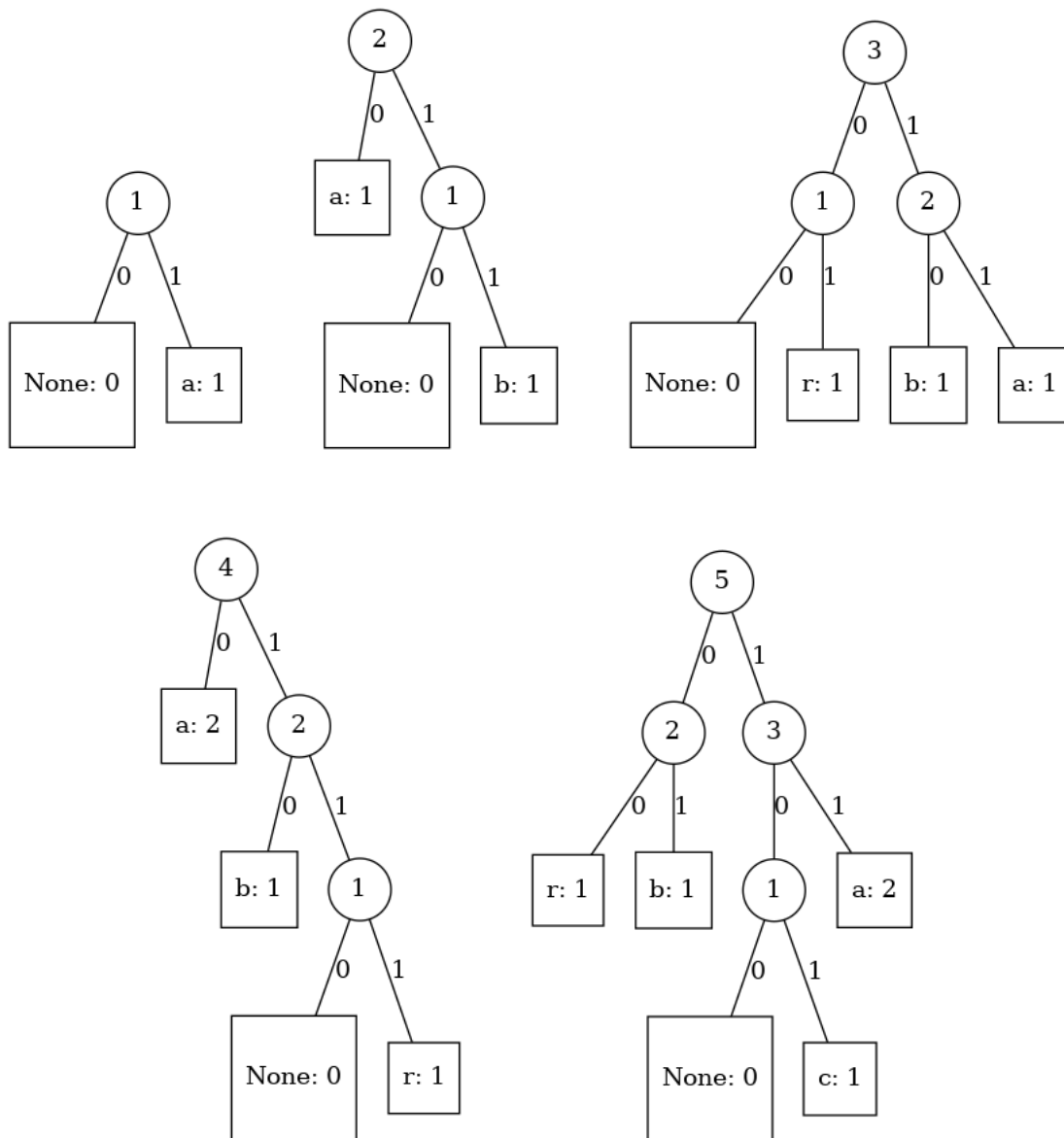
Wykonałem test dla słowa "abracadabra" i zapisałem struktury w postaci zwizualizowanej.

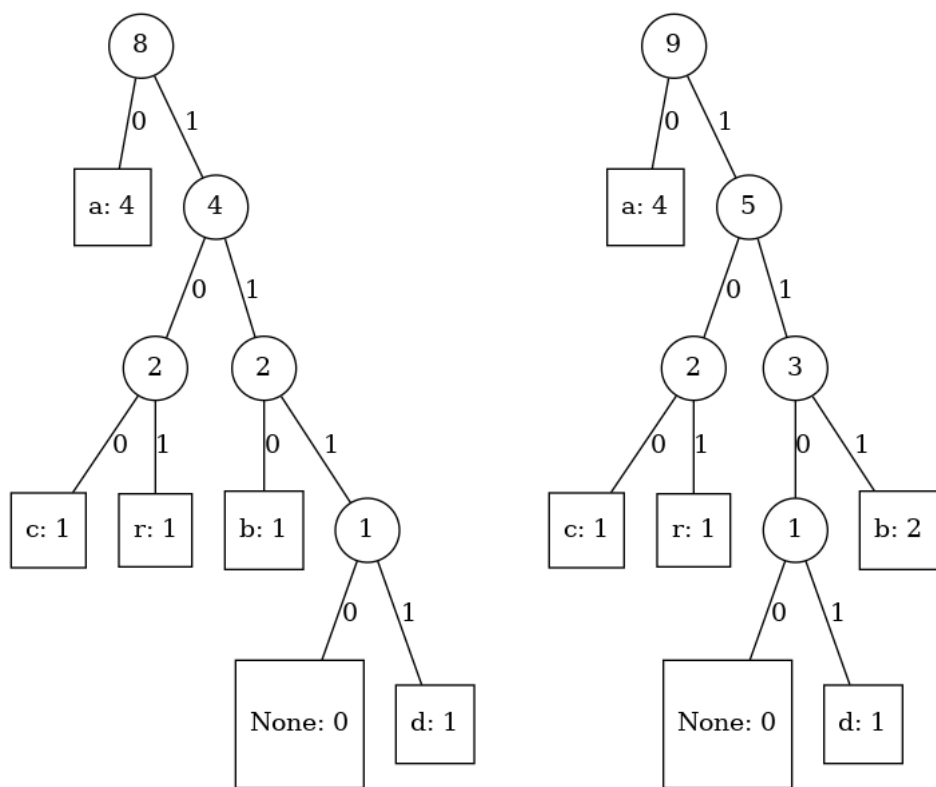
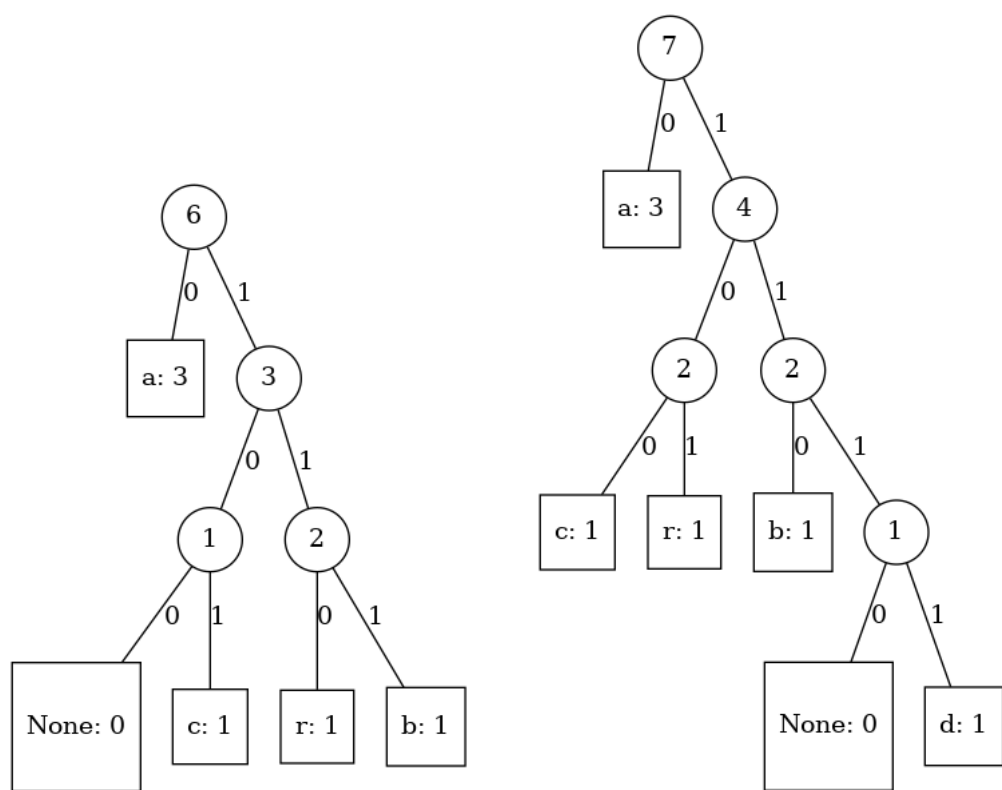
5.1 Statyczne drzewo Huffmana

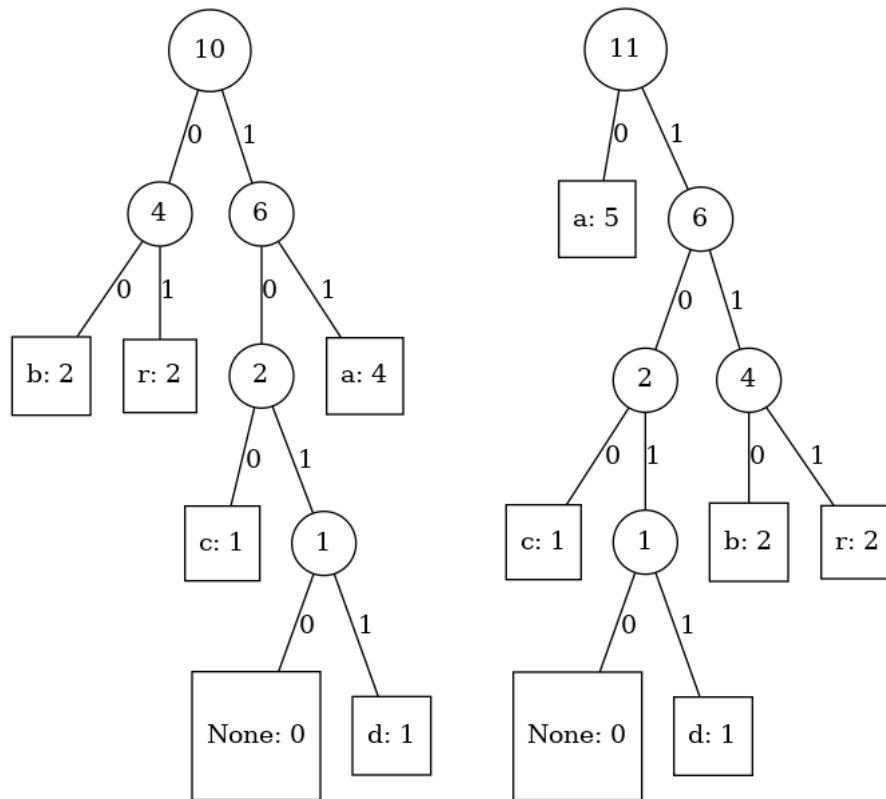


5.2 Dynamiczne drzewo Huffmana

Na kolejnych grafach będzie widać, w jaki sposób powstaje drzewo Huffmana po dodaniu każdej kolejnej litery.







Zauważmy, że w każdym momencie zachowany jest warunek z lematu opisanego przy implementacji klasy List.

6 Wnioski

Obydwa algorytmy dobrze spełniają swoją rolę tj. współczynnik kompresji znajduje się w granicach 70 – 75% w zależności od danych. Kompresja za pomocą dynamicznego drzewa Huffmana trwa dłużej, jednak jest bardziej elastyczna – dodanie kolejnej litery jest możliwe bez określania od nowa kodowania dla wszystkich liter.