

Algorytmy tekstowe – Dopasowanie wzorców

Tomasz Gargula

1 Implementacja algorytmów

1.1 Algorytm naiwny

```
1 def naive(data, pattern):
2     n = len(data)
3     m = len(pattern)
4     result = []
5
6     for s in range(n - m + 1):
7         if pattern == data[s:s+m]:
8             result.append(s)
9
10    return result
```

1.2 Automat skończony

Do działania automatu skończonego potrzebne jest wygenerowanie tablicy przejścia.

```
1 def transition_table(pattern):
2     Q = range(len(pattern) + 1)
3     Sigma = set(pattern)
4     delta = []
5
6     for q in Q:
7         delta.append({})
8         for letter in Sigma:
9             k = min(len(Q), q + 2) - 1
10            text = pattern[:q] + letter
11            while not text.endswith(pattern[:k]):
12                k -= 1
13
14            delta[q][letter] = k
15
16    return delta
```

Zaimplementowano dwie wersje algorytmów (z wykonanym preprocessingiem i bez), aby ułatwić korzystanie z funkcji na późniejszym etapie.

```
1 def fa(data, pattern):
2     delta = transition_table(pattern)
3     return fa_preprocessed(data, delta)
4
5
6 def fa_preprocessed(data, delta):
7     q = 0
8     result = []
9
10    for s, letter in enumerate(data):
```

```

11         q = delta[q].get(letter, 0)
12         if q + 1 == len(delta):
13             result.append(s - q + 1)
14
15     return result

```

1.3 Algorytm Knutha-Morrisa-Pratta

Poniżej zamieszczono kod z funkcją, która wyznacza funkcję przejścia niezbędną do prawidłowego działania algorytmu Knutha-Morrisa-Pratta.

```

1 def prefix_function(pattern):
2     m = len(pattern)
3     pi = [0] * m
4     k = 0
5
6     for q in range(1, m):
7         while k > 0 and pattern[k] != pattern[q]: # check back
8             k = pi[k - 1]
9         if pattern[k] == pattern[q]: # letters match
10             k += 1
11         pi[q] = k
12
13     return pi

```

Analogicznie jak w przypadku automatu skończonego, zaimplementowano dwie wersje algorytmu.

```

1 def kmp(data, pattern, pi=None):
2     pi = prefix_function(pattern)
3     return kmp_preprocessed(data, pattern, pi)
4
5
6 def kmp_preprocessed(data, pattern, pi):
7     n = len(data)
8     m = len(pattern)
9     q = 0 # number of matching symbols
10    result = []
11    for i, letter in enumerate(data):
12        while q > 0 and pattern[q] != letter:
13            q = pi[q - 1]
14        if pattern[q] == letter:
15            q += 1
16        if q == m:
17            result.append(i - m + 1)
18            q = pi[q - 1]
19
20    return result

```

2 Testy porównujące szybkość działania

Czas działania algorytmu był mierzony poprzez różnicę czasów po i przed wykonaniem funkcji. W tym celu użyłem funkcji `time` z biblioteki `time`.

```

1 from time import time
2
3 def test(f, *args):
4     start = time()
5     f(*args)
6     return time() - start

```

3 Wyszukiwanie wystąpień wzorca

Aby zademonstrować działanie algorytmów napisany został algorytm, znajdujący wszystkie wystąpienia wzorca "art" w *załączonej ustawie*. Ten dokument został umieszczony w wersji tekstowej w folderze data.

```
1 from pathlib import Path
2
3 from algorithms.pattern_matching import fa, kmp, naive
4
5 if __name__ == '__main__':
6     text = Path('data', 'act.txt').open('r').read()
7     pattern = 'art'
8
9     for f in (naive, fa, kmp):
10         print(f'{f.__name__}:')
11         print(f(text, pattern))
12         print()
```

4 Porównanie szybkości działania

4.1 Wyszukiwanie wzorca "art" w ustawie

```
1 from pathlib import Path
2 from pprint import pprint
3
4 from algorithms.pattern_matching import fa, kmp, naive
5 from algorithms.time_test import test
6
7 if __name__ == '__main__':
8     text = Path('data', 'act.txt').open('r').read()
9     pattern = 'art'
10    result = {}
11
12    for f in (naive, fa, kmp):
13        result[f.__name__] = test(f, text, pattern)
14
15    pprint(result)
```

Wynik działania algorytmu:

```
{'fa': 0.03156280517578125, 'kmp': 0.018553495407104492, 'naive': 0.023760318756103516}
```

4.2 Propozycja danych, dla których czas działania algorytmu naiwnego jest co najmniej pięciokrotnie dłuższy od czasu działania automatu skończonego i algorytmu Knutha-Morrisa-Pratta

Zaproponowany tekst to: ('a' * 10000 + 'b' + 'a' * 100 + 'c') * 100

Zaproponowany wzorzec to: 'a' * 7000 + 'b'

```
1 from pprint import pprint
2
3 from algorithms.pattern_matching import fa_preprocessed as fa
4 from algorithms.pattern_matching import kmp_preprocessed as kmp
5 from algorithms.pattern_matching import naive, prefix_function, transition_table
6 from algorithms.time_test import test
7
8 if __name__ == '__main__':
9     text = ('a' * 10000 + 'b' + 'a' * 100 + 'c') * 100
10    pattern = 'a' * 7000 + 'b'
```

```

11     result = {}
12
13     delta = transition_table(pattern)
14     pi = prefix_function(pattern)
15
16     result['naive'] = test(naive, text, pattern)
17     result['fa'] = test(fa, text, delta)
18     result['kmp'] = test(kmp, text, pattern, pi)
19
20     print(result['naive'] / result['fa'] > 5 and result['naive'] / result['kmp'] > 5)
21
22     pprint(result)

```

Wynik działania algorytmu:

```
True {'fa': 0.1463453769683838, 'kmp': 0.13521957397460938, 'naive': 0.7533395290374756}
```

4.3 Propozycja wzorca, dla którego czas obliczenia tablicy przejścia automatu skończonego będzie co najmniej pięciokrotnie dłuższy niż czas potrzebny na utworzenie funkcji przejścia w algorytmie Knutha-Morrisa-Pratta

Zaproponowany wzorec to: 'Wstawaj samuraju!'

```

1  from pprint import pprint
2
3  from algorithms.pattern_matching import prefix_function, transition_table
4  from algorithms.time_test import test
5
6  if __name__ == '__main__':
7
8      pattern = 'Wstawaj samuraju!'
9
10     result = {}
11
12     result['transition_table'] = test(transition_table, pattern)
13     result['prefix_function'] = test(prefix_function, pattern)
14
15     print(result['transition_table'] / result['prefix_function'] > 5)
16     pprint(result)

```

Wynik działania algorytmu:

```
True {'prefix_function': 3.0994415283203125e-06, 'transition_table': 0.0003154277801513672}
```