

Санкт-Петербургский Национальный Исследовательский
Университет Информационных Технологий, Механики и Оптики

Факультет инфокоммуникационных технологий

Лабораторная работа №7 "Алгоритмы поиска подстрок"
Вариант 7

Выполнили:
Голованов Д.И.,
Шарыпов Е.А.,
Стафеев И.А.
Проверил
Мусаев А.А.

Санкт-Петербург,
2024

СОДЕРЖАНИЕ

Стр.

ВВЕДЕНИЕ	3
1 Задача 1	4
2 Задача 2	10
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ.....	13

ВВЕДЕНИЕ

В данной лабораторной работе необходимо решить следующие задачи:

1. Заполнить массив 500 числами Фибоначчи написанными слитно. Используя алгоритмы поиска подстрок (наивный, Рабина-Карпа, Бойера-Мура, Кнута-Морриса-Пратта), посчитать количество наиболее часто встречающихся двузначных чисел в образовавшейся строке. Сравнить изученные алгоритмы поиска подстрок.
2. Дан реферат на тему «Логика». Определить количество плагиата (в % от общего количества символов в реферате) в тексте реферата, взяв за основу статью из Википедии «Логика». За плагиат считать любые 3 совпавших слова, идущих подряд. Обосновать выбранный алгоритм поиска.

1 Задача 1

Ниже приведены скриншоты кода для теста алгоритмов 1.1, алгоритма наивного поиска 1.2, алгоритма Рабина-Карпа 1.3, алгоритма Бойера-Мура 1.4 и алгоритма Кнута-Морриса-Пратта 1.5.

```

1 from naive import naive
2 from kmp import kmp
3 from rk import rk
4 from bm import bm
5 import time
6 nums = [1, 1]
7
8 for i in range(500-2):
9     nums.append(nums[-1] + nums[-2])
10
11 haystack = ''.join(str(num) for num in nums)
12
13 def test(f):
14     ans = (-1, -1)
15     for i in range(10, 100):
16         n = len(f(haystack, str(i)))
17         if n > ans[1]:
18             ans = (i, n)
19     return ans
20
21 def measure(f):
22     correct_ans = (71, 297)
23     start = time.time()
24     ans = test(f)
25     if ans != correct_ans:
26         raise Exception
27     end = time.time()
28     return end - start
29
30 algs = {
31     "naive": naive,
32     "rk": rk,
33     "kmp": kmp,
34     "bm": bm
35 }
36
37 for name, fn in algs.items():
38     print(f'{name}: {int(measure(fn) * 1000)}ms')

```

Рисунок 1.1 — Код для теста алгоритмов

```

1 def naive(haystack, needle):
2     result = []
3     for i in range(len(haystack) - len(needle) + 1):
4         if haystack[i:i+len(needle)] == needle:
5             result.append(i)
6     return result

```

Рисунок 1.2 — Наивный поиск подстроки

```

1 q = 101
2 d = 256
3
4 def rk(haystack, needle):
5     n = len(needle)
6
7     needle_h = 0
8     for sym in needle:
9         needle_h = (needle_h * d + ord(sym)) % q
10
11     window_h = 0
12     for sym in haystack[:n]:
13         window_h = (window_h * d + ord(sym)) % q
14
15     ans = []
16     if window_h == needle_h:
17         ans.append(0)
18
19     k = d ** (n - 1)
20     for l in range(len(haystack) - n):
21         old_sym = haystack[l]
22         new_sym = haystack[l + n]
23
24         window_h = d * (window_h - k * ord(old_sym)) + ord(new_sym)
25         window_h %= q
26
27         if window_h == needle_h and haystack[l + 1:l + n + 1] == needle:
28             ans.append(l + 1)
29
30     return ans

```

Рисунок 1.3 — Алгоритм Рабина-Карпа

```

1 N = 256
2
3 def bad_char_heuristic(string, size):
4     bad_chars = [-1] * N
5
6     for i, sym in enumerate(string[:size]):
7         bad_chars[ord(sym)] = i
8
9     return bad_chars
10
11
12 def bm(haystack, needle):
13     m = len(needle)
14     n = len(haystack)
15     ans = []
16
17     bad_chars = bad_char_heuristic(needle, m)
18
19     s = 0
20     while(s <= n - m):
21         j = m - 1
22
23         while j >= 0 and needle[j] == haystack[s+j]:
24             j -= 1
25
26         if j < 0:
27             ans.append(s)
28             s += (m-bad_chars[ord(haystack[s+m])] if s+m < n else 1)
29         else:
30             s += max(1, j-bad_chars[ord(haystack[s+j])])
31
32     return ans

```

Рисунок 1.4 — Алгоритм Бойера-Мура

```

1  def kmp(haystack, needle):
2      text = f'{needle}${haystack}'
3      prefix = [0] * len(text)
4      ans = []
5
6      i, j = 1, 0
7
8      while i < len(text):
9          if text[i] == text[j]:
10             j += 1
11             prefix[i] = j
12             if j == len(needle):
13                 ans.append(j)
14             i += 1
15             continue
16
17             if j != 0:
18                 j = prefix[j - 1]
19             continue
20
21         i += 1
22
23     return ans

```

Рисунок 1.5 — Алгоритм Кнута-Морриса-Пратта

На рисунке 1.6 приведен скриншот результатов измерения скорости выполнения алгоритмов для исходных данных.


```
aisd python test.py  
naive: 189ms  
rk: 322ms  
kmp: 202ms  
bm: 313ms
```

Рисунок 1.6 — Результаты измерения

Наивный алгоритм оказался самым быстрым в данном случае из-за маленькой длины шаблона и дешевого сравнения. Алгоритм Рабина-Карпа оказался довольно медленным, потому что пересчитывание хэша окна в случае шаблона из двух символов занимает столько же времени, сколько подсчёт нового хэша. Если бы мы сохраняли хэши для всех окон при поиске различных чисел, то этот алгоритм был бы значительно быстрее. Алгоритм КМП работает со скоростью наивного из-за небольшой длины шаблона. Алгоритм Бойера-Мура довольно медленный из-за небольшого алфавита, потому что эвристика плохих символов в этом случае работает плохо.

2 Задача 2

В данной задаче мы применяем алгоритм, похожий на алгоритм Рабина-Карпа. Сначала хэшируются все слова исходного текста, эти слова объединяются в тройки, для каждой тройки считается хэш на основе хэшей входящих в неё слов. Потом алгоритм итерирует по тройкам слов из статьи, и каждая тройка сопоставляется с увиденными в исходном тексте на основе хэшей сначала тройки, потом слов, потом содержания слов.

На рисунках 2.1, 2.2 и 2.3 приведен код программы, состоящий из хэш-функций для строк и для слов, функция-итераторов по словам и строкам слов и функция для вычисления процента плагиата.

```
1 import time
2 d = 256
3 k = d ** 2
4 q = 101
5 def hash_w(word):
6     h = 0
7     for s in word:
8         n = ord(s)
9         if n > 1000: # if russian
10             n -= 1040
11
12         h = (h * d + n) % q
13     return h
14
15 def hash_t(t):
16     return (t[0][0] * k + t[1][0] * d + t[2][0]) % q
17
18 def roll_t(h, old, new):
19     return ((h - old[0] * k) * d + new[0]) % q
20
21 def len_t(triplet):
22     return len(triplet[0][1]) + len(triplet[1][1]) + len(triplet[2][1])
23
24 def eq_t(t1, t2):
25     eq_h = t1[0][0] == t2[0][0] and t1[1][0] == t2[1][0] and t1[2][0] == t2[2][0]
26     if not eq_h:
27         return False
28     return t1[0][1] == t2[0][1] and t1[1][1] == t2[1][1] and t1[2][1] == t2[2][1]
```

Рисунок 2.1 — Хэш функции для строк и для слов

```

30 def words(text):
31     word_start = None
32     for i, s in enumerate(text):
33         if word_start is None and s.isalpha():
34             word_start = i
35             continue
36
37         if word_start is not None and not s.isalpha():
38             content = text[word_start:i]
39             word_h = hash_w(content)
40             word_start = None
41             yield word_h, content
42
43     if word_start != None:
44         yield hash_w(text[word_start:]), text[word_start:]
45
46 def triplets(text):
47     iter_words = words(text)
48     w1, w2, w3 = next(iter_words), next(iter_words), next(iter_words)
49
50     t1 = (w1, w2, w3)
51     triplet_h = hash_t(t1)
52     yield triplet_h, t1
53
54     last_words = [w2, w3]
55
56     for word in iter_words:
57         # triplet_h = roll_t(triplet_h, last_words[0], word)
58         triplet_h = hash_t((*last_words, word))
59         yield triplet_h, (*last_words, word)
60
61         last_words[0] = last_words[1]
62         last_words[1] = word

```

Рисунок 2.2 — Функции-итераторы по словам и по тройкам

```

64 # word: (hash, content)
65 def hash_triplets(text):
66     table = [[] for _ in range(q)] # triplet_hash: (word, word, word)
67     for h, triplet in triplets(text):
68         # if h not in table:
69         #     table[h] = []
70         table[h].append(triplet)
71     return table
72
73 def plagiat(source, text):
74     source_triplets = hash_triplets(source)
75     ans = 0
76     last = -5
77     for i, (h, t1) in enumerate(triplets(text)):
78         # if h not in source_triplets:
79         #     continue
80         for t2 in source_triplets[h]:
81             if eq_t(t1, t2):
82                 if i - last == 1:
83                     ans += len(t1[2][1])
84                     continue
85                 if i - last == 2:
86                     ans += len(t1[2][1]) + len(t1[1][1])
87                     continue
88                 ans += len_t(t1)
89                 last = i
90     return ans / len(text)
91
92 text = open('text.txt').read()
93 wiki = open('wiki.txt').read()
94
95 s = time.time()
96 p = plagiat(wiki, text)
97 e = time.time()
98
99 print(f'{int(p * 100)}%: {int((e - s) * 1000)}ms')

```

Рисунок 2.3 — Подсчет плагиата между двумя текстами

На рисунке 2.4 приведен результат выполнения программы.



```

aisd python plagiat.py
31%: 27ms

```

Рисунок 2.4 — Результат выполнения алгоритма

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Википедия. Статья "Логика" для антиплагиата [Электронный ресурс]: [сайт]. - URL: <https://ru.wikipedia.org/wiki/%D0%9B%D0%BE%D0%B3%D0%B8%D0%BA%D0%B0> (дата обращения: 22.03.2024).
2. Википедия. Алгоритм Бойера-Мура [Электронный ресурс]: [сайт]. - URL: https://en.wikipedia.org/wiki/Boyer%E2%80%93Moore_string-search_algorithm (дата обращения: 22.03.2024).
3. Википедия. Алгоритм Рабина-Карпа [Электронный ресурс]: [сайт]. - URL: https://en.wikipedia.org/wiki/Rabin%E2%80%93Karp_algorithm (дата обращения: 22.03.2024).
4. Википедия. Алгоритм Кнута-Морриса-Пратта [Электронный ресурс]: [сайт]. - URL: https://en.wikipedia.org/wiki/Knuth%E2%80%93Morris%E2%80%93Pratt_algorithm (дата обращения: 22.03.2024).