

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной  
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №1 по курсу «Дискретный анализ»

Студент: А. О. Ларченко  
Преподаватель: А. А. Кухтичев  
Группа: М8О-206Б  
Дата:  
Оценка:  
Подпись:

Москва, 2024

## Лабораторная работа №4

**Задача:** Необходимо реализовать один из стандартных алгоритмов поиска образцов для указанного алфавита.

**Вариант алгоритма:** Поиск одного образца при помощи алгоритма Апостолико-Джанкарло.

**Вариант алфавита:** Числа в диапазоне от 0 до  $2^{32}-1$  .

Запрещается реализовывать алгоритмы на алфавитах меньшей размерности, чем указано в задании.

# 1 Описание

Требуется написать реализацию алгоритма Апостолико-Джанкарло для поиска одного образца в тексте.

Алгоритма Апостолико-Джанкарло построен на основе алгоритма Бойера-Мура, но является его улучшенной версией. Он не меняет то, что уже использовалось в алгоритме Бойера-Мура (проход по образцу справа-налево, использование эвристик плохого символа и хорошего суффикса для эффективного сдвига при несовпадении), а добавляет к этому эффективный сдвиги, которые позволяют нам избежать лишних сравнений. [1]

Основная идея заключается в использовании большей полезной информации из правила хорошего суффикса. Для реализации алгоритма будем использовать дополнительный массив  $M$  длиной равной длине исходного текста, а также массив  $N$ , взятый из правила хорошего суффикса, который показывает, что  $N_i(P)$  - длина наибольшего суффикса  $P[1..i]$ , совпадающего с суффиксом  $P$ . [1]

Также как и в алгоритме Бойера-Мура мы идем по тексту  $T$  указателем  $j$ , а по образцу  $P$  указателем  $i$ , но теперь при каждом сдвиге смотрим на значения  $N_i$  и  $M_j$ . Тут возможно 4 различных случая:

1)  $M_j$  не определено или  $N_i == 0 == M_j$ . Тогда если  $P_i == T_j$  и  $i > 0$ , то декрементируем  $i$  и  $j$ , если  $i == 0$ , то запоминаем вхождение образца и  $M_h$  присваиваем значение  $m$ , где  $h$  - правая граница текста, к которой в текущий момент приложен правый конец образца, а  $m$  - длина образца. Иначе - сдвиг по правилу Бойера-Мура, а  $M_h$  присваиваем значение  $h - j$ .

2)  $M_j < N_i$ . Эта ситуация означает что если подставить к  $T_j$  часть совпадет, т.к. совпадает суффикс. Поэтому уменьшаем  $i$  и  $j$  на значение  $M_j$ .

3)  $M_j == N_i$ . Такой же смысл как и в пункте 2), поэтому уменьшаем  $i$  и  $j$  на значение  $M_j$ . Но при этом появляется условие, что если  $N_i == i$ , то мы фиксируем вхождение.

4)  $M_j > N_i$ . Эта ситуация означает что есть совпадение суффиксу, но следующий символ обязательно не совпадет (если это не начало паттерна). Поэтому если  $N_i == i$ , то мы фиксируем вхождение, иначе  $M_h$  присваиваем значение  $h - j$  и делаем сдвиг по правилу Бойера-Мура.

## 2 Исходный код

Т.к. алгоритм Апостолико-Джанкарло основывается на алгоритме Бойера-Мура, и использует нотации хорошего суффикса и плохого символа, то удобно сразу инициализировать эти нотации. Это позволит нам получать сдвиги за  $O(1)$ . Поэтому создадим класс *TPatern*.

```
1 template<class T>
2 class TPatern{
3     private:
4         vector<T> pattern_data;
5         unordered_map<T, vector<uint64_t>> bad_symbol_pos;
6         vector<uint64_t> good_suf_pos;
7
8         void bad_symbols_init();
9         void good_suffix_init();
10
11         uint64_t bad_symbol_shift(uint64_t cur_idx, T mismatch_val);
12         uint64_t good_suffix_shift(uint64_t cur_idx);
13     public:
14         vector<uint64_t> N_array;
15
16         TPatern(vector<T> input_array){
17             pattern_data = move(input_array);
18             bad_symbols_init();
19             good_suffix_init();
20         }
21         uint64_t size(){
22             return pattern_data.size();
23         }
24         T operator [] (uint64_t idx) const{
25             return this->pattern_data[idx];
26         }
27         T& operator [] (const uint64_t idx){
28             return this->pattern_data[idx];
29         }
30         uint64_t get_shift(uint64_t cur_pattern_idx, T mismatch_val);
31
32         void print_bad_symbol_array();
33         void print_good_suffix_array();
34         ~TPatern(){};
35 };
```

Функция	Описание
void bad_symbols_init()	Инициализирует правило плохого символа (заполняет словарь unordered_map<T, vector<uint64_t> bad_symbol_pos)
void good_suffix_init	Инициализирует правило хорошего суффикса (заполняет массив vector<uint64_t> N_array и vector<uint64_t> good_suf_pos)
uint64_t bad_symbol_shift(uint64_t cur_idx, T mismatch_val)	Возвращает сдвиг по правилу плохого символа для текущей позиции
uint64_t good_suffix_shift(uint64_t cur_idx)	Возвращает сдвиг по правилу хорошего суффикса для текущей позиции
TPatern(vector<T> input_array)	Вызывает инициализацию эвристик
uint64_t get_shift(uint64_t cur_pattern_idx, T mismatch_val)	Возвращает наибольший сдвиг для текущей позиции

Фнкция vector<T> pattern\_parser() отвечает за парсинг образца, а функция vector<TVal<T> text\_parser(uint64\_t start\_num\_line, uint64\_t start\_num\_val) за парсинг текста. Эти функции реализованы с использованием шаблона, что делает их универсальными и позволяет производить поиск образца разных типов данных без переписывания кода.

```

1  template<class T>
2  vector<T> pattern_parser(){
3      string str_input_data;
4      getline(cin, str_input_data);
5      istringstream input_stream(str_input_data);
6      vector<T> v_output_data;
7      T cur_val;
8      while(input_stream>>cur_val){
9          v_output_data.push_back(cur_val);
10     }
11     return v_output_data;
12 }
13
14 template<class T>
15 vector<TVal<T>> text_parser(uint64_t start_num_line, uint64_t start_num_val){
16     vector<TVal<T>> v_output_data;
17     string str_line;
18     uint64_t cnt_line = start_num_line;
19     while(getline(cin, str_line)){
20         istringstream input_stream(str_line);
21         T cur_val;
22         uint64_t cnt_val = start_num_val;
23         while(input_stream>>cur_val){

```

```

24         v_output_data.push_back(pair(cur_val, pair(cnt_line, cnt_val)));
25         cnt_val++;
26     }
27     cnt_line++;
28 }
29
30 return v_output_data;
31 }

```

Реализация основного тела алгоритма представлена в в функции *main*.

```

1  while(r_border<text_data.size()){
2      uint64_t i_t=r_border, i_p=p_size-1;
3      while(i_p>0 and text_data[i_t].first == pattern[i_p]){
4          if(M_vector[i_t]==uint64_t(-1) or (M_vector[i_t]==0 and M_vector[i_t]==
5              pattern.N_array[i_p])){
6              if (text_data[i_t].first == pattern[i_p]){
7                  i_p--;
8                  i_t--;
9              } else{
10                 M_vector[r_border] = r_border - i_t;
11                 break;
12             }
13         } else if (M_vector[i_t]!=uint64_t(-1)){
14             if (M_vector[i_t]<pattern.N_array[i_p]){
15                 i_p-=M_vector[i_t];
16                 i_t-=M_vector[i_t];
17             } else if(M_vector[i_t]==pattern.N_array[i_p]){
18                 if (pattern.N_array[i_p]==i_p){
19                     M_vector[r_border] = r_border-i_t;
20                 }
21                 i_p-=M_vector[i_t];
22                 i_t-=M_vector[i_t];
23             } else{
24                 if (pattern.N_array[i_p]==i_p){
25                     M_vector[r_border] = r_border-i_t;
26                     i_p-=M_vector[i_t];
27                     i_t-=M_vector[i_t];
28                 } else {
29                     M_vector[r_border] = r_border - i_t;
30                     break;
31                 }
32             }
33         }
34     }
35     if (i_p==0 and text_data[i_t].first == pattern[i_p]){
36         v_answ.push_back(text_data[i_t].second);
37         r_border++;
38     } else{

```

```
38 ||         r_border+= pattern.get_shift(i_p, text_data[i_t].first);
39 ||     }
40 || }
```

### 3 Консоль

```
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ g++ -std=c++20 main.cpp -o  
main  
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ cat <test  
11 45 11 45 90  
0011 45 011 0045 11 45 90    11  
45 11 45 90  
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./main <test  
1,3  
1,8arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$
```



## 4 Тест производительности

Тест производительности представляет из себя следующее: алгоритм Апостолико-Джанкарло сравнивается с наивным алгоритмом поиска подстроки в строке. В тест производительности также включен препроцессинг. Для поиска зависимостей было подготовлен 21 тест, которые условно можно разделить на 3 группы: тестирование зависимости скорости от размера паттерна, от размера текста и от частоты выпадения паттерна в тексте. Ниже представлены пара таких тестов.

```
// Тестирование зависимости скорости от частоты выпадения паттерна в тексте
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./benchmark <tests/01.t
Probability = 0.1
Pattern size = 7354
Text size = 1001817
Apostolico_Giancarlo algorithm time: 0.016757281 s
Dumm algorithm time: 0.007159800 s
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./benchmark <tests/05.t
Probability = 0.98
Pattern size = 7354
Text size = 1004459
Apostolico_Giancarlo algorithm time: 0.014491209 s
Dumm algorithm time: 0.007348549 s

// Тестирование зависимости скорости от размера паттерна
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./benchmark <tests/06.t
Probability = 0.59
Pattern size = 3
Text size = 1000004
Apostolico_Giancarlo algorithm time: 0.053340360 s
Dumm algorithm time: 0.005611640 s
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./benchmark <tests/07.t
Probability = 0.59
Pattern size = 2252
Text size = 1000839
Apostolico_Giancarlo algorithm time: 0.012560133 s
Dumm algorithm time: 0.006874890 s
arsenii@PC-Larcha14:~/Documents/C_pp_uk/DA/Lab_4$ ./benchmark <tests/10.t
Probability = 0.59
Pattern size = 8999
Text size = 1005222
Apostolico_Giancarlo algorithm time: 0.014650538 s
```

Dumm algorithm time: 0.006812823 s

// Тестирование зависимости скорости от размера текста

Probability = 0.59

arsenii@PC-Larcha14:~/Documents/C\_pp\_uk/DA/Lab\_4\$ ./benchmark <tests/11.t

Pattern size = 7354

Text size = 14846

Apostolico\_Giancarlo algorithm time: 0.005261117 s

Dumm algorithm time: 0.000103084 s

arsenii@PC-Larcha14:~/Documents/C\_pp\_uk/DA/Lab\_4\$ ./benchmark <tests/16.t

Probability = 0.59

Pattern size = 7354

Text size = 511085

Apostolico\_Giancarlo algorithm time: 0.009006845 s

Dumm algorithm time: 0.003601026 s

arsenii@PC-Larcha14:~/Documents/C\_pp\_uk/DA/Lab\_4\$ ./benchmark <tests/21.t

Probability = 0.59

Pattern size = 7354

Text size = 1006971

Apostolico\_Giancarlo algorithm time: 0.014268859 s

Dumm algorithm time: 0.007304772 s

Если не брать тот факт, что алгоритм Апостолико-Джанкарло проигрывает на всех тестах наивному алгоритму, что очень странно (т.к. скорость наивного алгоритма -  $O(n*m)$ , а теоретическая скорость алгоритм Апостолико-Джанкарло -  $O(m)$ , т.к. в отличие от алгоритма Бойера-Мура (на основе которого построен алгоритм Апостолико-Джанкарло), который в худшем случае работает за  $O(n*m)$ , алгоритм Апостолико-Джанкарло использует всю полезную информацию о сдвигах.  $n$  - длина паттерна, а  $m$  - длина текста), то можно изучить закономерности. Например, что при увеличении вероятности выпадения паттерна, время алгоритма уменьшается, тогда как скорость наивного алгоритма повышается. При увеличении длины паттерна происходит аналогичный процесс. А вот при увеличении размера текста время выполнения обоих алгоритмом повышается.

У меня есть теория, что такие плохие результаты по сравнению с наивным алгоритмом получены из-за искусственно-составленных неидеальных тестов. Т.к. по условию задачи мы используем алфавит из  $2^{32} - 1$  элементов, вероятность случайного выпадения одинаковых чисел в паттерне, длиной от 3 до 10000 из этого диапазона мала, следовательно основные нотации, которые использует алгоритм (правило хорошего суффикса и правило плохого символа) не играют большой роли и не помогают ускорить алгоритм, в то время как их препроцессинг, который, стоит отметить, занимает

$\text{const} * O(n)$ , требует время. Моей ошибкой в данном случае является игнорирование повышения вероятности выпадения одиноковых символов в паттерне.

## 5 Выводы

В этой лабораторной работе я познакомился с алгоритмами поиска образца в тексте. Получил много новых знаний об алгоритме Бойера-Мура, в частности об его эвристиках (правилах плохого символа и хорошего суффикса, а также об их сильных реализациях) и об его улучшенном варианте - алгоритме Апостолико-Джанкарло, а также реализовал эти два алгоритма.

Данная лабораторная сложностей у меня не вызвала, кроме тестирования производительности алгоритма, которое сложно назвать удачным и связано это, как я полагаю, с плохими искусственными данными. В реальной жизни, при работе с алфавитами меньшей мощности алгоритм Апостолико-Джанкарло показывает себя гораздо лучше и работает за  $O(m)$ .

## Список литературы

[1] Гасфилд Дэн

Строки, деревья и последовательности в алгоритмах: Информатика и вычислительная биология / Пер. с англ. И. В. Романовского. — СПб.: Невский Диалект; БХВ-Петербург, 2003. — 654 с: ил.