

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Дискретный анализ»

Студент: П. А. Мохляков
Преподаватель: Н. С. Капралов
Группа: М8О-208Б
Дата:
Оценка:
Подпись:

Москва, 2021

Лабораторная работа №4

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

Вариант дерева: Поиск одного образца-маски: в образце может встречаться «джокер» (представляется символом ? — знак вопроса), равный любому другому символу. При реализации следует разбить образец на несколько, не содержащих «джокеров», найти все вхождения при помощи алгоритма Ахо-Корасик и проверить их относительное месторасположение.

Вариант алгоритма: Последовательность букв английского алфавита длиной не более 256 символов

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

1 Описание

Алогоритм Ахо-Корасик - это алшоритм поиска подстроки в строке. В отсличии от других алгоритмов, алогоритм Ахо-Корасик позволяет деталь поиск сразу несколько паттернов одноврененно, то есть за один проход. Благодаря этому скорость поиска множества паттернов намного выше, чем в других алгоритмах.

Рассмотрим данный алгоритм на примере моего исходного кода.

2 Исходный код

Данный алгоритм можно разбить на несколько основных частей:

- Создание trie из паттернов.
- Создание суффиксных ссылок.
- Поиск в тексте.

Добавление паттерна в trie происходит поэлементно. Стартовой позицией является Root, мы смотрим, есть ли у элемента ребенок в виде нашего введенного элемента. Если он существует, то мы просто переходим в него, иначе создаем, добавляем в таблицу детей и переходим. Дойдя до конца паттерна мы говорим, что это последний элемент и добавляем его позицию с учетом джокеров.

```
1 void Push(const std::vector<long long> &pattern, int pos){
2     node_ptr cur = Root;
3     for(auto elem_pattern: pattern){
4         node_ptr child = cur->ChildPtr(elem_pattern);
5         if(child == nullptr){
6             child = new TNode(elem_pattern);
7             child->Parent = cur;
8             cur->Childs.insert({elem_pattern,child});
9         }
10        cur = child;
11    }
12    cur->Last = true;
13    cur->Jockers.push_back(pos);
14 }
```

Далее происходит создание ссылок на суффиксы и концы префиксов, для этого нужно обойти trie в ширину. При поиске суффикса мы переходим в родителя, и далее переходим по суффиксам родителя до тех пор пока не у суффикса не будет в качестве ребенка наш элемент, в таком случае это суффикс нашего элемента, иначе идем пока не дойдем до нулевого указателя, в таком случае суффикс нашего элемента это Root. Для проверки на концы префиксов мы смотрим является ли наш суффикс концом паттерна или же указывает ли он другой конец префикса.

```
1 void BorSuf(){
2     //BFS
3     std::queue<node_ptr> queue;
4     queue.push(Root);
5     while(!queue.empty()){
6         node_ptr elem = queue.front();
7         queue.pop();
8         for(auto child: elem->Childs)
9             {
```

```

10     queue.push(child.second);
11 }
12 //~BFS
13 //Suf
14 if(elem == Root) //not for Root
15 continue;
16
17 node_ptr parent = elem->Parent;
18 parent = parent->Suffix;
19 while(parent != nullptr && parent->ChildPrt(elem->Sym) == nullptr)
20     parent = parent->Suffix;
21
22 if(parent == nullptr)
23     elem->Suffix = Root;
24 else
25     elem->Suffix = parent->ChildPrt(elem->Sym);
26 //~Suf
27 //SubLast
28 if(elem->Suffix->Last)
29     elem->SubLast = elem->Suffix;
30 else if(elem->Suffix->SubLast != nullptr)
31     elem->SubLast = elem->Suffix->SubLast;
32 //~SubLast
33 }//while
34 }

```

При поиске мы смотрим на символ и переходим в ребенка с этим символом, если такого не существует, то переходим по суффиксу пока не сможем перейти или пока не упрямся в Root. В случае если мы достигли элемента конца паттерна или элемента со ссылок на префиксы конца паттерна, то мы добавляем единицу в вектор вхождений в позиции текста минус позиция паттерна с учетом джокеров. После поиска мы смотрим на вектор вхождений, в точках где значение равно количеству паттернов и есть начальные точки вхождений паттернов.

```

1 void ThisLast(node_ptr Cur, std::vector<int> &vect_incl, int pos){
2     if(Cur->SubLast != nullptr)
3         ThisLast(Cur->SubLast, vect_incl, pos);
4     for(auto jockey: Cur->Jockers){
5         if(pos-jockey >= 0)
6             ++vect_incl[pos-jockey];
7     }
8 }
9
10 void Find(std::vector<SText> &text, std::vector<int> &pos_incl){
11     node_ptr cur = Root;
12     int i=0;
13     for(auto elem: text){
14         while(cur->ChildPrt(elem.Sym)==nullptr && cur->Suffix!=nullptr)
15             cur = cur->Suffix;

```

```

16         if(cur->ChildPrt(elem.Sym)!=nullptr)
17             cur = cur->Childs.at(elem.Sym);
18         if(cur->Last)
19             ThisLast(cur,pos_incl,i);
20         else if(cur->SubLast != nullptr)
21             ThisLast(cur->SubLast,pos_incl,i);
22
23         ++i;
24     }
25 }

```

Листинг:

```

1  #pragma once
2
3  #include <vector>
4  #include <unordered_map>
5  #include <queue>
6  #include <iostream>
7
8  struct SText
9  {
10     long long Sym;
11     int str;
12     int word;
13 };
14
15
16 class TAhoKorasik{
17 protected:
18 class TNode{
19     public:
20     using node_ptr = TNode*;
21     long long Sym; //Data
22     bool Last;
23     std::vector<int> Jockers;
24     node_ptr Parent; //Ptrs
25     node_ptr Suffix;
26     node_ptr SubLast;
27     std::unordered_map<long long,node_ptr> Childs;
28
29     TNode(long long inSym = -1,bool inLast = false):Sym(inSym),Last(inLast){
30         Parent = nullptr;
31         Suffix = nullptr;
32         SubLast = nullptr;
33     }
34
35     ~TNode(){
36         for(auto Child: Childs){
37             delete(Child.second);

```

```

38     }
39 }
40
41 node_ptr ChildPrt(long long &inSym){
42     std::unordered_map<long long,node_ptr>::iterator child = Childs.find(inSym);
43     if(child == Childs.end())
44         return nullptr;
45     else
46         return child->second;
47 }
48 };//TNode
49 using node_ptr = TNode::node_ptr;
50
51 node_ptr Root;
52
53 public:
54
55 TAhoKorasik(){
56     Root = new TNode();
57 }
58 ~TAhoKorasik(){
59     delete(Root);
60 }
61
62 void Push(const std::vector<long long> &pattern, int pos){
63     node_ptr cur = Root;
64     for(auto elem_pattern: pattern){
65         node_ptr child = cur->ChildPrt(elem_pattern);
66         if(child == nullptr){
67             child = new TNode(elem_pattern);
68             child->Parent = cur;
69             cur->Childs.insert({elem_pattern,child});
70         }
71         cur = child;
72     }
73     cur->Last = true;
74     cur->Jockers.push_back(pos);
75 }
76
77 void BorSuf(){
78     //BFS
79     std::queue<node_ptr> queue;
80     queue.push(Root);
81     while(!queue.empty()){
82         node_ptr elem = queue.front();
83         queue.pop();
84         for(auto child: elem->Childs)
85         {
86             queue.push(child.second);

```

```

87     }
88     //~BFS
89     //Suf
90     if(elem == Root) //not for Root
91         continue;
92
93     node_ptr parent = elem->Parent;
94     parent = parent->Suffix;
95     while(parent != nullptr && parent->ChildPrt(elem->Sym) == nullptr)
96         parent = parent->Suffix;
97
98     if(parent == nullptr)
99         elem->Suffix = Root;
100    else
101        elem->Suffix = parent->ChildPrt(elem->Sym);
102    //~Suf
103    //SubLast
104    if(elem->Suffix->Last)
105        elem->SubLast = elem->Suffix;
106    else if(elem->Suffix->SubLast != nullptr)
107        elem->SubLast = elem->Suffix->SubLast;
108    //~SubLast
109 } //while
110 }
111
112 void ThisLast(node_ptr Cur, std::vector<int> &vect_incl, int pos){
113     if(Cur->SubLast != nullptr)
114         ThisLast(Cur->SubLast, vect_incl, pos);
115     for(auto jockey: Cur->Jockers){
116         if(pos-jockey >= 0)
117             ++vect_incl[pos-jockey];
118     }
119 }
120
121 void Find(std::vector<SText> &text, std::vector<int> &pos_incl){
122     node_ptr cur = Root;
123     int i=0;
124     for(auto elem: text){
125         while(cur->ChildPrt(elem.Sym)==nullptr && cur->Suffix!=nullptr)
126             cur = cur->Suffix;
127         if(cur->ChildPrt(elem.Sym)!=nullptr)
128             cur = cur->Childs.at(elem.Sym);
129         if(cur->Last)
130             ThisLast(cur, pos_incl, i);
131         else if(cur->SubLast != nullptr)
132             ThisLast(cur->SubLast, pos_incl, i);
133
134         ++i;
135     }

```



```
136 || }  
137 ||  
138 || };//TAhoKorasik
```

3 Консоль

```
pavel@DESKTOP-VBSMFB3:~/Projects/mai/2_course/DA/LB4/second$ cat ../test
1 ? 02
0001
1
02
2
pavel@DESKTOP-VBSMFB3:~/Projects/mai/2_course/DA/LB4/second$ ./solution <../test
1,1
2,1
```

4 Тест производительности

Сравнение производительности будет производиться с `std::find`, представленным в стандартной библиотеке шаблонов C++.

	Длина паттерна	Количество строк	Количество слов в строке
test1.txt	10	10	100
test2.txt	1000	10000	100
test3.txt	10 где джокеры и символы через один	10000	10

```
pavel@DESKTOP-VBSMFB3:~/solution$ make
pavel@DESKTOP-VBSMFB3:~/solution$ ./solution <test1.txt
Aho-AhoKorasik: 1ms
std::find: 1ms
pavel@DESKTOP-VBSMFB3:~/solution$ ./solution <test2.txt
Aho-AhoKorasik: 5ms
std::find: 12ms
pavel@DESKTOP-VBSMFB3:~/solution$ ./solution <test3.txt
Aho-AhoKorasik: 7
std::find: 11
```

5 Выводы

Выполнив вторую лабораторную работу по курсу «Дискретный анализ», я реализовал алгоритм поиска подстроки в строке Ахо-Корасик.

Данный алгоритм отличается быстротой поиска множества паттернов, за счет того, что поиски всех паттернов происходят одновременно за один проход по тексту. Тем не менее при высокой энтропии текста и малом количестве паттернов этот алгоритм может быть сравним по скорости с наивным алгоритмом.