

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №5 по курсу «Дискретный анализ»

Студент: А. О. Знай
Преподаватель: С. А. Михайлова
Группа: М8О-301Б-21
Дата: 25.09.2023
Оценка:
Подпись:

Москва, 2023

Лабораторная работа №5

Задача: Необходимо реализовать поиск паттерна в тексте с использованием суффиксного дерева.

Вариант алгоритма: Найти образец в тексте используя статистику совпадений.

1 Описание

Суффиксное дерево — бор, содержащий все суффиксы некоторой строки (и только их). Позволяет выяснять, входит ли строка w в исходную строку t , за время $O(|w|)$, где $|w|$ — длина строки w .

2 Описание программы

Построение суффиксного дерева будем осуществлять с помощью алгоритма Укконена за $O(n)$. Мы должны добавить все суффиксы каждого префикса строки, но делать это используя несколько эвристик: 1) Мы должны использовать линейное количество памяти, поэтому на ребрах хранятся не символы строки, а 2 числа - позиция самого левого и самого правого ее символов в тексте. 2) Автоматическое продление всех листов за $O(1)$ благодаря увеличению переменной текущего конца на единицу. 3) Использование суффиксных ссылок: если есть суффикс $ха$, то есть и суффикс $а$, который мы можем пропустить перейдя по суффиксной ссылке между двумя вершинами 4) Переход по счетчику внутри ребер

В данном варианте для поиска строки в тексте используется статистика совпадений. Ее смысл заключается в следующем: создается массив ms размером с текст, где в i ячейке хранится число k , которое означает, что начиная с i позиции в тексте совпадут k символов какой-то подстроки образца. Для заполнения такого массива необходимо: 1) построить суффиксное дерево не для текста, а для паттерна 2) считывать символы текста до тех пор пока есть совпадения в дереве и увеличиваем k , если произошла неудача, то записываем число в массив, откатываемся до ближайшей внутренней вершины, переходим по суффиксной ссылке и накатываем $k-1$ символов для текущего символа в тексте и продолжаем искать совпадения 3) пройтись по массиву и вывести все индексы, в которых лежат k , равные длине паттерна

Сложность алгоритма поиска статистикой совпадений - $O(m)$, где m - длина текста.

[1]

Код

```
1 #include <string>
2 #include <iostream>
3 #include <vector>
4 #include <stack>
5 #include <memory>
6
7 struct SuffixTree {
8     struct Edge {
9         int left;
10         std::shared_ptr<int> right;
11         int nodeTarget;
12
13         Edge(int left, std::shared_ptr<int> right, int node):
14             left(left),
15             right(right),
16             nodeTarget(node) {}
17     };
18
19     std::vector< std::vector<Edge> > Tree;
20     std::string Text;
21     std::vector<int> VSuffix;
22     std::vector<int> LenToNode;
23
24     int FindEdge(int node, char ch) {
25         for (int i = 0; i < Tree[node].size(); ++i) {
26             if (Text[Tree[node][i].left] == ch) {
27                 return i;
28             }
29         }
30         return -1;
31     };
32
33     SuffixTree(const std::string & str): Text(str),
34         Tree(str.size() * 2),
35         VSuffix(str.size() * 2),
36         LenToNode(str.size() * 2) {
37         int left = 0;
38         std::shared_ptr<int> end(new int);
39         *end = 0;
40
41         bool newNode = false;
42         int newNodeNumber = 0;
43
44         int eqFromRoot = 0;
45         int eqCur = 0;
46
47         int curNode = 0;
48         int curEdge = -1;
```

```

49
50     for (int i = 0; i < str.size(); ++i) {
51         ++*end;
52
53         int curNewNode = -1;
54         int lastNewNode = 0;
55
56         while (*end > left) {
57             if (curEdge == -1) {
58                 int nextEdge = FindEdge(curNode, Text[left + eqFromRoot + eqCur])
59                     ;
60
61                 if (nextEdge == -1) {
62                     newNode = true;
63                 }
64                 else {
65                     eqCur = 1;
66                 }
67
68                 curEdge = nextEdge;
69             }
70             else if (
71                 Text[Tree[curNode][curEdge].left + eqCur] ==
72                 Text[left + eqFromRoot + eqCur ]
73             ) {
74                 ++eqCur;
75             }
76             else {
77                 newNode = true;
78             }
79
80             if (eqCur > 0 and Tree[curNode][curEdge].left + eqCur == *Tree[
81                 curNode][curEdge].right) {
82                 curNode = Tree[curNode][curEdge].nodeTarget;
83                 curEdge = -1;
84
85                 eqFromRoot += eqCur;
86                 eqCur = 0;
87             }
88             if (newNode) {
89                 ++newNodeNumber;
90
91                 if (eqCur == 0 and curEdge == -1) {
92                     Tree[curNode].push_back(Edge(left + eqFromRoot, end,
93                         newNodeNumber));
94                 }
95                 else {
96                     curNewNode = newNodeNumber;
97                     Edge tempCurEdge = Tree[curNode][curEdge];

```

```

95
96         std::shared_ptr<int> newRight(new int);
97         *newRight = Tree[curNode][curEdge].left + eqCur;
98         Tree[curNode][curEdge].right = newRight;
99
100        Tree[curNode][curEdge].nodeTarget = newNodeNumber;
101        tempCurEdge.left = *newRight;
102        Tree[newNodeNumber].push_back(tempCurEdge);
103
104        ++newNodeNumber;
105        Tree[newNodeNumber - 1].push_back(Edge(left + eqCur +
106                                                eqFromRoot, end, newNodeNumber));
107
108        LenToNode[newNodeNumber - 1] = LenToNode[curNode] + *Tree[
109            curNode][curEdge].right - Tree[curNode][curEdge].left;
110
111        if (lastNewNode > 0) {
112            VSuffix[lastNewNode] = curNewNode;
113        }
114        lastNewNode = curNewNode;
115    }
116
117    int nextNode = VSuffix[curNode];
118    int nextEdge = -1;
119    int nextEqFromRoot = LenToNode[nextNode];
120    int nextEqCur = eqCur + eqFromRoot - nextEqFromRoot - 1;
121
122    while (nextEqCur > 0) {
123        nextEdge = FindEdge(nextNode, Text[left + nextEqFromRoot +
124            1]);
125
126        int r = *Tree[nextNode][nextEdge].right;
127        int l = Tree[nextNode][nextEdge].left;
128
129        if (r - 1 <= nextEqCur) {
130            nextEqFromRoot += r - 1;
131            nextEqCur -= r + 1;
132            nextNode = Tree[nextNode][nextEdge].nodeTarget;
133            nextEdge = -1;
134        }
135        else {
136            break;
137        }
138    }
139
140    if (nextEdge != -1) {
141        eqCur = nextEqCur;
142    }
143    else {

```

```

141         eqCur = 0;
142     }
143
144     eqFromRoot = nextEqFromRoot;
145     curNode = nextNode;
146     curEdge = nextEdge;
147     left++;
148     newNode = false;
149 }
150
151 }
152 }
153 }
154
155 std::vector<size_t> MatchStatistics(const std::string &text) {
156     size_t len = text.size();
157     std::vector<size_t> ms(len);
158
159     int curNode = 0;
160     int curEdge = FindEdge(curNode, text[0]);
161     int eqFromRoot = 0;
162     int eqCur = 0;
163
164     std::stack<int> path;
165     path.push(0);
166
167     for (int i = 0; i < len; i++) {
168         int nextNode = curNode;
169
170         while(VSuffix[nextNode] == 0 and !path.empty()) {
171             nextNode = path.top();
172             path.pop();
173         }
174         curNode = VSuffix[nextNode];
175
176         path.push(curNode);
177         eqCur = std::max(0, eqFromRoot + eqCur - 1 - LenToNode[curNode]);
178         eqFromRoot = LenToNode[curNode];
179         curEdge = -1;
180
181         while (eqCur >= 0) {
182             int nextEdge = FindEdge(curNode, text[i + eqFromRoot]);
183
184             if(nextEdge != -1) {
185                 int edgeLen = *Tree[curNode][nextEdge].right - Tree[curNode][
                    nextEdge].left;
186
187                 if (eqCur >= edgeLen) {
188                     eqCur -= edgeLen;

```

```

189         eqFromRoot += edgeLen;
190         curNode = Tree[curNode][nextEdge].nodeTarget;
191         nextEdge = -1;
192     }
193     else{
194         curEdge = nextEdge;
195         break;;
196     }
197 }
198 else break;
199 }
200 while (
201     curEdge != -1 and i + eqFromRoot + eqCur < len and text[i + eqCur +
202         eqFromRoot]
203     ==
204     Text[Tree[curNode][curEdge].left + eqCur]
205 ) {
206     ++eqCur;
207     if (Tree[curNode][curEdge].left + eqCur == *Tree[curNode][curEdge].right
208         ) {
209         int nextNode = Tree[curNode][curEdge].nodeTarget;
210         int nextEdge = FindEdge(nextNode, text[i + eqFromRoot + eqCur]);
211         int edgeLen = eqCur;
212         eqFromRoot += edgeLen;
213         eqCur = 0;
214         curNode = nextNode;
215         path.push(curNode);
216         curEdge = nextEdge;
217     }
218     ms[i] = eqFromRoot + eqCur;
219 }
220 return ms;
221 };
222
223 int main() {
224     std::ios_base::sync_with_stdio(false);
225     std::cin.tie(nullptr);
226     std::cout.tie(nullptr);
227
228     std::string pattern;
229     std::string text;
230
231     std::cin >> pattern;
232     pattern += "$";
233     SuffixTree tree(pattern);
234
235     std::cin >> text;

```



```

236 |     std::vector<size_t> ms(text.size());
237 |     ms = tree.MatchStatistics(text);
238 |
239 |     for (int i = 0; i < ms.size(); ++i) {
240 |         if (ms[i] == pattern.size() - 1) {
241 |             std::cout << i + 1 << "\n";
242 |         }
243 |     }
244 | }

```

3 Консоль

```
artemznaj@MacBook-Air-Artem lab5 % ./a.out
ba
baobabbaobaoba
1
4
7
10
13
artemznaj@MacBook-Air-Artem lab5 % ./a.out
aba
qababaz
2
4
artemznaj@MacBook-Air-Artem lab5 % ./a.out
bao
bacbaobabbaobabbac
4
10
artemznaj@MacBook-Air-Artem lab5 % ./a.out
baobab
baobabaobab
1
6
```

4 Тест производительности

Тест производительности представляет собой сравнение реализованного мной алгоритма статистики совпадений и STL find. Сравнение производится путем запуска поиска 400000 паттернов в тексте. Результат работы:

```
artemznaj@MacBook-Air-Artem lab5 % ./a.out <test.txt  
ms: 86481861 ms  
STL find: 202546125 ms
```

5 Выводы

Выполнив пятую лабораторную работу, мною была изучена такая структура данных, как суффиксное дерево. Научился строить его за $O(n)$, используя алгоритм Укконена, и искать строку в тексте (за $O(m)$, где m - длина текста) с использованием приложения к суффиксным деревьям - статистику совпадений.

Несмотря на очевидный плюс в виде сложности, у алгоритма Укконена есть и минусы, например, большие затраты по памяти относительно размера текста, также для простых задач поиска подстроки в строке лучше использовать другие алгоритмы.

Список литературы

- [1] String Algorithms and Algorithms in Computational Biology - Gusfield