

**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: Д. А. Тарпанов
Преподаватель: Н. С. Капралов
Группа: М8О-307Б-19
Дата:
Оценка:
Подпись:

Москва, 2021

Лабораторная работа №5

Задача:

Необходимо реализовать алгоритм Укконена построения суффиксного дерева за линейное время. Построив такое дерево для некоторых из выходных строк, необходимо воспользоваться полученным суффиксным деревом для решения своего варианта задания.

Алфавит строк: строчные буквы латинского алфавита (т.е. от a до z).

Вариант:

Найти самую длинную общую подстроку двух строк.

Формат входных данных

Две строки.

Формат результата

На первой строке нужно распечатать длину максимальной общей подстроки, затем перечислить все возможные варианты общих подстрок этой длины в порядке лексикографического возрастания без повторов.

1 Описание

Требуется реализовать алгоритм Укконена за время $O(n)$. Идея заключается в том, чтобы взять наивный алгоритм, работающий за $O(n^3)$ и оптимизировать его до линии. Наивный алгоритм будет добавлять все суффиксы каждого префикса строки. Эвристика 1: пусть в суффиксном дереве есть строка xa (x – первый символ, a – оставшаяся строка). Тогда a тоже будет в дереве, т.к. является суффиксом xa . Если для строки xa существует некоторая вершина u , то вершина u существует и для строки a . Ссылка из u в v называется суффиксной ссылкой. Она позволяет не проходить каждый раз по дереву из корня. Для построения суффиксных ссылок достаточно хранить номер последней созданной вершины. Если на этой фазе мы создаем еще одну новую вершину, то нужно построить суффиксную ссылку из предыдущей в текущую. Это улучшение позволяет сократить временную сложность до $O(n^2)$.

Эвристика 2: для ускорения до линии необходимо оптимизировать хранение данных в дереве. Вместо строк будем хранить индексы начала и конца, причем переменная с индексом конца будет общей для всех вершин дерева, и её инкрементирование будет проходить за $O(1)$.

При использовании вышеописанных эвристик временная сложность сократится до $O(n)$.

Для нахождения наибольшей общей подстроки двух строк конкатенируем строки, добавив между ними символ-разделитель. Построим суффиксное дерево для получившейся строки. Потом необходимо обойти дерево и отметить вершины, принадлежащие одной первой строке, вершины, принадлежащие второй строке, и вершины, принадлежащие обоим строкам. Тогда наибольшей общей подстрокой будет самая длинная строка, принадлежащая обоим строкам.

2 Исходный код

В файле SuffTree.cpp содержится класс суффиксного дерева, в файле main.cpp содержится функция main, осуществляющая ввод дерева и вызов функции поиска наибольшей общей подстроки.

Листинг main.cpp

```
1 | #include "SuffTree.cpp"
2 |
3 | int main() {
4 |     std::string a, b;
5 |     std::cin >> a >> b;
6 |     if (a.empty() || b.empty()) {
7 |         exit(0);
8 |     }
9 |     TSuffTree tree(a, b);
10 |    tree.LongestCommonSubstring();
11 | }
```

Листинг SuffTree.cpp

```
1 | #pragma once
2 |
3 | #include <iostream>
4 | #include <vector>
5 | #include <string>
6 | #include <algorithm>
7 | #include <unordered_map>
8 | #include <set>
9 |
10 | class TSuffTree {
11 | private:
12 |     class TNode {
13 |     private:
14 |         std::unordered_map<char, TNode*> children;
15 |         TNode *link;
16 |         int start, *end;
17 |         int ind;
18 |     public:
19 |         friend TSuffTree;
20 |         TNode(TNode *_link, int _start, int *_end, int _ind) : link(_link), start(
21 |             _start),
22 |             end(_end), ind(_ind) {}
23 |         TNode(TNode *_link, int _start, int *_end) : TNode(_link, _start, _end, -1) {}
24 |     };
25 | private:
```

```

25 void DeleteHelper(TNode* node) {
26     if (node == nullptr) {
27         return;
28     }
29     for (auto it : node->children) {
30         DeleteHelper(it.second);
31     }
32     if (node->ind < -1) {
33         delete node->end;
34     }
35     delete node;
36 }
37 int EdgeLength(TNode *node) {
38     return *(node->end) - (node->start) + 1;
39 }
40 void ExtendTree(int ind) {
41     internalNode = nullptr;
42     leafEnd++;
43     remainSuff++;
44     while (remainSuff > 0) {
45         if (activeLength == 0) {
46             activeEdge = ind;
47         }
48         auto findSymb = activeNode->children.find(text[activeEdge]);
49         if (findSymb == activeNode->children.end()) {
50             activeNode->children.insert(std::make_pair(text[activeEdge], new TNode(
51                 root, ind, &leafEnd, ind - remainSuff + 1)));
52             if (internalNode != nullptr) {
53                 internalNode->link = activeNode;
54                 internalNode = nullptr;
55             }
56             else {
57                 TNode *next = findSymb->second;
58                 int edgeLen = EdgeLength(next);
59                 if (activeLength >= edgeLen) {
60                     activeEdge += edgeLen;
61                     activeLength -= edgeLen;
62                     activeNode = next;
63                     continue;
64                 }
65                 if (text[next->start + activeLength] == text[ind]) {
66                     if (internalNode != nullptr && activeNode != root) {
67                         internalNode->link = activeNode;
68                     }
69                     activeLength++;
70                     break;
71                 }
72                 TNode *split = new TNode(root, next->start, new int(next->start +
73                     activeLength - 1));

```

```

72         activeNode->children[text[activeEdge]] = split;
73         next->start += activeLength;
74         split->children.insert(std::make_pair(text[ind], new TNode(root, ind, &
75             leafEnd, ind - remainSuff + 1)));
76         split->children.insert(std::make_pair(text[next->start], next));
77         if (internalNode != nullptr) {
78             internalNode->link = split;
79         }
80         internalNode = split;
81     }
82     remainSuff--;
83     if (activeNode == root && activeLength > 0) {
84         activeLength--;
85         activeEdge++;
86     } else if (activeNode != root) {
87         activeNode = activeNode->link;
88     }
89 }
90 void PrintHelper(TNode* p) {
91     static int level = 0;
92     level++;
93     if (p->start == -1) {
94         printf("\\_<root>");
95     }
96     else {
97         printf("_<%d, %d>[%d]", p->start, *p->end, p->ind);
98         //std::cout << text.substr(p->start, *p->end) << "[" << p->ind << "]";
99     }
100    for (int i = 0; i < (int)p->children.size(); i++) {
101        if (p->children[i]) {
102            printf("\n");
103            for (int j = 0; j < level; j++)
104                printf(" ");
105            printf("\\");
106            PrintHelper(p->children[i]);
107        }
108    }
109    level--;
110 }
111
112 int doTraversal(TNode *n, int labelHeight, int* maxHeight, std::vector<int> &
113     substringStartIndex) {
114     if(n == nullptr) {
115         return 0;
116     }
117     int ret = -1;
118     if(n->ind < 0) {
119         for (auto it : n->children) {

```

```

119         TNode *temp = it.second;
120         if(temp != nullptr) {
121             ret = doTraversal(temp, labelHeight + EdgeLength(temp), maxHeight,
122                             substringStartIndex);
123             if(n->ind == -1) {
124                 n->ind = ret;
125             }
126             else if((n->ind == -2 && ret == -3) || (n->ind == -3 && ret == -2)
127                    || n->ind == -4) {
128                 n->ind = -4;
129                 if(*maxHeight < labelHeight) {
130                     *maxHeight = labelHeight;
131                     substringStartIndex.clear();
132                     substringStartIndex.push_back(*(n->end) - labelHeight + 1);
133                 }
134                 else if(*maxHeight == labelHeight && !substringStartIndex.empty()
135                        && substringStartIndex.back() != *(n->end) - labelHeight +
136                        1) {
137                     substringStartIndex.push_back(*(n->end) - labelHeight + 1);
138                 }
139             }
140         }
141         else if(n->ind > -1 && n->ind < size1)
142             return -2;
143         else if(n->ind >= size1)
144             return -3;
145         return n->ind;
146     }
147
148     TNode *root = new TNode(nullptr, -1, new int(-1));
149     TNode *internalNode = nullptr;
150
151     std::string text;
152
153     TNode *activeNode = nullptr;
154     int activeEdge = -1;
155     int activeLength = 0;
156     int remainSuff = 0;
157     int leafEnd = -1;
158     int size1;
159 public:
160     void Print() {
161         PrintHelper(root);
162     }
163     TSuffTree(std::string &str1, std::string &str2) {

```

```

164     size1 = str1.size();
165     text = str1 + '#' + str2 + "$";
166     BuildSuffTree();
167 };
168 void LongestCommonSubstring() {
169     int maxHeight = 0;
170     std::vector<int> substringStartIndex;
171     doTraversal(root, 0, &maxHeight, substringStartIndex);
172     std::cout << maxHeight << "\n";
173     std::vector<std::string> ans;
174     std::string temp;
175     for (int i = 0; i < substringStartIndex.size(); ++i) {
176         int k = 0;
177         temp.clear();
178         for (k = 0; k < maxHeight; k++) {
179             temp += text[k + substringStartIndex[i]];
180         }
181         ans.push_back(temp);
182     }
183     std::stable_sort(ans.begin(), ans.end());
184     for(auto i:ans) {
185         std::cout << i << "\n";
186     }
187 }
188 void BuildSuffTree() {
189     activeNode = root;
190     for (int i = 0; i < text.length(); ++i) {
191         ExtendTree(i);
192     }
193 }
194 ~TSuffTree() {
195     DeleteHelper(root);
196 }
197 };

```


3 Консоль

```
kng@kng-Legion-Y540-15IRH:~/CLionProjects/DISKRAN5/cmake-build-debug ./DISKRAN5
xabay
xabcbay
3
bay
xab
```

4 Тест производительности

Для изучения производительности, сравним время работы программы на тестах разных размеров со строками длиной 10^4 , 10^5 , 10^6

```
1. $10^4$  
/home/kng/CLionProjects/DISKRAN5/cmake-build-debug/DISKRAN5  
251279ms  
Process finished with exit code 0  
  
2. $10^5$  
2791727ms  
Process finished with exit code 0  
  
3. $10^6$  
/home/kng/CLionProjects/DISKRAN5/cmake-build-debug/DISKRAN5  
34305311ms  
Process finished with exit code 0
```

Видно, что рост времени работы практически линейный.

5 Выводы

Во время выполнения работы, я познакомился с такой структурой данных, как суффиксное дерево. Его применение на практике выглядит весьма специфическим, в чем я убедился, посмотрев варианты лабораторных работ. Алгоритм Укконена оказался весьма сложным как для понимания, так и для построения. Сам поиск наибольшей общей подстроки в сравнении с ним оказался простой задачей.

Список литературы