

**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: М. А. Волков  
Преподаватель: А. А. Кухтичев  
Группа: М8О-207Б-19  
Дата: 6 апреля 2021 г.  
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.

# 1 Описание

Требуется реализовать алгоритм Укконена построения суффиксного дерева за линейное время.

Алгоритм Укконена в самой простой реализации имеет сложность  $O(n^3)$ , так как добавляет каждый суффикс каждого префикса строки в отличие от добавления всех суффиксов строки. Основная идея в том, чтобы оптимизировать его и получить сложность  $O(n)$ .

Заметим, что проще будет продлевать суффиксы на один символ. В некоторых случаях при продлении можно идти не по всем символам, а сразу по ребру дерева, что значительно уменьшает число шагов.

Пусть в суффиксном дереве есть строка  $x\alpha$  ( $x$  — первый символ строки,  $\alpha$  — оставшаяся строка), тогда  $\alpha$  тоже будет в суффиксном дереве, потому что  $\alpha$  является суффиксом  $x\alpha$ . Если для строки  $x\alpha$  существует некоторая вершина  $u$ , то существует и вершина  $u$  для  $\alpha$ . Ссылка из  $u$  в  $v$  называется суффиксной ссылкой.

Суффиксные ссылки позволяют не проходить каждый раз по дереву из корня. Для построения суффиксных ссылок достаточно хранить номер последней созданной вершины при продлении. Если на этой же фазе мы создаём ещё одну новую вершину, то нужно построить суффиксную ссылку из предыдущей в текущую.

Сложность алгоритма с использованием продления суффиксов и суффиксных ссылок  $O(n^2)$ . Подробное доказательство изложено в [1].

Для ускорения до  $O(n)$  нужно уменьшить объём потребляемой памяти. Будем в каждом ребре дерева хранить не подстроку, а только индекс начала и конца подстроки.

У исходной строки  $n$  суффиксов и будет создано не более  $n$  внутренних вершин, в среднем продление суффиксов работает за  $O(1)$ .

При использовании всех вышеописанных эвристик получим временную и пространственную сложность  $O(n)$ .

Для нахождения минимального лексикографического разреза строки  $s$  построим суффиксное дерево от удвоенной строки и найдём лексикографически минимальный путь длины  $|s|$  в дереве. Сложность  $O(n)$ .

## 2 Исходный код

Для того, чтобы написать дерево, использовался класс под названием *TBoris*. Главные методы в этом классе тут *Build*, *SetText* и *Colorize*, которые строят суфф дерево, вставляется текст для анализа и поиск ответа соответственно. Метод *Colorize* был так назван, потому что я, опираясь на свой же придуманный метод поиска ответа, красил вершины цветом: красный цвет – принадлежность вершины к 1 слову, синий – ко 2-ому слову.

```
1  #ifndef SUFF_TREE_LAB5_MAIN_H
2  #define SUFF_TREE_LAB5_MAIN_H
3  const char FANTOM_CHAR1 = '%';
4  const char FANTOM_CHAR2 = '$';
5
6  class TBoris {
7      //--variables
8  public:
9      std::shared_ptr<int> end;
10 private:
11     std::string texts;
12     int globID = -1;
13     int one2two = -1;
14
15
16     //-----simply nodes-----
17     struct TBorNode {...};
18
19 protected:
20     TBorNode* root;
21
22     //----functions
23 public:
24     TBoris() {
25         end = std::make_shared<int>(0);
26         root = new TBorNode (root, globID, -1, end);
27         root->url = root;
28     };
29
30     ~TBoris() {
31         delete root;
32     };
33
34     void SetText(const std::string &str1, const std::string &str2) {
35         texts = str1 + FANTOM_CHAR1 + str2 + FANTOM_CHAR2;
36         one2two = str1.size();
37     }
38
39     //-----ITERATOR-----
```

```

40 class TIterator {...};
41
42 void Build() {
43     int activeLen = 0;
44     int edge = 0;
45     TIterator Node(root);
46     bool splitFlag = true;
47
48     for (*end = 0; *end < texts.size(); ++*end) {
49         Node.SetPrevNull();
50         while (globID <= *end) {
51             TBorNode* checkNode = (Node.Next(texts[edge]));
52             while (checkNode != nullptr && activeLen > *checkNode->last - checkNode->begin)
53             {
54                 activeLen = activeLen - 1 - (*checkNode->last - checkNode->begin);
55                 Node = checkNode;
56                 edge = *end - activeLen;
57                 checkNode = Node.Next(texts[edge]);
58             }
59             if (Node.IsPrevNotNull() && activeLen == 0) {
60                 Node.GetPrevNode()->url = Node.GetActiveNode();
61                 Node.SetPrevNull();
62             }
63
64             if (activeLen <= 0 && !Node.FindPath(texts[*end])) {
65                 splitFlag = true;
66
67                 Node.Insert(std::pair<char, TBorNode*>(texts[*end], new TBorNode(root, globID
68                     , *end, end)));
69                 if (Node.IsPrevNotNull()) {
70                     Node.GetPrevNode()->url = Node.GetActiveNode();
71                     Node.SetPrevNull();
72                 }
73                 Node.GoThrowURL();
74                 Node.SetPrevNull();
75                 edge = 0;
76             }
77
78             else if (Node.Next(texts[edge]) != nullptr &&
79                 activeLen > 0 &&
80                 texts[*end] != texts[Node.Next(texts[edge])->begin + activeLen]) {
81                 splitFlag = true;
82                 {
83                     TIterator tmp(Node.Next(texts[edge]));
84                     tmp.Split(root, globID, texts, end, *end, Node.Next(texts[edge])->begin +
85                         activeLen, Node.GetActiveNode());

```

```

86     {
87         TIterator tmp(Node.Next(texts[edge]));
88         if (Node.IsPrevNotNull()) {
89             Node.GetPrevNode()->url = tmp.GetActiveNode();
90             Node.SetPrevNull();
91         }
92
93         if (Node.GetActiveNode() == root) {
94             Node.SetPrev(tmp.GetActiveNode());
95             ++edge;
96             --activeLen;
97         }
98         else {
99             Node.GoThrowURL();
100             Node.SetPrev(tmp.GetActiveNode());
101         }
102     }
103 }
104
105 else {
106     if (splitFlag) {
107         splitFlag = false;
108         edge = *end;
109     }
110     ++activeLen;
111     break;
112 }
113 }
114 }
115 --*end;
116 }
117
118 private:
119 void ColorizeHelp(TIterator &node, std::set<std::string> &ans, int &ansCount, int
    len) {
120     len = len + (*node.GetActiveNode()->last - node.GetActiveNode()->begin) + 1;
121     if (node.GetActiveNode()->leaf) {
122         if (node.GetActiveNode()->begin <= one2two) {
123             node.ColorRed();
124             return;
125         }
126         node.ColorBlue();
127         return;
128     }
129
130     for (const auto &item : node.GetActiveNode()->next) {
131         TIterator nextNode(item.second);
132         ColorizeHelp(nextNode, ans, ansCount, len);
133     }

```

```

134     if (nextNode.GetBlue()) {
135         node.ColorBlue();
136     }
137     if (nextNode.GetRed()) {
138         node.ColorRed();
139     }
140 }
141
142 if (node.BothColored()) {
143     if (len > ansCount) {
144         ans.clear();
145         ansCount = len;
146         ans.insert(node.GetString(texts, *node.GetActiveNode()->last - len + 1, *node.
            GetActiveNode()->last));
147     }
148     else if (len == ansCount) {
149         ans.insert(node.GetString(texts, *node.GetActiveNode()->last - len + 1, *node.
            GetActiveNode()->last));
150     }
151 }
152 }
153
154 public:
155     void Colorize() {
156         std::set<std::string> ans;
157         int ansCount = 0;
158         TIterator node(root);
159         for (const auto &item : node.GetActiveNode()->next) {
160             TIterator nextNode(item.second);
161             ColorizeHelp(nextNode, ans, ansCount, 0);
162         }
163
164         std::cout << ansCount << std::endl;
165         if (!ans.empty()) {
166             for (const auto &item : ans) {
167                 std::cout << item << std::endl;
168             }
169         }
170     }
171
172 };
173
174 #endif //SUFF_TREE_LAB5_MAIN_H

```

Для удобства передвижения по суффиксному дереву, мной был придуман итератор с нужными для меремещения методами такими, как *Insert* и *Split*. Данная структура данных мне помогла локально изменять свой код, при этом изменяя глобальную логику моей работы. Например я обнаружил ошибку в разделении вершины. Мне нужно было поменять метод *Split* и все остальное работает хорошо.

```
1  class TIterator {
2  private:
3      TBorNode* prevNode;
4      TBorNode* activeNode;
5
6  public:
7      TIterator() {
8          activeNode = nullptr;
9          prevNode = nullptr;
10     }
11
12     explicit TIterator(TBorNode* node) {
13         activeNode = node;
14         prevNode = nullptr;
15     }
16
17     ~TIterator() {
18         activeNode = nullptr;
19         prevNode = nullptr;
20     };
21
22     TBorNode* Next(const char &c) {
23         if (activeNode == nullptr) {
24             return nullptr;
25         }
26
27         auto tmp = activeNode->GetNodeElem(c);
28         if (tmp != nullptr) {
29             return tmp;
30         }
31         return nullptr;
32     }
33
34     bool FindPath(const char &c) {
35         if (activeNode == nullptr) {
36             return false;
37         }
38
39         auto tmp = activeNode->GetNodeElem(c);
40         if (tmp != nullptr) {
41             return true;
42         }
43         return false;
44     }
```



```

44     }
45
46     void GoThrowURL() {
47         if (activeNode->url == nullptr) {
48             throw std::logic_error("Trying to call nullptr in URL");
49         }
50         prevNode = activeNode;
51         activeNode = activeNode->url;
52     }
53
54     void Insert(const std::pair<char, TBorNode*> &tmp_pair) {
55         activeNode->Insert(tmp_pair);
56     }
57
58     void Split(TBorNode* _root, int &_globID, const std::string &_texts,
59             const std::shared_ptr<int> &_end,
60             const int &begin, const int &splitter, TBorNode* Node) {
61         Node->next.erase(_texts[activeNode->begin]);
62         Node->Insert(std::pair<char, TBorNode*>(_texts[activeNode->begin],
63             new TBorNode (_root, _globID,
64                                     activeNode
65                                     ->
66                                     begin
67                                     ,
68                                     splitter
69                                     -
70                                     1)
71                                     ));
72
73         TBorNode* tmp = Node->GetNodeElem(_texts[activeNode->begin]);
74         --_globID;
75         tmp->Insert(std::pair<char, TBorNode*>(_texts[begin],
76             new TBorNode(_root, _globID,
77                 begin,
78                 _end)))
79         ;
80
81         activeNode->begin = splitter;
82         tmp->Insert(std::pair<char, TBorNode*>(_texts[activeNode->begin],
83             activeNode));
84     }
85 };

```

Также была сделана структурка небольшая, в которую я помещаю всю информацию о вершине.

```
1  struct TBorNode {
2      //variables
3      int begin;
4      std::shared_ptr<int> last;
5      bool leaf = true;
6      TBorNode* url;
7      std::unordered_map<char, TBorNode*> next;
8      int id;
9      bool red = false;
10     bool blue = false;
11
12     //constructors
13     TBorNode(TBorNode* root, int &globID, const int &_begin,
14             const std::shared_ptr<int> &end) {
15         id = globID++;
16         last = end;
17         begin = _begin;
18         url = root;
19     }
20
21     TBorNode(TBorNode* root, int &globID, const int &_begin, const int &end) {
22         id = globID++;
23         last.reset(new int(end));
24         begin = _begin;
25         url = root;
26         leaf = false;
27     }
28
29     TBorNode() = delete;
30
31     //destructors
32     ~TBorNode(){
33         url = nullptr;
34         for (auto item : next){
35             delete item.second;
36         }
37     };
38
39     void PrintNode(const std::string &text) {
40         for (int i = begin; i <= *last; ++i) {
41             std::cout << text[i];
42         }
43     }
44
45     void Insert(const std::pair<char, TBorNode*> &tmp_pair) {
46         next.insert(tmp_pair);
47         leaf = false;
```

```

48     }
49
50     TBorNode* GetNodeElem(const char &c) {
51         if (!next.empty()) {
52             if (next.find(c) == next.end()) {
53                 return nullptr;
54             }
55             return next.find(c)->second;
56         }
57         return nullptr;
58     }
59
60 };

```

### 3 Тест производительности

Сравним наинный алгоритм, который сортирует все циклические разрезы строки, с алгоритмом, использующим суффиксное дерево.

Тесты состоят из строк длины 100, 1000, 10000.

test1.txt

Naive time 0.1 ms

My time 0.2 ms

test2.txt

Naive time 0.54 ms

My time 0.20 ms

test3.txt

Naive time 3.484 ms

My time 0.171 ms

Видно, что на маленькой строке константа в сложности алгоритма Укконена даёт о себе знать, а на больших строках квадратичная сложность не может соревноваться с линейной.

## 4 Выводы

Во время выполнения данной лабораторной работы, я, проверяя свою программу на похожем задании на сайте *codeforces*, понял, что алгоритм Укконена работает очень быстро. Данный алгоритм даже не занимает много места в оперативной памяти при непосредственном выполнении.

К сожалению написание данного алгоритма достаточно сложное и поэтому большинство программистов выбирают алгоритм КМП. Или же, решая мою задачу, используют динамический метод программирования.

## Список литературы

- [1] *Алгоритм Укконена — Викиконспекты*  
URL: [https://neerc.ifmo.ru/wiki/index.php?title=Алгоритм\\_Укконена](https://neerc.ifmo.ru/wiki/index.php?title=Алгоритм_Укконена) (дата обращения: 19.03.2021).
- [2] *Visualization of Ukkonen's Algorithm*  
URL: <http://brenden.github.io/ukkonen-animation/> (дата обращения: 18.03.2021).