

**Московский авиационный институт
(национальный исследовательский университет)**

**Институт информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

Лабораторная работа №5 по курсу «Дискретный анализ»

Студент: Е. А. Кондратьев
Преподаватель: А. А. Кухтичев
Группа: М8О-206Б-19
Дата:
Оценка:
Подпись:

Москва, 2021

Лабораторная работа №5

Задача: Вариант №5

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 написать реализацию алгоритма Укконена для построения суффиксного дерева, затем с помощью него построить суффиксный массив, посредством которого будет происходить поиск образцов в тексте.

Основная идея алгоритма Укконена заключается в том, чтобы построить суффиксное дерево за линейное время.

Для этого нужно, во-первых, использовать линейное количество памяти, поэтому в рёбрах (вершинах) будем хранить два числа — позиции самого левого и самого правого символов в исходном тексте.

Во-вторых, по мере создания дерева будем создавать и использовать суффиксные ссылки, основой для которых является факт, если какая-либо строка уже была добавлена в дерево, то и все ее суффиксы тоже присутствуют в дереве. То есть суффиксной ссылкой будет являться ссылка от вершины $ха$ в вершину $а$, где $х$ — первым символом строки, $а$ — оставшаяся подстрока (возможно пустая). Использование суффиксных ссылок позволяет пропускать ненужные сравнения.

Суффиксный массив для строки s в i -ой позиции хранит индекс начала i -го (в лексикографическом порядке) суффикса строки s . Будем строить данный массив путём обхода в глубину суффиксного дерева в лексикографическом порядке.

Сложность $O(t)$, где t — длина текста. Поиск образцов в суффиксном массиве будем осуществлять бинарным поиском (поскольку суффиксы отсортированы лексикографически). Сложность $O(p \cdot \log(t))$, где t , p — длины текста и образца соответственно.

Таким образом, за $O(t)$ времени и памяти мы строим суффиксное дерево с помощью алгоритма Укконена, затем за $O(t)$ создаём суффиксный массив, в котором за суммарное время $O(\log(t) \cdot (p_1 + p_2 + \dots + p_m))$ ищем m образцов. Итоговая сложность всех операций $O(t + \log(t) \cdot (p_1 + p_2 + \dots + p_m))$.

Код

В `suffix_tree.h` опишем классы суффиксного дерева и его вершин. В вершине будем хранить итераторы на начала и конец ребра, суффиксную ссылку, указатели на следующие вершины будем хранить с `std::map` с ключом в виде первого символа строки.

```
1 namespace NSuffixTrees {
2
3     class TSuffixArray;
4
5     class TNode {
6     public:
7         TNode(std::string::iterator begin, std::string::iterator end);
8         ~TNode() = default;
9
10        std::string::iterator begin;
11        std::string::iterator end;
12
13        std::map< char, TNode* > to;
14        TNode* suffixLink;
15    };
16
17    class TSuffixTree {
18    public:
19        TSuffixTree(std::string str);
20        ~TSuffixTree();
21        friend TSuffixArray;
22
23    private:
24        std::string text;
25        TNode* root;
26
27        TNode* needSufLink;
28        TNode* activeNode;
29        int remainder;
30        int activeLen;
31        std::string::iterator activeEdge;
32
33        void TreeExtend(std::string::iterator toAdd);
34        void DeleteTree(TNode* node);
35
36        int GetEdgeLen(TNode* node, std::string::iterator pos) const;
37        void AddSuffixLink(TNode* node);
38        void DFS(TNode* node, std::vector<size_t>& result, size_t depth) const;
39    };
40 }
```

В `suffix_tree.cpp` напомним реализации методов классов `TNode` и `TSuffixTree`. `TreeExtend` представляет собой шаг алгоритма Укконена по добавления очередной буквы `toAdd` в дерево. `AddSuffixLink` используется для создания суффиксных ссылок. `DFS` используется для обхода в глубину и заполнения суффиксного массива `result`.

```
1 namespace NSuffixTrees {
2     /// TNode
3     TNode::TNode(std::string::iterator begin, std::string::iterator end) :
4         begin(begin),
5         end(end),
6         suffixLink(0)
7     {}
8 }
```

```

8
9  //// TSuffixTree
10 TSuffixTree::TSuffixTree(std::string str) :
11     text(str),
12     root(new TNode(text.end(), text.end())),
13     remainder(0)
14 {
15     activeEdge = text.begin();
16     activeNode = needSufLink = root->suffixLink = root;
17     activeLen = 0;
18
19     for (std::string::iterator it = text.begin(); it != text.end(); ++it) {
20         TreeExtend(it);
21     }
22 }
23
24 TSuffixTree::~TSuffixTree() {
25     DeleteTree(root);
26 };
27
28 void TSuffixTree::DeleteTree(TNode* node) {
29     for (std::map< char, TNode* >::iterator it = node->to.begin(); it != node->to.end()
30         ; ++it) {
31         DeleteTree(it->second);
32     }
33     delete node;
34 }
35
36 int TSuffixTree::GetEdgeLen(TNode* node, std::string::iterator pos) const {
37     return std::min(node->end, pos + 1) - node->begin;
38 }
39
40 void TSuffixTree::TreeExtend(std::string::iterator toAdd) {
41     needSufLink = root;
42     ++remainder;
43
44     while (remainder) {
45         if (!activeLen) activeEdge = toAdd;
46
47         TNode *next = NULL;
48         std::map< char, TNode* >::iterator it = activeNode->to.find(*activeEdge);
49         if (it != activeNode->to.end()) next = it->second;
50
51         if (!next) {
52             TNode* leaf = new TNode(toAdd, text.end());
53             activeNode->to[*activeEdge] = leaf;
54             AddSuffixLink(activeNode);
55         } else {
56             if (activeLen >= GetEdgeLen(next, toAdd)) {
57                 activeEdge += GetEdgeLen(next, toAdd);
58                 activeLen -= GetEdgeLen(next, toAdd);
59                 activeNode = next;
60                 continue;
61             }
62
63             if (*(next->begin + activeLen) == *toAdd) {
64                 ++activeLen;
65                 AddSuffixLink(activeNode);
66                 break;
67             }

```

```

68     TNode* split = new TNode(next->begin, next->begin + activeLen);
69     TNode* leaf = new TNode(toAdd, text.end());
70     activeNode->to[*activeEdge] = split;
71
72     split->to[*toAdd] = leaf;
73     next->begin += activeLen;
74     split->to[*next->begin] = next;
75     AddSuffixLink(split);
76 }
77
78 --remainder;
79 if (activeNode == root && activeLen > 0) {
80     --activeLen;
81     activeEdge = toAdd - remainder + 1;
82 } else {
83     if (activeNode->suffixLink) {
84         activeNode = activeNode->suffixLink;
85     } else {
86         activeNode = root;
87     }
88 }
89 }
90 }
91
92 void TSuffixTree::AddSuffixLink(TNode* node) {
93     if (needSufLink != root) needSufLink->suffixLink = node;
94     needSufLink = node;
95 }
96
97 void TSuffixTree::DFS(TNode* node, std::vector<size_t>& result, size_t depth) const
98 {
99     if (node->to.empty()) {
100         result.push_back(text.size() - depth);
101         return;
102     }
103     for (std::map<char, TNode*>::iterator it = node->to.begin(); it != node->to.end();
104         ++it) {
105         DFS(it->second, result, depth + it->second->end - it->second->begin);
106     }
107 }

```

suffix_array.h опишем класс суффиксного массива с конструктором от суффиксного дерева.

```

1 namespace NSuffixTrees {
2     class TSuffixTree;
3
4     class TSuffixArray {
5     public:
6         TSuffixArray(const TSuffixTree& tree);
7         ~TSuffixArray() = default;
8
9         std::vector<size_t> Find(const std::string& pattern);
10    private:
11        std::string text;
12        std::vector<size_t> array;
13    };
14 }

```

В suffix_array.cpp напишем соответствующие реализации.

```
1 namespace NSuffixTrees {
2   TSuffixArray::TSuffixArray(const TSuffixTree& tree) :
3     text(tree.text),
4     array()
5   {
6     tree.DFS(tree.root, array, 0);
7   }
8
9   std::vector<size_t> TSuffixArray::Find(const std::string& pattern) {
10     std::pair<std::vector<size_t>::iterator, std::vector<size_t>::iterator> range(array
11       .begin(), array.end());
12     for (size_t i = 0; i < pattern.size() && range.first != range.second; ++i) {
13       range = equal_range(range.first, range.second, std::numeric_limits<size_t>::max()
14         , [this, &pattern, &i] (size_t index1, size_t index2) -> bool {
15         if (index1 == std::numeric_limits<size_t>::max()) {
16           return pattern[i] < text[i + index2];
17         } else {
18           return text[i + index1] < pattern[i];
19         }
20       });
21     }
22     std::vector<size_t> result(range.first, range.second);
23     std::sort(result.begin(), result.end());
24     return result;
25   }
26 }
```

2 Консоль

```
du@Du$ ./wrapper.sh [2021-06-17 15:53:13] [INFO] Compiling... g++ -std=c++17 -O3
-pedantic -Wall -Wno-unused-variable -c main.cpp -o main.o g++ -std=c++17 -O3 -
pedantic -Wall -Wno-unused-variable -c suffix_tree.cpp -o suffix_tree.o g++ -std=c++17
-O3 -pedantic -Wall -Wno-unused-variable -c suffix_array.cpp -o suffix_array.o g++ -
std=c++17 -O3 -pedantic -Wall -Wno-unused-variable main.o suffix_tree.o suffix_array.o
-o solution [2021-06-17 15:53:15] [INFO] Executing tests/t01.t... OK [2021-06-17 15:53:15]
[INFO] Executing tests/t02.t... OK [2021-06-17 15:53:15] [INFO] Executing tests/t03.t...
OK [2021-06-17 15:53:15] [INFO] No failed tests, hooray du@Du$ cat tests/t02.t wowCanYouFindMeHere
Here You Can du@Du$ ./solution <tests/t02.t 1: 16,22,29 2: 7,36 3: 4,33
```


3 Выводы

Выполнив пятую лабораторную работу по курсу «Дискретный анализ», я познакомился с суффиксными деревьями и суффиксными массивами, а также как работать с ними. Пригодился опыт прошлых лабораторных работ с сбалансированными деревьями: перемещение по дереву и создание новых вершин не представляло трудности. В процессе отладки программы на тестах вскрылись многие баги и ошибки, путём исправления которых получилось доработать изначальный алгоритм.

Список литературы

- [1] *Поисковик - Google.*
URL: <https://www.google.com/>
- [2] *Сайт с подробной документацией библиотек C++*
URL: <https://en.cppreference.com/>
- [3] *Алгоритм Укконена — ИТМО Вики.*
URL: <https://neerc.ifmo.ru/>