

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. А. Кабанов
Преподаватель: С. А. Михайлова
Группа: М8О-301Б-22
Дата:
Оценка:
Подпись:

Москва, 2025

Лабораторная работа №5

Задача: Поиск в известном тексте неизвестных заранее образцов.

Найти в заранее известном тексте поступающие на вход образцы.

Формат ввода: Текст располагается на первой строке, затем, до конца файла, следуют строки с образцами.

Формат вывода: Для каждого образца, найденного в тексте, нужно распечатать строчку, начинающуюся с последовательного номера этого образца и двоеточия, за которым, через запятую, нужно перечислить номера позиций, где встречается образец в порядке возрастания.

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, α - оставшаяся строка), тогда α тоже будет в суффиксном дереве, потому что α является суффиксом $x\alpha$. Если для строки $x\alpha$ существует некоторая вершина u , то существует и вершина u для α . 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 $O(n^2)$.

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

Реализуем суффиксное дерево и алгоритм, описанные в предыдущем пункте. Были произведены некоторые изменения в коде, связанные с реализацией на C++: к примеру, вместо хранения индексов строки хранятся итераторы.

```
1 #include <bits/stdc++.h>
2
3 struct Node;
4 struct STree;
5
6 struct Node {
7     std::map<char, Node *> to;
8     std::string::iterator begin;
9     std::string::iterator end;
10    Node *suffixLink;
11
12    Node(std::string::iterator begin, std::string::iterator end)
13        : begin(begin), end(end), suffixLink(nullptr) {}
14
15    ~Node() {};
16 };
17
18 struct STree {
19     std::string text;
20     Node *root;
21     Node *activeNode;
22     Node *lastAdded;
23     std::string::iterator activeEdge;
24     int32_t activeLength;
25     int32_t remainder;
26
27     STree() = default;
28     STree(std::string &str)
29         : text(str), activeLength(0), activeEdge(text.begin()), remainder(0) {
30         root = new Node(text.end(), text.end());
31         lastAdded = root;
32         activeNode = root;
33         root->suffixLink = root;
34
35         for (std::string::iterator it = text.begin(); it != text.end(); ++it) {
36             addLetter(it);
37         }
38     }
39
40     ~STree() { destroy(root); }
41
42     void addLetter(std::string::iterator i) {
43         lastAdded = root;
```

```

44 ++remainder;
45 while (remainder > 0) {
46     activeEdge = (activeLength == 0) ? i : activeEdge;
47     std::map<char, Node *>::iterator it =
48     activeNode->to.find(*activeEdge);
49     Node *next;
50     if (it == activeNode->to.end()) {
51         Node *leaf = new Node(i, text.end());
52         activeNode->to[*activeEdge] = leaf;
53         lastAdded->suffixLink =
54         (lastAdded != root) ? activeNode : lastAdded->suffixLink;
55         lastAdded = activeNode;
56     } else {
57         next = it->second;
58         if (checkEdge(i, next)) {
59             continue;
60         }
61         if (*(next->begin + activeLength) == *i) {
62             ++activeLength;
63             lastAdded->suffixLink =
64             (lastAdded != root) ? activeNode : lastAdded->suffixLink;
65             lastAdded = activeNode;
66             break;
67         }
68         Node *split = new Node(next->begin, next->begin + activeLength);
69         Node *leaf = new Node(i, text.end());
70         activeNode->to[*activeEdge] = split;
71         split->to[*i] = leaf;
72         next->begin += activeLength;
73         split->to[*next->begin] = next;
74         lastAdded->suffixLink =
75         (lastAdded != root) ? split : lastAdded->suffixLink;
76         lastAdded = split;
77     }
78     --remainder;
79     if (activeNode == root && activeLength > 0) {
80         --activeLength;
81         activeEdge = i - remainder + 1;
82     } else {
83         activeNode = activeNode->suffixLink;
84     }
85 }
86 }
87
88 void searchLeaves(Node *node, std::vector<int32_t> &answer, int32_t loc) {
89     if (node->end == text.end()) {
90         answer.push_back(text.size() - loc + 1);
91     } else {
92         Node *child;

```

```

93     for (auto it = node->to.begin(); it != node->to.end(); ++it) {
94         child = it->second;
95         searchLeaves(child, answer, loc + child->end - child->begin);
96     }
97 }
98 }
99
100 std::vector<int32_t> search(std::string &str) {
101     std::vector<int32_t> answer;
102     int32_t loc = 0;
103     Node *current = root;
104     if (str.length() > text.length()) {
105         return answer;
106     }
107     for (std::string::iterator pos = str.begin(); pos != str.end(); ++pos) {
108         auto path = current->to.find(*pos);
109         if (path == current->to.end()) {
110             return answer;
111         }
112         current = path->second;
113         loc += current->end - current->begin;
114         for (std::string::iterator p = current->begin;
115              p != current->end && pos != str.end(); ++p, pos++) {
116             if (*p != *pos) {
117                 return answer;
118             }
119         }
120         if (pos == str.end()) {
121             break;
122         }
123         --pos;
124     }
125     searchLeaves(current, answer, loc);
126     sort(answer.begin(), answer.end());
127     return answer;
128 }
129
130 void destroy(Node *node) {
131     for (auto it = node->to.begin(); it != node->to.end(); ++it) {
132         destroy(it->second);
133     }
134     delete node;
135 }
136
137 bool checkEdge(std::string::iterator pos, Node *node) {
138     int32_t edgeLength = (pos + 1 < node->end) ? (pos + 1 - node->begin)
139     : (node->end - node->begin);
140     if (activeLength >= edgeLength) {
141         activeEdge += edgeLength;

```

```

142     activeLength -= edgeLength;
143     activeNode = node;
144     return true;
145 }
146 return false;
147 }
148 };
149
150 int main(int argc, char *argv[]) {
151     std::ios_base::sync_with_stdio(false);
152     std::cin.tie(nullptr);
153
154     std::string text;
155     std::cin >> text;
156     text += "$";
157
158     STree tree(text);
159     int32_t k{0};
160     std::vector<int32_t> answer;
161     std::string pattern;
162     while (std::cin >> pattern) {
163         if (pattern.empty()) {
164             break;
165         }
166         ++k;
167         answer = tree.search(pattern);
168         if (!answer.empty()) {
169             std::cout << k << ": ";
170             for (std::size_t i = 0; i < answer.size(); ++i) {
171                 std::cout << answer[i] << ((i == answer.size() - 1) ? "\n" : ", ");
172             }
173         }
174     }
175     return 0;
176 }

```

3 Консоль

```
[anton@home lab5-suffix-tree]$ make
g++ -O2 -lm -fno-stack-limit -std=c++20 -x c++ solution.cpp -o executable
[anton@home lab5-suffix-tree]$ ./executable
abcdabc
abcd
bcd
bc
1: 1
2: 2
3: 2,6
```


4 Тест производительности

Тест производительности представляет из себя следующее: сравнение написанного мной алгоритма за $O(n)$ и наивного алгоритма за $O(n^3)$.

```
[anton@home lab7-dp]$ make
g++ main.cpp -o main
g++ naive.cpp -o naive

...

./main <tests/1.in
Suffix tree time 0.178 ms
./naive <tests/1.in
Naive time 0.116 ms
./main <tests/2.in
Suffix tree time 1.680 ms
./naive <tests/2.in
Naive time 3.180 ms
./main <tests/3.in
Suffix tree time 17.795 ms
./naive <tests/3.in
Naive time 81.510 ms
./main <tests/4.in
Suffix tree time 18.455 ms
./naive <tests/4.in
Naive time 1334.365 ms
```

Видно, что построение суффиксного дерева по алгоритму Укконена и дальнейший поиск по нему выигрывают у наивного алгоритма.

5 Выводы

Выполнив седьмую лабораторную работу по курсу «Дискретный анализ», я вспомнил структуры данных, связанные со строками, и реализовал алгоритм Укконена, ознакомился с приложениями суффиксного дерева.

Суффиксное дерево позволяет быстро искать множество шаблонов в тексте, чего не могут другие алгоритмы. Но в повседневных задачах чаще требуется найти один шаблон в тексте, где лучше использовать более простые алгоритмы.

Список литературы

- [1] *Алгоритм Укконена - Викиконспекты*
URL: https://neerc.ifmo.ru/wiki/index.php?title=Алгоритм_Укконена (дата обращения: 10.12.2024).
- [2] *Visualization of Ukkonen's Algorithm*
URL: <http://brenden.github.io/ukkonen-animation/> (дата обращения: 10.12.2024).