

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. Терво
Преподаватель: А. А. Кухтичев
Группа: М8О-207Б
Дата:
Оценка:
Подпись:

Москва, 2022

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 задания: найти самую длинную общую подстроку двух строк.

1 Описание

Требуется реализовать алгоритм Укконена построения суффиксного дерева за линейное время.

В простейшей реализации метод построения суффиксного дерева работает за $O(n^3)$. В нём каждый суффикс продлевается независимо от остальных. Алгоритм Укконена является усовершенствованием данного метода, и работает за линейное время — $O(n)$.

В работе алгоритм использует суффиксные ссылки. Говоря простым языком, суффиксная ссылка — это ссылка из вершины с путевой меткой $x\alpha$ в вершину с меткой α , то есть в вершину с такой же меткой, но без первого символа. Более корректное определение приведено в [1]. Суффиксные ссылки нужны для того, чтобы быстро перейти в следующий суффикс и выполнить в нём необходимые операции.

Для ускорения продления суффиксов и уменьшения расхода памяти в каждом узле хранятся индекс начала и конца подстроки, причём конец подстроки хранится в виде указателя на переменную, в которой он записан. Это сделано для того, чтобы не продлевать каждый суффикс в отдельности, а добавлять по одному символу сразу ко всем.

Для нахождения максимальной общей подстроки двух строк A и B построим дерево для строки $A\$B*$, где $\$$ и $*$ — знаки, не входящие в алфавиты обеих строк. Затем найдём самую глубокую вершину, в одном поддереве которой заканчивается строка A , а в другом — строка B . Путевая метка этой вершины и будет максимальной общей подстрокой.

2 Исходный код

Опишем класс дерева и вспомогательную структуру TEdge для хранения информации о ребре. Нам понадобятся конструктор от строки и методы нахождения наиболее длинной общей подстроки.

```
1  class TSuffixTree {
2      private:
3          struct TEdge {
4              int Left;
5              std::shared_ptr<int> Right;
6              int IdTo;
7
8              TEdge(int l, std::shared_ptr<int> r, int id) : Left(l), Right(r), IdTo(id) {}
9          };
10
11          std::string DataString;
12          std::vector< std::vector< TSuffixTree::TEdge > > Data;
13          std::vector<int> SuffixPtr;
14          std::vector<int> PathSize;
15      public:
16          TSuffixTree(const std::string & s);
17          void Print();
18          void Find(const char, const char);
19          int RecMarker(const char, const char, const int, const int, int &, std::vector<
20              int> &);
21          void RecCatchStrings(const int, const int, const int, const std::string, std::
22              vector<int> &, std::vector<std::string> &);
23      };
24 }
```

В таком случае main.cpp получается достаточно небольшой.

```
1  #include "suffix_tree.hpp"
2
3  const char SENTINEL1 = '$';
4  const char SENTINEL2 = '*';
5
6  int main() {
7      std::string s;
8      std::cin >> s;
9      std::string t = s + SENTINEL1;
10     std::cin >> s;
11     t = t + s + SENTINEL2;
12     TSuffixTree st(t);
13     st.Find(SENTINEL1, SENTINEL2);
14     return 0;
15 }
```

Самым сложным является реализация конструктора дерева от строки, являющаяся

собственно реализацией алгоритма Укконена. Дерево имеет максимум $2 * n$ вершин и $2 * n - 1$ рёбер, поэтому его удобно хранить списками смежности.

```

1  TSuffixTree::TSuffixTree(const std::string & s) : DataString(s), Data(s.size() * 2),
    SuffixPtr(s.size() * 2), PathSize(s.size() * 2) {
2      bool newVertex = false;
3      std::shared_ptr<int> end(new int);
4      *end = 0;
5      int curVertex = 0;
6      int l = 0;
7      int vertexId = 0;
8      int curEq = 0;
9      int passEq = 0;
10     int curEdgeId = -1;
11     for (size_t i = 0; i < s.size(); ++i) {
12         ++*end;
13         int lastCreatedVertexId = 0;
14         int nextCreatedVertexId = -1;
15         while (l < *end) {
16             if (curEdgeId == -1) {
17                 int nextEdgeId = -1;
18                 for (size_t j = 0; j < Data[curVertex].size(); ++j) {
19                     if (DataString[Data[curVertex][j].Left] == DataString[l + passEq + curEq])
20                         {
21                             nextEdgeId = j;
22                             curEq = 1;
23                             break;
24                         }
25                     if (nextEdgeId == -1) {
26                         newVertex = true;
27                     }
28                     curEdgeId = nextEdgeId;
29                 } else {
30                     if (DataString[Data[curVertex][curEdgeId].Left + curEq] == DataString[l +
31                         curEq + passEq]) {
32                         ++curEq;
33                     } else {
34                         newVertex = true;
35                     }
36                 }
37                 if (curEq > 0 and Data[curVertex][curEdgeId].Left + curEq == *Data[curVertex][
38                     curEdgeId].Right) {
39                     curVertex = Data[curVertex][curEdgeId].IdTo;
40                     curEdgeId = -1;
41                     passEq = passEq + curEq;
42                     curEq = 0;
43                 }
44                 if (newVertex) {
45                     ++vertexId;

```

```

44     if (curEq == 0 and curEdgeId == -1) {
45         Data[curVertex].push_back(TSuffixTree::TEdge(1 + passEq, end, vertexId));
46     } else {
47         nextCreatedVertexId = vertexId;
48         TEdge curEdge = Data[curVertex][curEdgeId];
49         std::shared_ptr<int> newRightBorder(new int);
50         *newRightBorder = Data[curVertex][curEdgeId].Left + curEq;
51         Data[curVertex][curEdgeId].Right = newRightBorder;
52         Data[curVertex][curEdgeId].IdTo = vertexId;
53         curEdge.Left = *newRightBorder;
54
55         Data[vertexId].push_back(curEdge);
56         ++vertexId;
57         Data[vertexId - 1].push_back(TSuffixTree::TEdge(1 + curEq + passEq, end,
58             vertexId));
59         PathSize[vertexId - 1] = PathSize[curVertex] + *Data[curVertex][curEdgeId].
60             Right - Data[curVertex][curEdgeId].Left;
61
62         if (lastCreatedVertexId > 0) {
63             SuffixPtr[lastCreatedVertexId] = nextCreatedVertexId;
64         }
65         lastCreatedVertexId = nextCreatedVertexId;
66     }
67     int nextCurVertex = SuffixPtr[curVertex];
68     int nextPassEq = PathSize[nextCurVertex];
69     int nextCurEq = curEq + passEq - nextPassEq - 1;
70     int nextEdgeId = -1;
71     while (nextCurEq > 0) {
72         for (size_t j = 0; j < Data[nextCurVertex].size(); ++j) {
73             if (DataString[Data[nextCurVertex][j].Left] == DataString[1 + 1 +
74                 nextPassEq]) {
75                 nextEdgeId = j;
76                 break;
77             }
78         }
79         int curRight = *Data[nextCurVertex][nextEdgeId].Right;
80         int curLeft = Data[nextCurVertex][nextEdgeId].Left;
81         if (nextEdgeId != -1 and curRight - curLeft <= nextCurEq) {
82             nextPassEq = nextPassEq + curRight - curLeft;
83             nextCurEq = nextCurEq - curRight + curLeft;
84             nextCurVertex = Data[nextCurVertex][nextEdgeId].IdTo;
85             nextEdgeId = -1;
86         } else {
87             break;
88         }
89     }
90     if (nextEdgeId != -1) {
91         curEq = nextCurEq;
92     } else {

```

```

90         curEq = 0;
91     }
92     curEdgeId = nextEdgeId;
93     curVertex = nextCurVertex;
94     passEq = nextPassEq;
95     ++l;
96     newVertex = false;
97 } else {
98     if (i < s.size() - 1) {
99         break;
100     }
101 }
102 }
103 }
104 }

```

3 Консоль

```
[alext@alext-pc solution]$ make
g++ -std=c++17 -pedantic -g -Wall -Wextra -Wno-unused-variable main.cpp -o
solution
[alext@alext-pc solution]$ ./solution
abcabcabc
abcax
4
abca
[alext@alext-pc solution]$ ./solution
abracadabra
xyzabredabuey
3
abr
dab
[alext@alext-pc solution]$ ./solution
aaaaabaabaaa
aa
2
aa
[alext@alext-pc solution]$ ./solution
abcdefg
abcdefg
7
abcdefg
```


4 Тест производительности

Сравним наивный алгоритм поиска наибольшей общей подстроки с поиском подстроки с помощью суффиксного дерева.

```
[alex@alex-pc solution]$ ./naive <tests/test10.txt
Naive algorithm time: 2us
[alex@alex-pc solution]$ ./naive <tests/test100.txt
Naive algorithm time: 82us
[alex@alex-pc solution]$ ./naive <tests/test1k.txt
Naive algorithm time: 6346us
[alex@alex-pc solution]$ ./naive <tests/test10k.txt
Naive algorithm time: 635613us
[alex@alex-pc solution]$ ./solution <tests/test10.txt
Ukkonen algorithm time: 33us
[alex@alex-pc solution]$ ./solution <tests/test100.txt
Ukkonen algorithm time: 46us
[alex@alex-pc solution]$ ./solution <tests/test1k.txt
Ukkonen algorithm time: 303us
[alex@alex-pc solution]$ ./solution <tests/test10k.txt
Ukkonen algorithm time: 2772us
```

Видно, что на маленьких строках наивный алгоритм выигрывает за счёт простоты построения небольшой матрицы и вычисления на ней. Но уже на строках длиной 100 видно, что он начинает проигрывать по времени алгоритму с использованием суффиксного дерева.

Также видно, что сложность наивного алгоритма больше, чем у алгоритма с суффиксным деревом: при увеличении размера входных данных на порядок время наивного алгоритма растёт на два порядка, а время суффиксного — только на один.

5 Выводы

Во время выполнения этой лабораторной работы я вспомнил определение суффиксного дерева. Изучил наивный алгоритм его построения и алгоритм Укконена.

Этот алгоритм работает достаточно быстро и эффективен для поиска множества шаблонов в строке. Но он имеет и минусы, в числе которых его сложность для понимания.

Для поиска одного шаблона в строке удобнее использовать другие алгоритмы, такие как КМП, алгоритм Бойера-Мура или Рабина-Карпа. Для решения же моей задачи — поиска наибольшей общей подстроки проще использовать метод динамического программирования.

Список литературы

- [1] *Алгоритм Укконена — Викиконспекты*
URL: https://neerc.ifmo.ru/wiki/index.php?title=Алгоритм_Укконена (дата обращения: 16.12.2021).