

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. Инютин
Преподаватель: А. А. Кухтичев
Группа: М8О-207Б-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 а до z).

Вариант: линейризовать циклическую строку, то есть найти минимальный в лексикографическом смысле разрез циклической строки.

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 $O(n^2)$. 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

Опишем класс дерева и вспомогательную структуру TEdge для хранения информации о ребре. Нам понадобится конструктор от строки и методы нахождения минимальной лексикографической подстроки заданной длины.

```
1 | #ifndef SUFFIX_TREE_HPP
2 | #define SUFFIX_TREE_HPP
3 |
4 | #include <memory>
5 | #include <string>
6 | #include <vector>
7 |
8 | class TSuffixTree {
9 | private:
10 |     struct TEdge {
11 |         int Left;
12 |         std::shared_ptr<int> Right;
13 |         int IdTo;
14 |
15 |         TEdge(int l, std::shared_ptr<int> r, int id) : Left(l), Right(r), IdTo(id) {}
16 |     };
17 |
18 |     std::string DataString;
19 |     std::vector< std::vector< TSuffixTree::TEdge > > Data;
20 |     std::vector<int> SuffixPtr;
21 |     std::vector<int> PathSize;
22 | public:
23 |     TSuffixTree(const std::string & s);
24 |     std::string LexMinString(const size_t n);
25 |     std::string LexMinString(const int id, const size_t n);
26 | };
27 |
28 | #endif /* SUFFIX_TREE_HPP */
```

При такой реализации решения занимает несколько строк.

```
1 | #include "suffix_tree.hpp"
2 | #include <iostream>
3 |
4 | const char SENTINEL = 'z' + 1;
5 |
6 | int main() {
7 |     std::string s;
8 |     std::cin >> s;
9 |     std::string t = s + s + SENTINEL;
10 |     TSuffixTree st(t);
11 |     std::cout << st.LexMinString(s.size()) << '\n';
12 |     return 0;
13 | }
```

Самое сложное — реализация конструктора дерева и поиска лексикографически минимального пути в нём. Суффиксное дерево имеет максимум $2 * n$ вершин и $2 * n - 1$ рёбер, поэтому достаточно выгодно хранить дерево списками смежности.

```

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

```

```

43     }
44     if (newVertex) {
45         ++createdVertexId;
46         if (curEq == 0 and curEdgeId == -1) {
47             Data[curVertexId].push_back(TSuffixTree::TEdge(1 + passEq, end,
48                 createdVertexId));
49         } else {
50             nextCreatedVertexId = createdVertexId;
51             TEdge curEdge = Data[curVertexId][curEdgeId];
52
53             std::shared_ptr<int> newRightBorder(new int);
54             *newRightBorder = Data[curVertexId][curEdgeId].Left + curEq;
55             Data[curVertexId][curEdgeId].Right = newRightBorder;
56             Data[curVertexId][curEdgeId].IdTo = createdVertexId;
57             curEdge.Left = *newRightBorder;
58
59             Data[createdVertexId].push_back(curEdge);
60             ++createdVertexId;
61             Data[createdVertexId - 1].push_back(TSuffixTree::TEdge(1 + curEq +
62                 passEq, end, createdVertexId));
63             PathSize[createdVertexId - 1] = PathSize[curVertexId] + *Data[
64                 curVertexId][curEdgeId].Right - Data[curVertexId][curEdgeId].
65                 Left;
66
67             if (lastCreatedVertexId > 0) {
68                 SuffixPtr[lastCreatedVertexId] = nextCreatedVertexId;
69             }
70             lastCreatedVertexId = nextCreatedVertexId;
71         }
72     }
73     int nextVertexId = SuffixPtr[curVertexId];
74     int nextEdgeId = -1;
75     int nextPassEq = PathSize[nextVertexId];
76     int nextCurEq = curEq + passEq - nextPassEq - 1;
77     while (nextCurEq > 0) {
78         for (size_t j = 0; j < Data[nextVertexId].size(); ++j) {
79             if (DataString[Data[nextVertexId][j].Left] == DataString[1 + 1 +
80                 nextPassEq]) {
81                 nextEdgeId = j;
82                 break;
83             }
84         }
85         int curRight = *Data[nextVertexId][nextEdgeId].Right;
86         int curLeft = Data[nextVertexId][nextEdgeId].Left;
87         if (nextEdgeId != -1 and curRight - curLeft <= nextCurEq) {
88             nextPassEq = nextPassEq + curRight - curLeft;
89             nextCurEq = nextCurEq - curRight + curLeft;
90             nextVertexId = Data[nextVertexId][nextEdgeId].IdTo;
91             nextEdgeId = -1;
92         } else {

```

```

87         break;
88     }
89 }
90 if (nextEdgeId != -1) {
91     curEq = nextCurEq;
92 } else {
93     curEq = 0;
94 }
95 passEq = nextPassEq;
96 curVertexId = nextVertexId;
97 curEdgeId = nextEdgeId;
98 ++l;
99 newVertex = false;
100 } else {
101     if (i < s.size() - 1) {
102         break;
103     }
104 }
105 }
106 }
107 }
108
109 std::string TSuffixTree::LexMinString(const size_t n) {
110     return LexMinString(0, n);
111 }
112
113 std::string TSuffixTree::LexMinString(const int id, const size_t n) {
114     char lexMinId = 0;
115     for (size_t i = 0; i < Data[id].size(); ++i) {
116         if (DataString[Data[id][i].Left] < DataString[Data[id][lexMinId].Left]) {
117             lexMinId = i;
118         }
119     }
120     size_t curEdgeLen = *Data[id][lexMinId].Right - Data[id][lexMinId].Left;
121     if (n <= curEdgeLen) {
122         return DataString.substr(Data[id][lexMinId].Left, n);
123     } else {
124         std::string res = DataString.substr(Data[id][lexMinId].Left, curEdgeLen);
125         return res + LexMinString(Data[id][lexMinId].IdTo, n - curEdgeLen);
126     }
127 }

```

3 Консоль

```
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ make
g++ -g -O2 -pedantic -std=c++17 -Wall -Wextra -Werror -c suffix_tree.cpp -o
suffix_tree.o
g++ -g -O2 -pedantic -std=c++17 -Wall -Wextra -Werror -c main.cpp -o main.o
g++ -g -O2 -pedantic -std=c++17 -Wall -Wextra -Werror suffix_tree.o main.o
-o solution
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ cat tests/1.in
xabcd
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ ./solution <tests/1.in
abcdx
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ cat tests/2.in
abracadabra
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ ./solution <tests/2.in
aabracadabr
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ cat tests/3.in
abaacaababaab
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ ./solution <tests/3.in
aababaababaac
```


4 Тест производительности

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

Тесты состоят из строк длины 100, 1000, 10000 и 50000.

```
engineerxl@engineerxl-GF63-Thin-9RCX:~/Study/DA/lab5$ make bench
```

```
./solution <tests/1.in  
Suffix tree time 0.178 ms  
./benchmark <tests/1.in  
Naive time 0.116 ms
```

```
./solution <tests/2.in  
Suffix tree time 1.680 ms  
./benchmark <tests/2.in  
Naive time 3.180 ms
```

```
./solution <tests/3.in  
Suffix tree time 17.795 ms  
./benchmark <tests/3.in  
Naive time 81.510 ms
```

```
./solution <tests/4.in  
Suffix tree time 18.455 ms  
./benchmark <tests/4.in  
Naive time 1334.365 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.

5 Выводы

Во время выполнения лабораторной работы я вспомнил префиксное дерево, способы хранения графов, реализовал алгоритм Укконена и ознакомился с приложениями суффиксного дерева.

Сложнее всего было понять, как алгоритм достигает линейной сложности и как работают суффиксные ссылки. Помогли рисунки в тетради и визуализатор [2].

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: алгоритм Кнута-Морриса-Пратта, Бойера-Мура, Рабина-Карпа.

Список литературы

- [1] *Алгоритм Укконена — Викиконспекты*
URL: 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).