

Московский авиационный институт
(национальный исследовательский университет)

Факультет компьютерных наук и прикладной математики

Кафедра вычислительной математики и программирования

Лабораторная работа №5 по курсу «Дискретный анализ»

Студент: О. А. Мезенин
Преподаватель: А. А. Кухтичев
Группа: М8О-306Б-21
Дата:
Оценка:
Подпись:

Москва, 2023

Лабораторная работа №5

Задача: Линеаризовать циклическую строку, то есть найти минимальный в лексикографическом смысле разрез циклической строки.

1 Описание

Первым шагом необходимо удвоить входную строку S , получив строку SS , и дописать терминальный символ $\$$. Затем построить суффиксное дерево получившийся строки. После чего нужно пройти дерево, двигаясь по дуге, которое начинается с наименьшей в лексикографическом смысле буквы (при этом терминальный символ $\$$ считается самым старшим символом). Двигаться по дереву нужно до тех пор, пока длина пути от корня не станет больше длины строки S . Получим номер K суффикса листа, в котором мы остановились. Ответом будет подстрока $SS[K : K + |\$|]$.

Для построения суффиксного дерева воспользуемся алгоритмом Укконена [1].

Пусть n — длина входной строки. Дерево строится за $O(n)$ и обходится за $O(n)$ времени. Итого алгоритм линеаризации циклической строки будет работать за $O(n)$ времени и потреблять $O(n)$ дополнительной памяти.

2 Исходный код

Создадим класс TSuffixTree – суффиксное дерево.

```
1 class TSuffixTree {
2     struct TNode {
3         size_t start;
4         int *end;
5         int suffixNumber;
6         TNode *suffixLink;
7         std::unordered_map<char, TNode *> children;
8
9         TNode(size_t start, int *end, int suffixNumber, TNode *suffixLink)
10            : start(start), end(end), suffixNumber(suffixNumber), suffixLink(suffixLink) {}
11     };
12     // methods
13 public:
14     TSuffixTree(std::string &text);
15     int FindMinimumLineSectionSuffixNumber(int n);
16 private:
17     void Build();
18     int Dfs(TNode *node, int n);
19     // variables
20 private:
21     std::string text;
22     TNode *root = nullptr;
23 };
```

В функции main читается строка S , создаётся строка SS , затем создаётся суффиксное дерево и вызывается функция для поиска номера суффикса минимального разреза строки.

При инициализации дерева к полученной строке добавляется терминальный символ \$, инициализируется корень и вызывается функция Build. Функция построения дерева с помощью алгоритма Укконена выглядит следующим образом:

```
1 void TSuffixTree::Build() {
2     /* builds a tree using Ukkonen's algorithm */
3     TNode *currentNode = root;
4     size_t j = 0;
5     size_t currentCharIndex = 0;
6     size_t currentEdgeLength = 0;
7     int *globalEnd = new int(-1);
8     for (size_t i = 0; i < text.length(); ++i) {
9         /* single phase algorithm (SPA) */
10        TNode *lastCreatedNode = nullptr;
11        (*globalEnd)++; // "add" a text[phase] to each leaf (rule 1)
12        for (; j <= i;) {
13            /* single extension algorithm (SEA) */
```

```

14     if (currentEdgeLength == 0) { // starts from root
15         currentCharIndex = i;
16     }
17
18     if (not currentNode->children.contains(text[currentCharIndex])) { // if
19         edges do not start with currentCharIndex (rules 1, 2)
20         currentNode->children[text[currentCharIndex]] = new TNode(i, globalEnd,
21             j, root);
22         if (lastCreatedNode) {
23             lastCreatedNode->suffixLink = currentNode;
24             lastCreatedNode = nullptr;
25         }
26     } else {
27         TNode *nextNode = currentNode->children[text[currentCharIndex]];
28         size_t nextEdgeLength = *(nextNode->end) - (nextNode->start) + 1;
29
30         if (currentEdgeLength >= nextEdgeLength) { // go down
31             currentCharIndex += nextEdgeLength;
32             currentEdgeLength -= nextEdgeLength;
33             currentNode = nextNode;
34             continue;
35         }
36
37         if (text[nextNode->start + currentEdgeLength] == text[i]) { // if edge
38             starts with text[phase] (rule 3)
39             ++currentEdgeLength;
40             if (lastCreatedNode and currentNode != root) {
41                 lastCreatedNode->suffixLink = currentNode;
42             }
43             break;
44         }
45
46         // insert a new node inside the edge (rule 2)
47         TNode *newNode = new TNode(nextNode->start, new int(nextNode->start +
48             currentEdgeLength - 1), -1, root);
49         currentNode->children[text[currentCharIndex]] = newNode;
50         nextNode->start += currentEdgeLength;
51         newNode->children[text[nextNode->start]] = nextNode;
52         newNode->children[text[i]] = new TNode(i, globalEnd, j, root);
53         if (lastCreatedNode) {
54             lastCreatedNode->suffixLink = newNode;
55         }
56         lastCreatedNode = newNode;
57     }
58
59     if (currentNode == root) {
60         if (currentEdgeLength) {
61             ++currentCharIndex;
62             --currentEdgeLength;

```

```

59         }
60     } else {
61         currentNode = currentNode->suffixLink;
62     }
63     ++j;
64 }
65 }
66 }

```

В этой функции реализуются все оптимизации и алгоритмы Single Phase Algorithm (SPA) и Single Extension Algorithm (SEA) [1].

Функции для поиска суффикса минимального разреза строки выглядят следующим образом:

```

1  int TSuffixTree::Dfs(TNode *node, int n) {
2      if (not node) {
3          return -1;
4      }
5      if ((*node->end) > n and node->suffixNumber != -1) {
6          return node->suffixNumber;
7      }
8      if (node->children.empty()) {
9          return -1;
10     }
11     char minChar = 'z'+1;
12     for (auto &key: node->children) {
13         if (key.first != '$' and key.first < minChar) {
14             minChar = key.first;
15         }
16     }
17     if (minChar == 'z'+1) {
18         minChar = '$';
19     }
20     return Dfs(node->children[minChar], n);
21 }
22
23 int TSuffixTree::FindMinimumLineSectionSuffixNumber(int n) {
24     return Dfs(root, n);
25 }

```

3 Консоль

```
aaprold@SAI:~/Documents/GitHub/MAI-DA/lab5$ make lab5
g++ -std=c++2a -pedantic -Wall -Wextra -Werror suffix_tree.cpp main.cpp -o
lab5.out
aprold@SAI:~/Documents/GitHub/MAI-DA/lab5$ ./lab5.out
xabcd
abcdx
aprold@SAI:~/Documents/GitHub/MAI-DA/lab5$ ./lab5.out
aabaacaa
aaaabaac
```

4 Тест производительности

Тесты производительности представляют из себя следующее: алгоритм линеаризации циклической строки с помощью суффиксного дерева будет сравниваться с примитивным алгоритмом: перебором всех строк, полученных циклическим смещением, для поиска наименьшей из таких. Скорость построения дерева также учитывается. Будет три теста со строками, состоящих из 10000, 100000 и 1000000 символов. Все строки случайные.

```
aprolld@SAI:~/Documents/GitHub/MAI-DA/lab5$ ./benchmark.out <tests/05.t
Suffix tree search time: 40782us
Primitive search time: 18709us
aprolld@SAI:~/Documents/GitHub/MAI-DA/lab5$ ./benchmark.out <tests/06.t
Suffix tree search time: 462362us
Primitive search time: 1252771us
aprolld@SAI:~/Documents/GitHub/MAI-DA/lab5$ ./benchmark.out <tests/07.t
Suffix tree search time: 3895521us
Primitive search time: 59833446us
```

Алгоритм линеаризации циклической строки с помощью суффиксного дерева проигрывает только на первом тесте, где длина строки равна 10000 символов. В остальных же тестах этот алгоритм оказался эффективней. Примитивный алгоритм работает за $O(n^3)$ времени, т.к. на каждом из n шагов создаётся новая строка за $O(n)$ и сравнивается с минимумом за $O(n)$ в худшем случае. Алгоритм линеаризации с помощью суффиксного дерева проиграл на первом тесте вероятно из-за большой константы.

5 Выводы

Выполнив пятую лабораторную работу по курсу «Дискретный анализ», разобрался в алгоритме Укконена для построения суффиксного дерева, а также изучил алгоритм линейаризации циклической строки с помощью суффиксного дерева.

Алгоритм линейаризации циклической строки с помощью суффиксного дерева (алгоритм Укконена) имеет линейную сложность и эффективен на строках большой длины.

Список литературы

[1] Гасфилд Дэн

Строки, деревья и последовательности в алгоритмах: Информатика и вычислительная биология / Пер. с англ. И.В.Романовского. — СПб.: Невский Диалект; БХВ-Петербург (ISBN 0-521-58519-8 (англ.); ISBN 5-7940-0103-8 ("Невский Диалект"); ISBN 5-94157-321-9 (БХВ-Петербург))

[2] *Реализации алгоритмов/Алгоритм Укконена*

URL: https://ru.wikibooks.org/wiki/Реализации_алгоритмов/Алгоритм_Укконена
(дата обращения: 23.09.2023).