

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: И. П. Моисеенков
Преподаватель: Н. С. Капралов
Группа: М8О-208Б-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 a до z).

Вариант 4: Линеаризовать циклическую строку, то есть найти минимальный в лексикографическом смысле разрез циклической строки.

Формат входных данных Некий разрез циклической строки.

Формат результата Минимальный в лексикографическом смысле разрез.

1 Описание

Согласно [1], суффиксное дерево для m -символьной строки S - это ориентированное дерево с корнем, имеющее ровно m листьев, занумерованных от 1 до m . Каждая внутренняя вершина, отличная от корня, имеет не меньше двух детей, а каждая дуга помечена непустой подстрокой строки S . Никакие две дуги, выходящие из одной и той же вершины, не могут иметь пометок, начинающихся с одного и того же символа. Главная особенность суффиксного дерева заключается в том, что для каждого листа i конкатенация меток дуг на пути от корня к листу i составляет суффикс строки S , который начинается в позиции i .

Для построения суффиксного дерева используется алгоритм Укконена. Каждую дугу дерева будем помечать не подстрокой, а индексами границ этой подстроки. Алгоритм состоит из нескольких фаз. Количество фаз совпадает с количеством символов в строке.

В структуре будем хранить активную вершину. Это положение, в которое мы будем вставлять следующий суффикс. Для однозначного определения активной вершины у нас будет три компоненты: сама активная вершина, активное ребро (ребро, по которому нужно перейти) и активная длина (количество элементов, через которое нужно перепрыгнуть на активном ребре).

Если мы вставили новый лист в дерево, то он останется листом навсегда. Поэтому удобно вместо правой границы этого узла использовать глобальную переменную end , которая будет увеличиваться на каждом шаге и всегда будет «указывать» на конец строки. Для обновления значения end нужно инкрементировать эту переменную при старте каждой фазы.

Введем понятие суффиксной ссылки. Пусть $x\alpha$ обозначает произвольную строку, где x - ее первый символ, а α - оставшаяся подстрока (возможно, пустая). Если для внутренней вершины v с путевой меткой $x\alpha$ существует другая вершина $s(v)$ с путевой меткой α , то указатель из v в $s(v)$ называется суффиксной ссылкой. Суффиксная ссылка нужна для быстрого перехода по дереву при вставке суффиксов.

Если при вставки мы не нашли вершины, начинающейся с вставляемого суффикса, то создаем новую вершину. Иначе спускаемся в эту вершину и смотрим, содержится ли наш суффикс целиком. Если да, то завершаем фазу, иначе создаем новую внутреннюю вершину и добавляем суффикс. После каждого создания узла необходимо обновить суффиксные ссылки, а после каждой фазы - активную вершину.

Сложность алгоритма Укконена - $O(n)$ по времени и по памяти.

Для нахождения минимального разреза циклической строки необходимо построить суффиксное дерево для удвоенной строки, выполнить его обход, всегда переходя по лексикографически минимальному ребру. Как только мы достигли листовой вершины, мы отсекаем вторую половину полученной строки и получаем ответ.

2 Исходный код

Создадим класс для вершины суффиксного дерева *TNode*. Его атрибуты - хеш-таблица с ребрами, суффиксная ссылка, позиции начала и конца подстроки и лексикографически минимальное ребро. Для поиска минимального разреза за линейное время будем находить для каждой вершины минимальное ребро при построении дерева.

Перейдем к классу самого дерева *TSuffixTree*. Его атрибуты - сама строка *text*, глобальная переменная для обозначения текущего конца строки *EndPos*, корень *Root*, активная вершина *ActiveNode*, активное ребро *ActiveEdge*, активная длина *ActiveLength* и количество суффиксов *Remaining*, которое нужно вставить. При создании объекта класса построение дерева выполняется автоматически. Для полноты дерева к концу строки добавляется сентинел (\$).

Выполнение одной фазы алгоритма описано в методе *Extend*. Там реализован алгоритм, описанный выше. Для нахождения лексикографически минимального разреза используется метод *FindLexMinString*, который по сути просто выполняет проход по дереву.

```
1 class TNode {
2 public:
3     TNode(int left = 0, int right = 0) : Start(left), End(right) {}
4     std::unordered_map<char, TNode*> Edges;
5     TNode* SuffixLink = nullptr;
6     int Start;
7     int End;
8     char LexMinEdge = 0;
9 };
10
11 class TSuffixTree {
12 public:
13     std::string Text;
14     int EndPos = -1;
15     TNode* Root = new TNode();
16
17     TNode* ActiveNode = Root;
18     int ActiveEdge = -1;
19     int ActiveLength = 0;
20
21     int Remaining = 0;
22
23     TSuffixTree(std::string &str);
24     ~TSuffixTree();
25     void Extend(int phase);
26     std::string FindLexMinString();
27     void Delete (TNode* node);
28 };
```

3 Консоль

```
mosik@LAPTOP-69S778GL:~/da_lab5$ ./solution
xabcd
abcdx
mosik@LAPTOP-69S778GL:~/da_lab5$ ./solution
abxabc
abcabx
```

4 Тест производительности

Сравним полученный алгоритм с наивным алгоритмом. Наивный алгоритм просто перебирает все возможные варианты разрезов и ищет среди них лексикографически наименьший.

Для сравнения я подготовил 3 теста. Первый тест - со строкой длины 500 символов, второй тест - 10^4 символов, третий - 10^6 .

```
mosik@LAPTOP-69S778GL:~/da_lab5$ ./suftree <test500
Time: 0.001716 ms
mosik@LAPTOP-69S778GL:~/da_lab5$ ./naive <test500
Time: 0.00017 ms
mosik@LAPTOP-69S778GL:~/da_lab5$ ./suftree <test1e4
Time: 0.0329882 ms
mosik@LAPTOP-69S778GL:~/da_lab5$ ./naive <test1e4
Time: 0.0042473 ms
mosik@LAPTOP-69S778GL:~/da_lab5$ ./suftree <test1e6
Time: 4.08355 ms
mosik@LAPTOP-69S778GL:~/da_lab5$ ./naive <test1e6
Time: 73.1969 ms
```

Видим, что на маленьких строках суффиксное дерево немного уступает наивному алгоритму. Это связано с тем, что построение дерева тоже требует времени. Однако на большой строке видно все преимущество суффиксных деревьев - выигрыш примерно в 18 раз. Это связано с тем, что наивный алгоритм работает за $O(n^2)$ - мы проверяем n разрезов, каждый из которых сравниваем с текущим минимальным (сравнение строк происходит за $O(n)$).

5 Выводы

Выполнив пятую лабораторную работу по курсу «Дискретный анализ», я познакомился с такой структурой данных как суффиксное дерево. Основное ее предназначение - поиск образцов в тексте. Поиск отличается от всеми известных алгоритмов КМП, Ахо-Корасик и др. тем, что мы предобрабатываем текст, а не образец.

Однако у суффиксных деревьев есть еще много применений. Благодаря ему мы можем за линейное время найти наибольшую общую подстроку, можем рассчитать статистику совпадений, построить суффиксный массив или найти минимальный лексикографический разрез.

Список литературы

- [1] Дэн Гасфилд. *Строки, деревья и последовательности в алгоритмах*. — БХВ-Петербург, 2003. Перевод с английского: И. В. Романовский. — 658 с. (ISBN 5-94157-321-9)