

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Дискретный анализ»

Студент: И. П. Попов
Преподаватель: А. А. Кухтичев
Группа: М8О-206Б
Дата:
Оценка:
Подпись:

Москва, 2022

Лабораторная работа №4

Задача: Необходимо реализовать один из стандартных алгоритмов поиска образцов для указанного алфавита.

Вариант алгоритма: Поиск одного образца при помощи алгоритма Бойера-Мура.

Вариант алфавита: Числа в диапазоне от 0 до $2^{32} - 1$.

Запрещается реализовывать алгоритмы на алфавитах меньшей размерности, чем указано в задании.

Формат ввода

Искомый образец задаётся на первой строке входного файла.

В случае, если в задании требуется найти несколько образцов, они задаются по одному на строку вплоть до пустой строки. Затем следует текст, состоящий из слов или чисел, в котором нужно найти заданные образцы.

Никаких ограничений на длину строк, равно как и на количество слов или чисел в них, не накладывается.

Формат вывода

В выходной файл нужно вывести информацию о всех вхождениях искомого образца в обрабатываемый текст: по одному вхождению на строку.

Для заданий, в которых требуется найти только один образец, следует вывести два числа через запятую: номер строки и номер слова в строке, с которого начинается найденный образец. В заданиях с большим количеством образцов, на каждое вхождение нужно вывести три числа через запятую: номер строки; номер слова в строке, с которого начинается найденный образец; порядковый номер образца.

Нумерация начинается с единицы. Номер строки в тексте должен отсчитываться от его реального начала (то есть, без учёта строк, занятых образцами).

Порядок следования вхождений образцов несущественен.

1 Описание

Интернет-ресурс [2] дает следующее описание алгоритму Бойера-Мура:

Алгоритм основан на трёх идеях.

1. Сканирование слева направо, сравнение справа налево

Совмещается начало текста (строки) и шаблона, проверка начинается с последнего символа шаблона. Если символы совпадают, производится сравнение предпоследнего символа шаблона и т. д. Если все символы шаблона совпали с наложенными символами строки, значит, подстрока найдена, и выполняется поиск следующего вхождения подстроки.

Если же какой-то символ шаблона не совпадает с соответствующим символом строки, шаблон сдвигается на несколько символов вправо, и проверка снова начинается с последнего символа.

Эти «несколько», вычисляются по двум эвристикам.

2. Эвристика стоп-символа

(Замечание: эвристика стоп-символа присутствует в большинстве описаний алгоритма Бойера — Мура, включая оригинальную статью Бойера и Мура, но не является необходимой для достижения оценки $O(n + m)$); , .)

3. Эвристика совпавшего суффикса

Неформально, если при чтении шаблона справа налево совпал суффикс S , а символ b , стоящий перед S в шаблоне (то есть шаблон имеет вид PbS), не совпал, то эвристика совпавшего суффикса сдвигает шаблон на наименьшее число позиций вправо так, чтобы строка S совпала с шаблоном, а символ, предшествующий в шаблоне данному совпадению S , отличался бы от b (если такой символ вообще есть). Формально, для данного шаблона $s[0..m-1]$ считается целочисленный массив $\text{suffshift}[0..m]$, в котором $\text{suffshift}[i]$ равно минимальному числу $j > 0$, такому что $s[i-j] \neq s[i-1]$ (если $i > 0$ и $i-j \geq 0$) и $s[i-j+k] = s[i-1+k]$ для любого $k > 0$, для которого выполняется $0 \leq i-j+k < m$ и $0 \leq i-1+k < m$ (для пояснения смотрите примеры ниже). Затем, если при чтении шаблона s справа налево совпало $k-1$ символов $s[m-1], s[m-2], \dots, s[m-k+1]$, а символ $s[m-k]$ не совпал, то шаблон сдвигается на $\text{suffshift}[m-k]$ символов вправо.

2 Исходный код

Основная логика работы программы состоит в следующем: двигаясь по входному тексту слева направо(начиная с символа, по индексу равному размеру входного паттерна), идет его сравнение с паттерном справа налево.

При выявлении несоответствия, происходит сдвиг по входному тексту на параметр *offset*. Данный параметр вычисляется с помощью двух эвристик, описанных выше. Применяется та эвристика, которая разрешит нам больший сдвиг.

Если паттерн пройден полностью, значит зафиксировано вхождение подстроки. И программа печатает номер строки и номер индекса в этой строке, где было зафиксировано вхождение.

Остановлюсь подробнее на подсчете сдвигов.

Код вычисления **эвристики плохого символа**:

```
1 class TBadLetterRule {
2 public:
3     TBadLetterRule(const vector<uint32_t>& pattern) :
4         patt_size(pattern.size()) {
5         for (int i = patt_size - 1; i >= 0; --i) {
6             rule[pattern[i]].push_back(i);
7         }
8     }
9
10    int Use(uint32_t letter, size_t idx_patt) const {
11        auto it = rule.find(letter);
12        if (it == rule.end()) {
13            return patt_size - idx_patt;
14        }
15        const vector<uint32_t>& idxs = it->second;
16        for (uint32_t i : idxs) {
17            if (i > idx_patt){
18                return 1;
19            }
20            else if (i < idx_patt) {
21                return idx_patt - i;
22            }
23        }
24        return patt_size;
25    }
26
27 private:
28     unordered_map<uint32_t, vector<uint32_t>> rule;
```

```

29 size_t patt_size;
30 };

```

Код вычисления эвристики хорошего суффикса

```

1 class TGoodSuffix {
2 public:
3     TGoodSuffix(const vector<uint32_t>& pattern) : patt_size(pattern.size())
4     {
5         rule = LFunction(pattern, gp_size);
6     }
7
8     int Use(size_t idx_patt) const {
9         if (idx_patt >= patt_size) {
10             return 1;
11         }
12         if (rule[idx_patt] == UNDEFINED) {
13             return patt_size - gp_size;
14         }
15         return patt_size - rule[idx_patt] - 1;
16     }
17
18     size_t Getgp_size() const {
19         return gp_size;
20     }
21 private:
22     vector<size_t> rule;
23     size_t patt_size;
24     size_t gp_size = 0;
25 };

```

Особого внимания заслуживает способ вычисления L-функции:

```

1 vector<size_t> ZFunction(const vector<uint32_t>& pattern) {
2     size_t n = pattern.size();
3     vector<size_t> res(n, 0);
4     size_t l = 0;
5     size_t r = 0;
6     for (size_t i = 1; i < n; ++i) {
7         if (i <= r) {
8             res[i] = min(r - i + 1, res[i - 1]);
9         }
10        while (i + res[i] < n && pattern[res[i]] == pattern[i + res[i]]) {
11            ++res[i];
12        }
13        if (i + res[i] - 1 > r) {
14            l = i;
15            r = i + res[i] - 1;
16        }
17    }
18    return res;

```

```

19 }
20
21 vector<size_t> NFunction(vector<uint32_t> pattern) {
22     reverse(pattern.begin(), pattern.end());
23     size_t n = pattern.size();
24     vector<size_t> z_func = ZFunction(move(pattern));
25     vector<size_t> res(n);
26     for (size_t i = 0; i < n; ++i) {
27         res[i] = z_func[n - i - 1];
28     }
29     return res;
30 }
31
32 const size_t UNDEFINED = -1;
33
34 vector<size_t> LFunction(const vector<uint32_t>& pattern, size_t& gp_size) {
35     gp_size = 0;
36     size_t n = pattern.size();
37     vector<size_t> res(n, UNDEFINED);
38     vector<size_t> n_func = NFunction(move(pattern));
39     for (size_t i = 0; i < n; ++i) {
40         size_t j = n - n_func[i];
41         if (j != n) {
42             res[j] = i;
43             if (i == n - j - 1) {
44                 gp_size = i + 1;
45             }
46         }
47     }
48     return res;
49 }

```

3 Консоль

tmp:

```

11 45 11 45 90
0011 45 011 0045 11 45 90    11
45 11 45 90

```

console output:

```

1,3
1,8

```

4 Тест производительности

Тест представлял из себя сравнение реализованного мной алгоритма Бойера-Мура с бинарным поиском, реализованным в STL посредством функции `upper_bound`.

В результате работы `benchmark.cpp` видны следующие результаты:

```
root@Lunidep:~/DA/da_lab4# ./wrapper.sh
[info][Tue May 17 11:35:24 MSK 2022] Stage #1. Compiling...
g++ -std=c++17 -pedantic -Wall -Wextra -Wno-unused-variable da_lab4.cpp -o
da_lab4
g++ -std=c++17 -pedantic -Wall -Wextra -Wno-unused-variable benchmark.cpp -o
benchmark
[info][Tue May 17 11:35:25 MSK 2022] Compiling OK
[info][Tue May 17 11:35:25 MSK 2022] Stage #2. Benchmark generating...
[info][Tue May 17 11:35:26 MSK 2022] Benchmark generating OK
[info][Tue May 17 11:35:26 MSK 2022] Stage #3. Benchmark results:
BM_search: 181 ms
find_search: 619 ms
[info][Tue May 17 11:35:28 MSK 2022] Benchmark OK
```

Из примера видно, реализованный мной алгоритм превосходит своего STL-соперника, потому что он обладает линейной сложностью, а бинарный поиск - операция за $O(\log(n))$.

5 Выводы

Выполнив вторую лабораторную работу по курсу «Дискретный анализ», мною были изучены различные алгоритмы поиска подстроки в строке и реализован один из них - алгоритм Бойера-Мура.

Реализация данного алгоритма очень удобно легла на ООП-парадигму языка C++, в частности методы подсчета сдвига в данной лабораторной были реализованы с её помощью.

Список литературы

- [1] Томас Х. Кормен, Чарльз И. Лейзерсон, Рональд Л. Ривест, Клиффорд Штайн. *Алгоритмы: построение и анализ, 2-е издание*. — Издательский дом «Вильямс», 2007. Перевод с английского: И. В. Красиков, Н. А. Орехова, В. Н. Романов. — 1296 с. (ISBN 5-8459-0857-4 (рус.))
- [2] Алгоритм Бойера—Мура. URL: <https://ru.wikipedia.org/wiki/Алгоритм->