

Московский авиационный институт
(национальный исследовательский университет)

Факультет компьютерных наук и прикладной математики

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Дискретный анализ»

Студент: Л. А. Постнов
Преподаватель: А. А. Кухтичев
Группа: М8О-206Б-22
Дата: 05.06.2024
Оценка:
Подпись:

Москва, 2024

Лабораторная работа №4

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

Вариант алгоритма: Поиск одного образца при помощи алгоритма Апостолико-Джанкарло.

Вариант алфавита: Слова не более 16 знаков латинского алфавита (регистронезависимые).

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

1 Описание

Основная идея алгоритма Апостолико-Джанкарло состоит в том, чтобы модифицировать алгоритм Бойера-Мура таким образом, чтобы опускать заведомо неприводящие к положительным результатам сравнения.

Про сложностьную оценку алгоритма сказано в [1]: Апостолико и Джанкарло предложили вариант алгоритма Бойера-Мура, который допускает замечательно простое доказательство линейной оценки наихудшего времени счета. В этом варианте никакой символ из T не участвует в сравнениях после его первого совпадения с каким-нибудь символом из P . Отсюда немедленно следует, что число сравнений не превзойдет $2m$. Каждое сравнение дает либо совпадение, либо несовпадение; последних может быть только T , так как при каждом несовпадении происходит ненулевой сдвиг P а совпадений - не больше T , так как никакой символ T не сравнивается после совпадения с символом из P . Мы покажем, что и остальная вычислительная работа (кроме сравнений) в этом методе линейно зависит от T .

Для этого мы будем держать массив, который равен по размерам тексту, где в каждой ячейке будем хранить число, которое означает, сколько позиций совпадет с суффиксом шаблона начиная с этой позиции при движении в сторону начала текста. Этот массив поможет нам опускать ненужные сравнения.

2 Исходный код

Рассмотрим функции, которые нужны для работы сдвигов по алгоритму Бойера-Мура

```
1  std::map<std::string, std::vector<int>> bad_symbol_prerocess(std::vector<std::string>
    &pattern) {
2      std::map<std::string, std::vector<int>> result;
3      for (int i = pattern.size() - 1; i >= 0; --i) {
4          result[pattern[i]].push_back(i);
5      }
6      return result;
7  }
8
9  std::vector<int> z_function(const std::vector<std::string> &s) {
10     int n = s.size();
11     std::vector<int> z(n);
12     int left = 0, right = 0;
13
14     for (int i = 1; i < n; ++i) {
15         if (i <= right)
16             z[i] = std::min(right - i + 1, z[i - left]);
17         while (i + z[i] < n && s[z[i]] == s[i + z[i]])
18             ++z[i];
19         if (i + z[i] - 1 > right) {
20             left = i;
21             right = i + z[i] - 1;
22         }
23     }
24
25     return z;
26 }
27
28 std::vector<int> n_function(std::vector<std::string> s) {
29     std::reverse(s.begin(), s.end());
30     std::vector<int> z = z_function(s), n(s.size());
31     for (int i = 1; i < z.size(); ++i) {
32         n[z.size() - i - 1] = z[i];
33     }
34     return n;
35 }
36
37 std::vector<int> l_fuction(const std::vector<int> &n) {
38     std::vector<int> l(n.size());
39     for (int i = 0; i < n.size(); ++i) {
40         if (n[i]) {
41             l[n.size() - n[i]] = i;
42         }
43     }
```

```

44     return 1;
45 }
46
47
48 int good_suffix(const std::vector<int> &l, int i) {
49     if (l.size() > i && l[i]) {
50         return l.size() - l[i];
51     }
52     return 0;
53 }
54
55
56 int bad_symbol(std::map<std::string, std::vector<int>> &table, const std::string &c,
57     int pos) {
58     if (!table[c].empty()) {
59         for (auto elem: table[c]) {
60             if (elem < pos) return pos - elem;
61         }
62     }
63     return 1;
64 }

```

Рассмотрим реализацию самой функции поиска:

```

1  std::vector<int> Search(std::vector<std::string> &text, std::vector<std::string> &
    pattern) {
2      auto N = n_function(pattern);
3      auto L = l_fuction(N);
4      auto table = bad_symbol_prerocess(pattern);
5
6      int h = pattern.size() - 1;
7      std::vector<int> M(text.size(), -1);
8      std::vector<int> positions;
9      while (h < text.size()) {
10         bool flag = true;
11         int position_to_stop = h - pattern.size();
12         std::string mismatched;
13         int i = pattern.size() - 1;
14         for (int j = h; j > position_to_stop;) {
15             if (M[j] == -1 || ((M[j] == 0) && (N[i] == 0))) {
16                 if (text[j] == pattern[i]) { // 1st case
17                     if (i > 0) {
18                         --i;
19                         --j;
20                     } else {
21                         break;
22                     }
23                 } else {
24                     mismatched = text[j];
25                     M[h] = h - j;

```

```

26         flag = false;
27         break;
28     }
29     } else if (M[j] < N[i] && M[j]) { // 2nd case
30         j -= M[j];
31         i -= M[j];
32     } else if (M[j] == N[i] && M[j]) { // 3rd case
33         if (i == N[i]) {
34             M[h] = h - j;
35             break;
36         }
37         i -= M[j];
38         j -= M[j];
39     } else if (M[j] > N[i]) { // 4th case
40         if (i == N[i]) {
41             M[h] = h - j;
42             break;
43         } else if (N[i] < i) {
44             mismatched = pattern[i - N[i] - 1];
45             M[h] = h - j;
46             flag = false;
47             break;
48         } else {
49             break;
50         }
51     } else { // 5th case
52         flag = false;
53         break;
54     }
55 }
56 if (flag) {
57     positions.push_back(h - pattern.size() + 1);
58     ++h;
59 } else {
60     int bad_symbol_res = bad_symbol(table, mismatched, i);
61     int good_suffix_res = good_suffix(L, h);
62     h += std::max(bad_symbol_res, std::max(1, good_suffix_res));
63 }
64 }
65 return positions;
66 }

```

3 Консоль

```
wednees@MSI:/mnt/c/Users/leoni/OneDrive/Desktop/study/Discrete_Analysis/lab4$  
./solution  
cat dog cat dog bird  
CAT dog CaT Dog Cat DOG bird CAT  
dog cat dog bird  
1,3  
1,8
```

4 Тест производительности

Тесты производительности представляют из себя следующее: наивный алгоритм поиска сравнивается с алгоритмом Апостолико-Джанкарло. Тестов будет три: на 10^4 , 10^5 и 10^6 строк.

```
wednees@MSI:/mnt/c/Users/leoni/OneDrive/Desktop/study/Discrete_Analysis/
lab4$ ./benchmark <test10000.txt
AG: 4ms
Naive: 6ms
wednees@MSI:/mnt/c/Users/leoni/OneDrive/Desktop/study/Discrete_Analysis/
lab4$ ./benchmark <test100000.txt
AG: 39ms
Naive: 78ms
wednees@MSI:/mnt/c/Users/leoni/OneDrive/Desktop/study/Discrete_Analysis/
lab4$ ./benchmark <test1000000.txt
AG: 429ms
Naive: 833ms
```

Как можно заметить, алгоритм Апостолико-Джанкарло работает быстрее, что связано с тем, что сложность наивного алгоритма составляет $O(n * m)$, в то время как сложность алгоритма Апостолико-Джанкарло $O(m + n)$, где n - длина исходного текста, а m - длина шаблона.

5 Выводы

Выполняя данную лабораторную работу по курсу «Дискретный анализ», я познакомился с алгоритмом Апостолико-Джанкарло, разобрался в нем и реализовал на языке программирования. Помимо этого алгоритма, потребовалось еще хорошо разобраться в алгоритме Бойера-Мура, так как алгоритм Апостолико-Джанкарло является его модификацией.

Помимо изучения алгоритмов, я также узнал про работу некоторых функций из стандартной библиотеки C++, предназначенных для работы со строчными данными, так как мой вариант представляет в качестве алфавита имеет не просто символы, как в обычном алфавите, а слова, и, как следствие, нужно было подумать, как организовать ввод таких данных.

Список литературы

- [1] Дэн Гасфилд. *Строки, деревья и последовательности в алгоритмах. Информатика и вычислительная биология*. Издательский дом «Невский Диалект БХВ-Петербург», Санкт-Петербург, 2003. Перевод с английского И. В. Романовского.