

Государственное образовательное учреждение высшего профессионального
образования
«Московский государственный технический университет имени Н. Э.
Баумана»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и система управления»
КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

ЛАБОРАТОРНАЯ РАБОТА №1
ПО КУРСУ «АНАЛИЗ АЛГОРИТМОВ»

Расстояние Левенштейна и Дамерау-Левенштейна

Выполнил: Сорокин А.П., гр. ИУ7-52Б
Преподаватели: Волкова Л.Л., Строганов Ю.В.

Москва, 2019г.

Оглавление

| | |
|---|-----------|
| Введение | 2 |
| 1 Аналитическая часть | 3 |
| 1.1 Задачи | 3 |
| 1.2 Описание алгоритмов | 3 |
| 2 Конструкторская часть | 5 |
| 2.1 Схемы алгоритмов | 5 |
| 2.2 Структуры данных | 5 |
| 3 Технологическая часть | 9 |
| 3.1 Требования к программному обеспечению | 9 |
| 3.2 Средства реализации | 9 |
| 3.3 Листинг кода | 9 |
| 3.4 Тесты | 12 |
| 4 Экспериментальная часть | 13 |
| 4.1 Примеры работы | 13 |
| 4.2 Исследование времени выполнения | 13 |
| Заключение | 14 |
| Список литературы | 15 |

Введение

Расстояние Левенштейна определяет минимальное количество операций, необходимых для превращения одной строки в другую. Задача определения такого минимума актуальна, так как она решает множество проблем в теории информации и компьютерной лингвистике, например:

- исправление ошибок в словах при вводе (при в поисковых ситсемах, базах данных, программах автоматического определения текста);
- сравнении текстовых файлов (к примеру, утилита diff);
- сравнение белков, генов и хромосом в биоинформатике.

1. Аналитическая часть

1.1 Задачи

Цель лабораторной работы: исследовать расстояния Левенштейна и Дамерау-Левенштейна. Для достижения этой цели были поставлены следующие задачи:

- изучить алгоритмы вычисления расстояний между строками;
- применить методы динамического программирования для матричной реализации алгоритмов;
- сравнить матричную и рекурсивную реализацию алгоритмов;
- оценить эффективность каждой из реализаций по времени и памяти.

1.2 Описание алгоритмов

При нахождении расстояния Левенштейна определяется минимального количество операций следующих видов:

- вставка (I - insert);
- удаление (D - delete);
- замена (R - replace);
- совпадение (M - match).

При нахождении расстояния Дамерау-Левенштейна добавляется операция транспозиции (T - transpose), или перестановки двух соседних символов.

Таким образом, если заданы две строки S_1 и S_2 с длинами m и n соответственно над некоторым алфавитом, то расстояние Левенштейна можно вычислить по следующей рекуррентной формуле:

$$D(S_1[1..m], S_2[1..n]) = \min(D(S_1[1..m-1], S_2[1..n]) + 1), \quad (1.1)$$

$$D(S_1[1..m], S_2[1..n-1]) + 1), \quad (1.2)$$

$$D(S_1[1..m-1], S_2[1..n-1]) + m(S_1[m], S_2[n])) \quad (1.3)$$

где

$$m(a, b) = \begin{cases} 0, a = b \\ 1, a \neq b \end{cases},$$

Соотношения в рекуррентной формуле отвечают за соответствующие разрешённые операции:

- (1.1) - вставка;
- (1.2) - удаление;
- (1.3) - замена или совпадение в зависимости от результата $m(a, b)$.

При вычислении расстояния Дамерау-Левенштейна в рекуррентную формулу вносится дополнительное соотношение в минимум:

$$D(S_1[1..m-2], S_2[1..n-2]) + 1 \quad (1.4)$$

Соотношение (1.4) вносится в качестве дополнительного аргумента минимума только при выполнении следующих условий:

- $m > 2, n > 2$;
- $S_1[m] = S_2[n-1]$;
- $S_1[m-1] = S_2[n]$.

Тривиальным случаем в рекуррентной формуле является случай, когда одна из строк пустая. В этом случае расстояние Левенштейна равно длине другой строки.

Расстояния Левенштейна и Дамерау-Левенштейна можно также вычислить, используя матрицу, в которой разрешённые операции определены следующим образом:

- движение по столбцам вправо - вставка;
- движение по строкам вверх - удаление;
- движение по диагонали - замена/совпадение.

2. Конструкторская часть

2.1 Схемы алгоритмов

2.2 Структуры данных

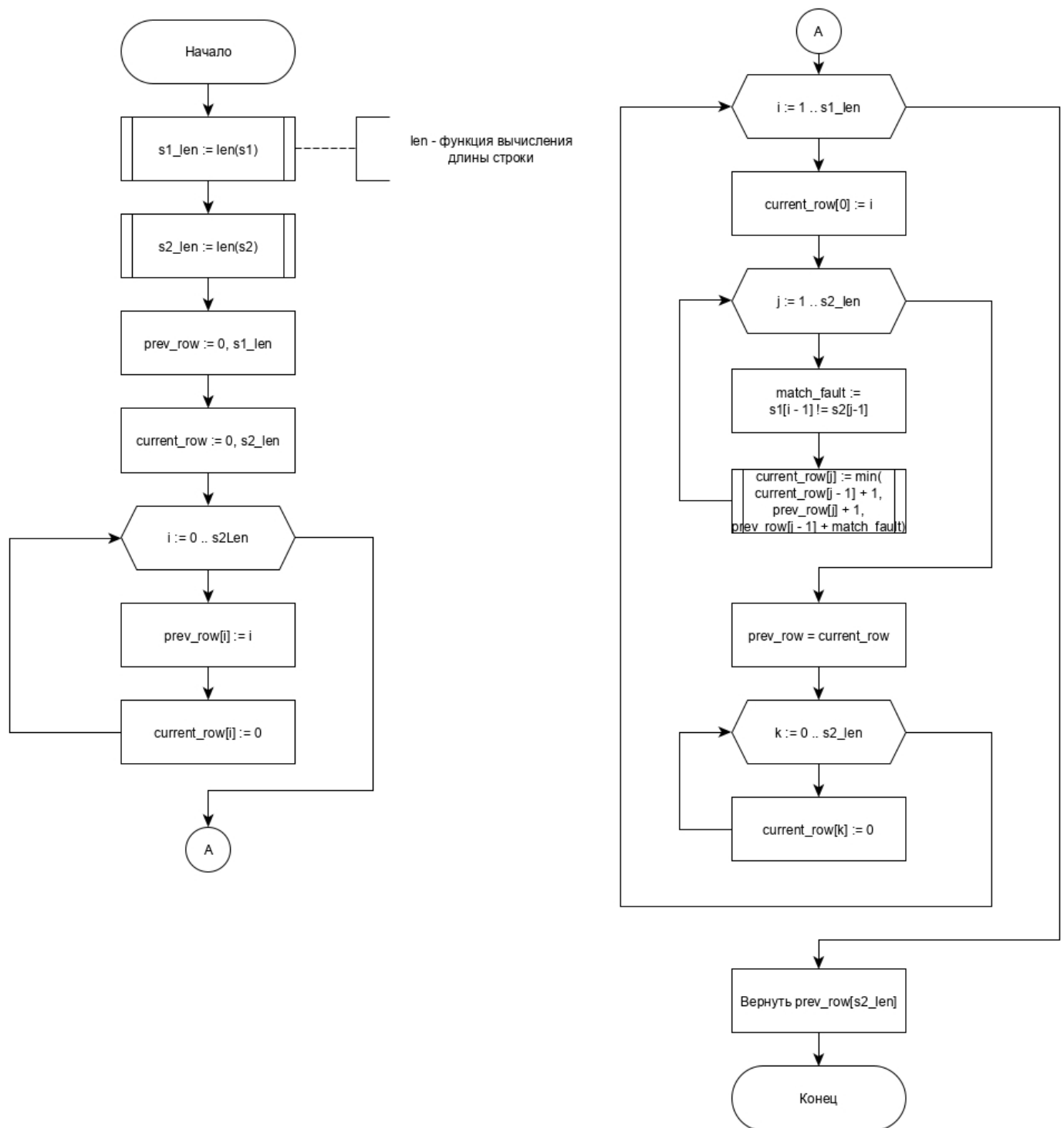


Рис. 2.1: Алгоритм Левенштейна

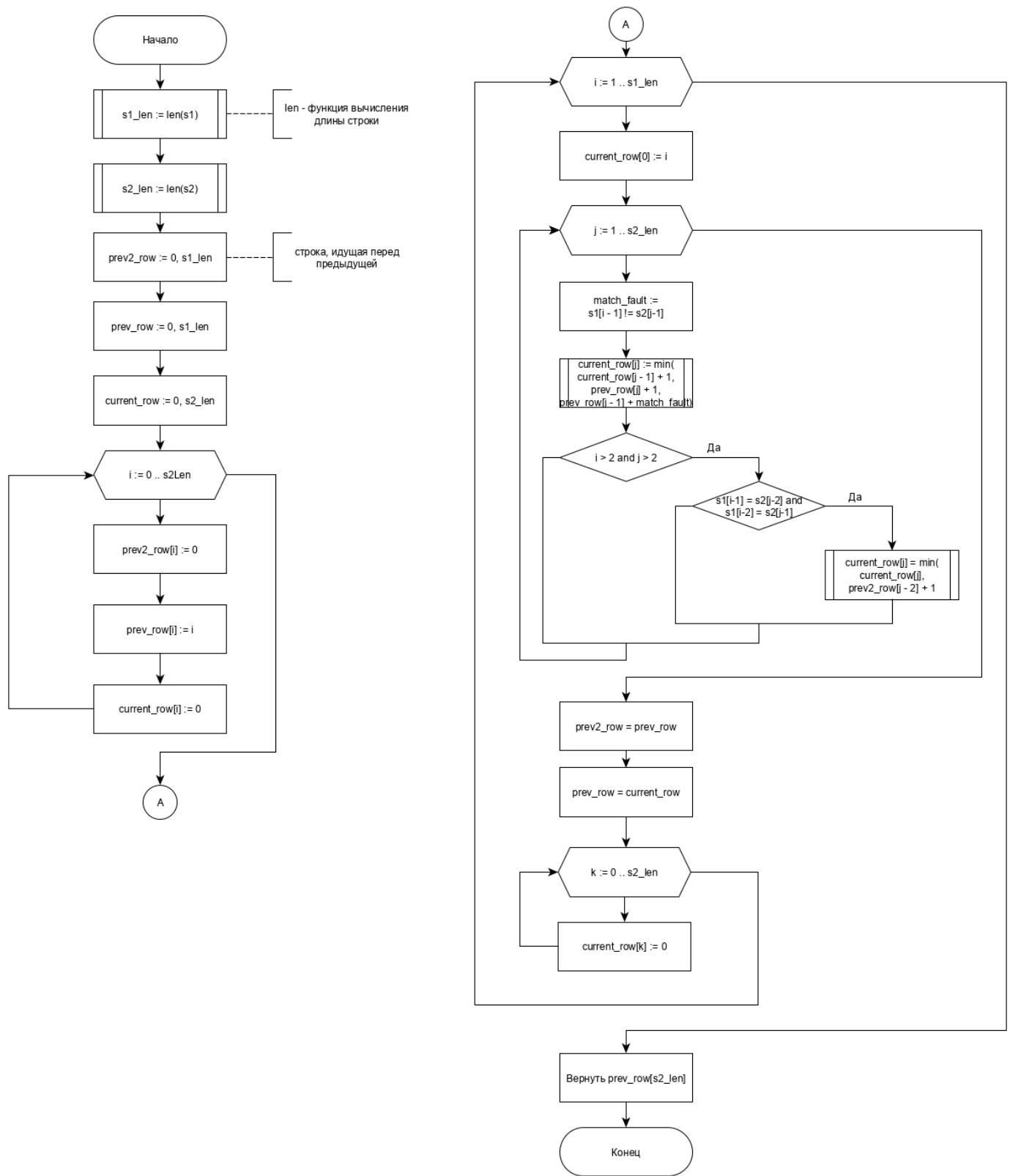


Рис. 2.2: Алгоритм Дамерау-Левенштейна (матричная реализация)

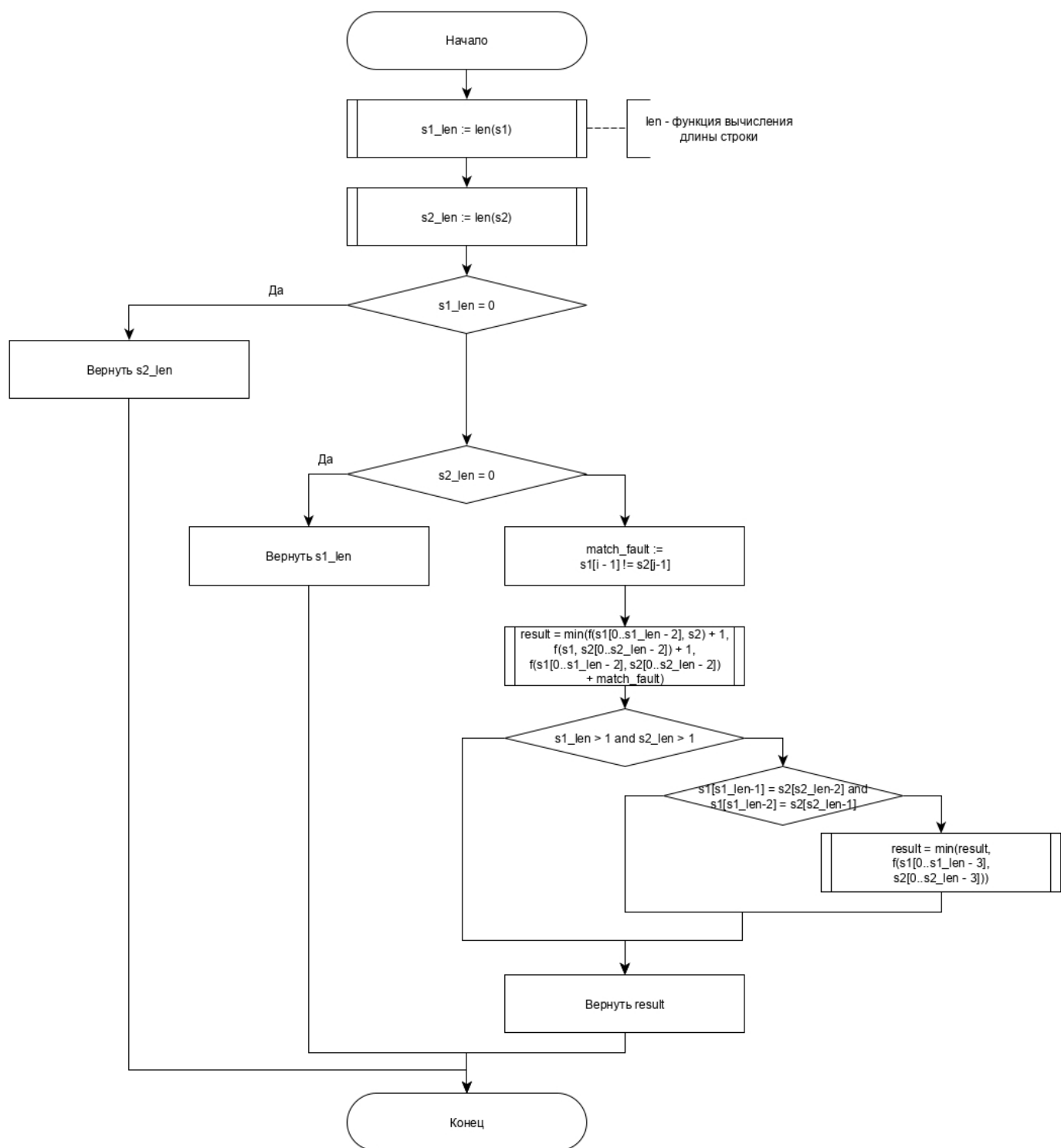


Рис. 2.3: Алгоритм Дамерау-Левенштейна (рекурсивная реализация)

3. Технологическая часть

3.1 Требования к программному обеспечению

3.2 Средства реализации

3.3 Листинг кода

Листинг 3.1: Расстояние Левенштейна (матрично)

```
1  def str_distance(s1, s2, to_print=False):
2      s1_len = len(s1)
3      s2_len = len(s2)
4
5      # initialization of first two rows
6      prev_row = [i for i in range(s2_len + 1)] # first row - [0, 1, ..., n]
7      current_row = [0] * (s2_len + 1)
8
9      if to_print:
10         print(prev_row)
11
12     for i in range(1, s1_len + 1): # row loop
13         # current row fill
14         current_row[0] = i
15         for j in range(1, s2_len + 1): # column loop
16             match_fault = int(s1[i - 1] != s2[j - 1]) # symbol
17                 match
18                 current_row[j] = min(current_row[j - 1] + 1, #
19                                     horizontal
20                                     prev_row[j] + 1, # vertical
21                                     prev_row[j - 1] + match_fault) # diagonal
22
23         if to_print:
24             print(current_row)
25
26         # row switching
27         prev_row = current_row
28         current_row = [0] * (s2_len + 1)
29
30     return prev_row[-1] # value in bottom right corner of table
```

Листинг 3.2: Расстояние Дамерау-Левенштейна (матрично)

```

1 def str_distance(s1, s2, to_print=False):
2     s1_len = len(s1)
3     s2_len = len(s2)
4
5     # initialization of first two rows
6     prev2_row = [0] * (s2_len + 1)
7     prev_row = [i for i in range(s2_len + 1)]
8     current_row = [0] * (s2_len + 1)
9
10    if to_print:
11        print(prev_row)
12
13    for i in range(1, s1_len + 1):          # row loop
14        # current row fill
15        current_row[0] = i
16        for j in range(1, s2_len + 1):      # column loop
17            match_fault = int(s1[i - 1] != s2[j - 1]) # if symbol matches
18            current_row[j] = min(current_row[j - 1] + 1, # horizontal
19                                prev_row[j] + 1,         # vertical
20                                prev_row[j - 1] + match_fault) # diagonal
21
22        # transposition check
23        if i > 2 and j > 2:
24            if s1[i - 1] == s2[j - 2] and s1[i - 2] == s2[j - 1]:
25                current_row[j] = min(current_row[j],
26                                      prev2_row[j - 2] + 1)
27
28        if to_print:
29            print(current_row)
30
31        # row switching
32        prev2_row = prev_row
33        prev_row = current_row
34        current_row = [0] * (s2_len + 1)
35
36    return prev_row[-1]                    # value in bottom right corner of table

```

Листинг 3.3: Расстояние Дамерау-Левенштейна (рекурсивно)

```

1 def str_distance(s1, s2):
2     s1_len = len(s1)
3     s2_len = len(s2)
4
5     if s1_len == 0:
6         return s2_len
7     if s2_len == 0:
8         return s1_len
9
10    match_fault = int(s1[-1] != s2[-1])
11
12    result = min(str_distance(s1[:-1], s2) + 1,
13                str_distance(s1, s2[:-1]) + 1,
14                str_distance(s1[:-1], s2[:-1]) + match_fault)

```

```
15 |
16 |     if s1_len > 1 and s2_len > 1:
17 |         if s1[-1] == s2[-2] and s1[-2] == s2[-1]:
18 |             result = min(result, str_distance(s1[:-2], s2[:-2]) + 1)
19 |
20 |     return result
```

| Строка 1 | Строка 2 | Ожидание | Результат |
|--------------|--------------|----------|-----------|
| <пустая> | <пустая> | 0 0 0 | 0 0 0 |
| <пустая> | а | 1 1 1 | 1 1 1 |
| а | <пустая> | 1 1 1 | 1 1 1 |
| а | а | 0 0 0 | 0 0 0 |
| а | б | 1 1 1 | 1 1 1 |
| азы | базы | 1 1 1 | 1 1 1 |
| компьютер | компьютер | 1 1 1 | 1 1 1 |
| данны | данные | 1 1 1 | 1 1 1 |
| email.ru | mail.ru | 1 1 1 | 1 1 1 |
| programmmer | programmer | 1 1 1 | 1 1 1 |
| mail.rus | mail.ru | 1 1 1 | 1 1 1 |
| ашибка | ошибка | 1 1 1 | 1 1 1 |
| алгоритм | алгорифм | 1 1 1 | 1 1 1 |
| копия | копии | 1 1 1 | 1 1 1 |
| укрсовой | курсовой | 2 1 1 | 2 1 1 |
| аглоритм | алгоритм | 2 1 1 | 2 1 1 |
| унивре | универ | 2 1 1 | 2 1 1 |
| курс | курсовой | 4 4 4 | 4 4 4 |
| курсовой | курс | 4 4 4 | 4 4 4 |
| курсовой | курсовик | 2 2 2 | 2 2 2 |
| код | закодировать | 9 9 9 | 9 9 9 |
| закодировать | код | 9 9 9 | 9 9 9 |
| ccoders | recoding | 5 5 5 | 5 5 5 |
| header | subheader | 3 3 3 | 3 3 3 |
| subheader | header | 3 3 3 | 3 3 3 |
| subheader | overheader | 4 4 4 | 4 4 4 |

Таблица 3.1: Функциональные тесты

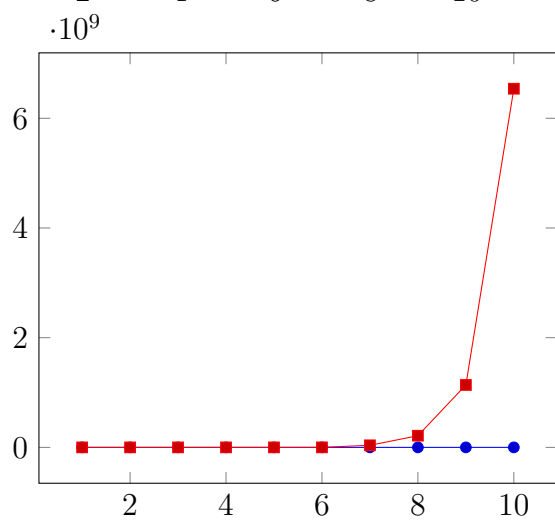
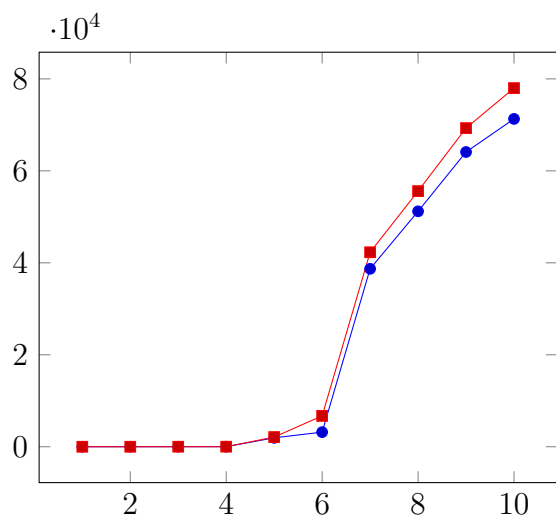
3.4 Тесты

Для проверки корректности работы были подготовлены следующие функциональные тесты:

4. Экспериментальная часть

4.1 Примеры работы

4.2 Исследование времени выполнения



Заключение

Алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна между строками были изучены и реализованы: были реализованы три варианта алгоритма для получения навыка динамического программирования.

Были исследованы затраты данных вариантов реализации по времени и памяти. Экспериментально было подтверждено, что рекурсивный вариант реализации алгоритма значительно проигрывает матричным вариантам при росте длины входных строк по обоим показателям.

Список литературы