

*Государственное образовательное учреждение высшего профессионального образования*

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

Рекуррентные соотношения: расстояние между строками

Москва 2018

## Содержание

Введение . . . . .	3
1 Аналитический раздел . . . . .	4
1.1 Описание Алгоритмов . . . . .	4
1.1.1 Расстояние Левенштейна . . . . .	4
1.1.2 Расстояние Дамерау – Левенштейна . . . . .	5
2 Конструкторский раздел . . . . .	7
2.1 Разработка алгоритмов . . . . .	7
2.1.1 Расстояние Левенштейна(обычный) . . . . .	7
2.1.2 Расстояние Дамерау – Левенштейна . . . . .	7
2.1.3 Расстояние Левенштейна(рекурсивный) . . . . .	7
3 Технологический раздел . . . . .	9
3.1 Листинг . . . . .	9
4 Исследовательский раздел . . . . .	11
4.1 Примеры работы . . . . .	11

## Введение

Целью работы является изучение и применение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна, а так же реализовать алгоритм Левенштейна в рекурсивном виде. Для достижения поставленной цели необходимо решить следующие задачи:

- Изучить алгоритмы Левенштейна и Дамера-Левенштейна нахождения расстояния между строками;
- Применить метод динамического программирования для матричной реализации указанных алгоритмов;
- Получить практические навыки реализации указанных алгоритмов: двух алгоритмов в матричной версии и алгоритма Левенштейна, реализованного рекурсивно;
- Провести сравнительный анализ линейной и рекурсивной реализаций алгоритма Левенштейна по затрачиваемым ресурсам(времени и памяти);
- Привести экспериментальное подтверждение различий во временной эффективности рекурсивной и нерекурсивной реализаций алгоритма Левенштейна при помощи разработанного ПО на материале замеров процессорного времени выполнения реализации на варьирующихся длинах строк;
- Описать и обосновать полученные результаты о выполненной работе;

# 1 Аналитический раздел

Перед теоритическим изложением алгоритмов, представленных в работе, требуется ввести понятия *редукционного расстояния* и *метода динамического программирования*.

*Редукционное расстояние*(*расстояние Эйнштейна*) – это минимальное количество редукционных операций, необходимых для преобразования одной строки в другую.

Есть следующие редукционные операции:

- Операции, вес которых - 1:
  - I - insert(вставка);
  - D - delete(удаление);
  - R - replace(замена);
- Операция, вес которой - 0:
  - M - match(совпадение);

Так минимальное расстояние между строками  $\min D(\text{'увлечение'}, \text{'развлечение'}) = 3$ , но чтобы найти это, требуется перебрать расстояния с разным выравниваем строк по отношению друг к другу.

		у	в	л	е	ч	е	н	и	е
р	а	з	в	л	е	ч	е	н	и	е
I	I	R	M	M	M	M	M	M	M	M

Проблема выравнивания решается рекуррентно через расстояния между подстроками фиксированной длины.

## 1.1 Описание Алгоритмов

Первым появившимся алгоритмом был алгоритм Левенштейна, который заложил фундамент в поиске расстояния между строками.

### 1.1.1 Расстояние Левенштейна

Расстояние Левенштейна имеет широкую область применения. Алгоритм используется в:

- поисковых системах, в базах данных, в автоматическом распознавание текста и речи для исправления ошибок и опечаток в слове;
- в утилитах для сравнения файлов(таких как diff);
- в биоинформатике для сравнения генов, хромосом и белков;

Алгоритм рекуррентно через расстояния между подстроками  $i$  и  $j$  находит расстояние между строками  $s1$  и  $s2$ . Проверяя, какое действие будет наиболее выгодным  $I(\text{insert})$ ,  $D(\text{delete})$ ,  $M(\text{match})$  или  $R(\text{replace})$ :

$$D(s1, s2) = D(s1[1..i], s2[1..j]) = \min \left( \begin{array}{l} D(s1[1..i], s2[1..j-1]) + 1, \\ D(s1[1..i-1], s2[1..j]) + 1, \\ D(s1[1..i-1], s2[1..j-1]) + \left[ \begin{array}{l} 0, \text{ если } s1[i] = s2[j], \\ 1, \text{ иначе} \end{array} \right] \end{array} \right) \quad (1.1)$$

, где  $i$  – длина подстроки строки  $s1$ , которая изначально равно длине  $s1$ ,  $j$  – длина подстроки строки  $s2$ , которая изначально равно длине  $s2$

Таким образом, применив метод динамического программирования мы разбиваем нашу задачу на небольшие подзадачи, которые достаточно легко можно решить. Данный метод удобно представлять в виде матрицы.

	$\lambda$	M	$\Gamma$
$\lambda$	0	1	2
M	0	0	1

<sup>1</sup>

В виде матрицы редукционные операции можно представить следующим образом:

- $I(\text{insert}) \rightarrow$ ;
- $D(\text{delete}) \downarrow$ ;
- $M(\text{match})$  or  $R(\text{replace}) \searrow$ ;

У данного алгоритма есть улучшенная версия, которую модифицировал Фредерик Дамерау.

### 1.1.2 Расстояние Дамерау – Левенштейна

Данный алгоритм применяется также как и обычный в:

- поисковых системах;
- биоинформатике(сравнение белков линейной структуры);

Причиной появления данного алгоритма было огромное количество ошибок ввода – ввод двух соседних символов не в том порядке. Отсюда появляется новая операция в дополнении к уже имеющимся:

- X - exchange(или T - transposition);

---

<sup>1</sup>(Представление  $D('M', 'MG')$ )

Отсюда, формула 1.1 переходит в:

$$D(s1[1..i], s2[1..j]) = \min \left( \begin{array}{c} D(s1[1..i], s2[1..j-1]) + 1, \\ D(s1[1..i-1], s2[1..j]) + 1, \\ D(s1[1..i-1], s2[1..j-1]) + \left[ \begin{array}{c} 0, \text{ если } s1[i] = s2[j], \\ 1, \text{ иначе} \end{array} \right], \\ D(s1[1..i-2], s2[1..j-2]) + 1 \end{array} \right), \quad (1.2)$$

Причем последняя операция(X) выполняется, если такая перестановка необходима и требуется, и это не совпадение.

## 2 Конструкторский раздел

Ниже представлены схемы алгоритмов – Дамерау - Левенштейна и двух реализаций Левенштейна(рекурсивный и обычный).

### 2.1 Разработка алгоритмов

#### 2.1.1 Расстояние Левенштейна(обычный)

На рисунке 2.1 представлена блок схема алгоритма Левенштейна.

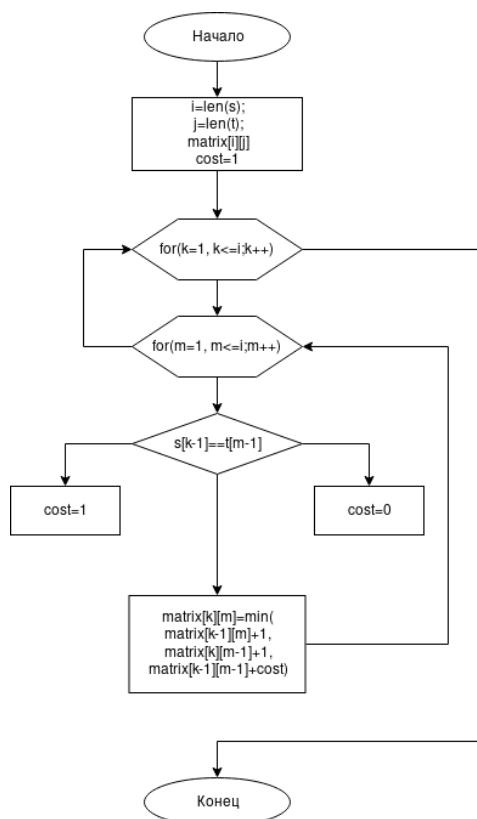


Рисунок 2.1 — Представлена схема алгоритма нахождения расстояния Левенштейна для матричной реализации

#### 2.1.2 Расстояние Дамерау – Левенштейна

На рисунке 2.2 представлена блок схема алгоритма Дамерау – Левенштейна.

#### 2.1.3 Расстояние Левенштейна(рекурсивный)

На рисунке 2.3 представлена блок схема рекурсивного алгоритма Левенштейна, в реализации которого используется функция *int match(char c, char d)*, которая возвращает 0, если символы  $c$  и  $d$  совпадают, иначе 1.

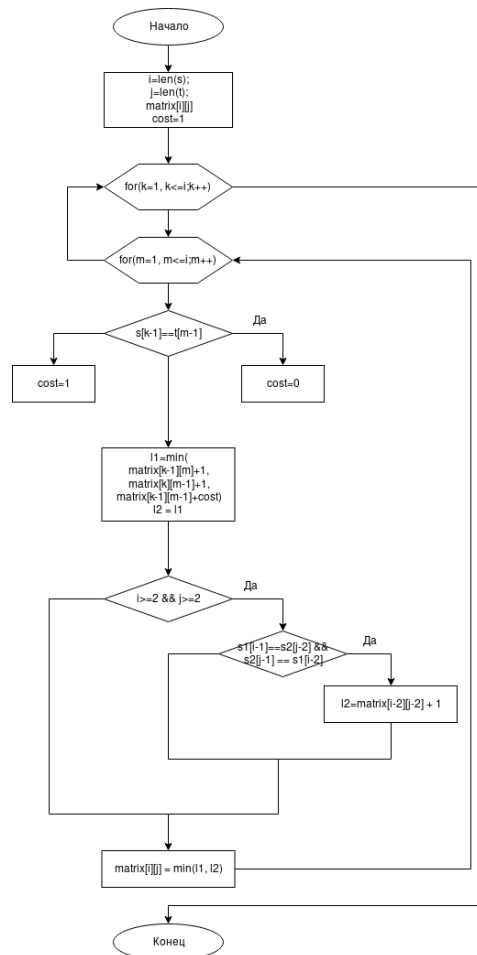


Рисунок 2.2 — Представлена схема алгоритма нахождения расстояния Дamerau – Левенштейна для матричной реализации

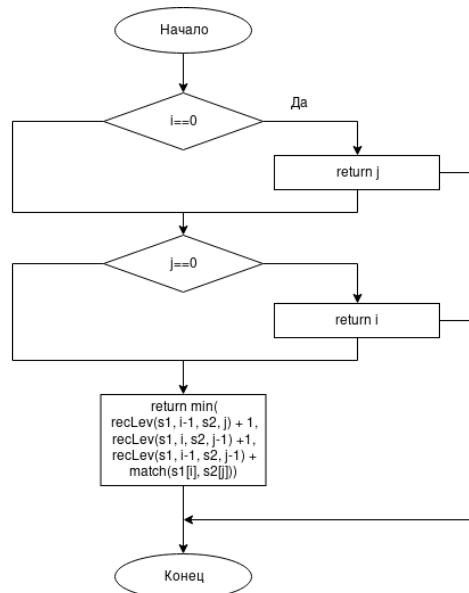


Рисунок 2.3 — Представлена схема алгоритма нахождения расстояния Левенштейна для рекурсивной реализации



### 3 Технологический раздел

Данные реализации разрабатывались на языке C, использовался компилятор gcc-8.2.1. Замер времени производится с помощью библиотеки *time.h*. Данная библиотека позволяет производить замер тиков, откуда можно получать время в секундах.

#### 3.1 Листинг

Алгоритм Левенштейна:

```
1 int levensteinDistance( char* s1, char* s2,
2                         int ( * matrix)[MAX_STRING_SIZE][MAX_STRING_SIZE] )
3 {
4     size_t len1 = strlen(s1);
5     size_t len2 = strlen(s2);
6
7     matrixInit(matrix, len1, len2);
8
9     int cost = 1;
10    for(int i = 1; i <= len1; i++)
11        for(int j = 1; j <= len2; j++) {
12            if(s1[i-1] == s2[j-1])
13                cost = 0;
14            else
15                cost = 1;
16            ( * matrix)[i][j] = min( 3,
17                                    ( * matrix)[i-1][j]+1,
18                                    ( * matrix)[i][j-1]+1,
19                                    ( * matrix)[i-1][j-1]+cost );
20        }
21 }
```

Листинг алгоритма Дамерау – Левенштейна:

```
1 int levensteinDistance( char* s1, char* s2,
2                         int ( * matrix)[MAX_STRING_SIZE][MAX_STRING_SIZE] )
3 {
4     size_t len1 = strlen(s1);
5     size_t len2 = strlen(s2);
6
7     matrixInit(matrix, len1, len2);
8
9     int cost = 1;
10    for(int i = 1; i <= len1; i++)
11        for(int j = 1; j <= len2; j++) {
12            int flag = 0;
13            if(s1[i-1] == s2[j-1])
14                cost = 0;
```

```

15         else
16             cost = 1;
17         int l1= min( 3,
18                     ( * matrix)[i-1][j]+1,
19                     ( * matrix)[i][j-1]+1,
20                     ( * matrix)[i-1][j-1]+cost );
21         int l2 = l1;
22         if (i>=2 && j>=2)
23             if (s1[i-1] == s2[j-2] && s2[j-1] == s1[i-2])
24                 l2 = ( * matrix)[i-2][j-2]+1;
25         ( * matrix)[i][j] = min(2, l1 , l2);
26     }
27 }

```

Листинг алгоритма Левенштейна в рекурсивной реализации:

```

1 int levensteinDistance( char* s1, int i, char* s2, int j )
2 {
3     if (i == 0) return j;
4     if (j == 0) return i;
5
6     return min( 3,levensteinDistance(s1, i-1, s2, j) + 1,
7                 levensteinDistance(s1, i, s2, j-1) + 1,
8                 levensteinDistance(s1, i-1, s2, j-1) + match(s1[i], s2[j
9     ]) );
10 }

```

Листинг функции match():

```

1 int match(char c, char d)
2 {
3     if (c == d) return 0;
4     else return 1;
5 }

```

## 4 Исследовательский раздел

Данный раздел посвящен тестированию и исследованию трех реализованных алгоритмов: алгоритмам Левенштейна, Дамерау – Левенштейна и рекурсивной реализации Левенштейна.

### 4.1 Примеры работы

Алгоритм Левенштейна:

Input s1: ОТАРА

Input s2: ТАРТАР

0 1 2 3 4 5 6

1 1 2 3 4 5 6

2 1 2 3 3 4 5

3 2 1 2 3 3 4

4 3 2 1 2 3 3

5 4 3 2 2 2 3

Input s1: MGTU

Input s2: MTGU

0 1 2 3 4

1 0 1 2 3

2 1 1 1 2

3 2 1 2 2

4 3 2 2 2

Алгоритм Дамерау – Левенштейна:

Input s1: MGTU

Input s2: MTGU

0 1 2 3 4

1 0 1 2 3

2 1 1 1 2

3 2 1 1 2

4 3 2 2 1

Input s1: KNUTH  
Input s2: TANENBAUM  
0 1 2 3 4 5 6 7 8 9  
1 1 2 3 4 5 6 7 8 9  
2 2 2 2 3 4 5 6 7 8  
3 3 3 3 3 4 5 6 6 7  
4 3 4 4 4 4 5 6 7 7  
5 4 4 5 5 5 5 6 7 8

Алгоритм Левенштейна(Рекурсивный):

Input s1: MGTU  
Input s2: MTGU  
Length is : 2  
Input s1: ОТАРА  
Input s2: ТАРТАР  
Length is : 3