



Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное  
учреждение высшего образования  
«Московский государственный технический университет имени  
Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)

---

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

## Отчет по лабораторной работе №1 по дисциплине "Анализ алгоритмов"

Тема Расстояние Левенштейна и Дamerau-Левенштейна

Студент Криков А. В.

Группа ИУ7-53Б

Преподаватели Волкова Л.Л., Строганов Ю.В.

# Оглавление

<b>Введение</b>	<b>2</b>
<b>1 Аналитическая часть</b>	<b>3</b>
1.1 Расстояние Левенштейна . . . . .	3
1.2 Расстояние Дамерау — Левенштейна . . . . .	4
<b>2 Конструкторская часть</b>	<b>5</b>
2.1 Разработка алгоритмов . . . . .	5
<b>3 Технологическая часть</b>	<b>9</b>
3.1 Выбор ЯП . . . . .	9
3.2 Требования к ПО . . . . .	9
3.3 Листинг кода . . . . .	9
<b>4 Исследовательская часть</b>	<b>14</b>
4.1 Технические характеристики . . . . .	14
4.2 Временные характеристики . . . . .	14
4.3 Характеристики по памяти . . . . .	15
<b>Заключение</b>	<b>17</b>

# Введение

В данной лабораторной работе будет рассмотрен алгоритм под названием "расстояние Левенштейна".

Данное расстояние показывает минимальное количество редакторских операций (вставки, замены и удаления), которые необходимы для перевода одной строки в другую. Это расстояние помогает определить схожесть двух строк.

Расстояние Левенштейна используется в компьютерной лингвистике для:

- исправления ошибок в слове
- сравнения текстовых файлов утилитой diff
- в биоинформатике для сравнения генов и белков

Однако кроме упомянутых трех ошибок (вставка лишнего символа, пропуск символа, замена одного символа другим), пользователь может нажимать на нужные клавиши не в том порядке. С этой проблемой поможет справиться расстояние Дameraу-Левенштейна. Данное расстояние задействует еще одну редакторскую операцию - транспозицию.

Целью данной работы является анализ и реализация алгоритма Дameraу-Левенштейна и Левенштейна.

Задачи лабораторной работы:

- изучение алгоритмов нахождения расстояния Левенштейна и Дameraу-Левенштейна;
- применение методов динамического программирования для реализации алгоритмов;
- получение практических навыков реализации алгоритмов Левенштейна и Дameraу — Левенштейна;
- сравнительный анализ алгоритмов на основе экспериментальных данных;
- подготовка отчета по лабораторной работе.

# 1 Аналитическая часть

Для преобразования одного слова в другое используются следующие операции:

- D - удаление
- I - вставка
- R - замена

Будем считать стоимость каждой вышеизложенной операции - 1. Введем понятие совпадения - M. Его стоимость будет равна - 0.

## 1.1 Расстояние Левенштейна

Расстояние Левенштейна между двумя строками a и b может быть вычислено по формуле:

$$D(i, j) = \begin{cases} 0 & i = 0, j = 0 \\ i & j = 0, i > 0 \\ j & i = 0, j > 0 \\ \min\{ \\ \quad D(i, j - 1) + 1 \\ \quad D(i - 1, j) + 1 & i > 0, j > 0 \\ \quad D(i - 1, j - 1) + m(a[i], b[j]) & (1.2) \\ \} \end{cases}, \quad (1.1)$$

где функция 1.2 определена как:

$$m(a, b) = \begin{cases} 0 & \text{если } a = b, \\ 1 & \text{иначе} \end{cases}. \quad (1.2)$$

## 1.2 Расстояние Дамерау — Левенштейна

В расстоянии Дамерау - Левенштейна задействуют еще одну операцию - транспозицию Т. Расстояние Дамерау — Левенштейна может быть найдено по формуле 1.3.

$$D_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ \\ \quad D_{a,b}(i, j - 1) + 1, \\ \quad D_{a,b}(i - 1, j) + 1, \\ \quad D_{a,b}(i - 1, j - 1) + m(a[i], b[j]), & \text{иначе} \\ \quad \left[ \begin{array}{ll} D_{a,b}(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a[i] = b[j - 1]; \\ & b[j] = a[i - 1] \\ & \infty, & \text{иначе} \end{array} \right. \\ \} \end{cases}, \quad (1.3)$$

## Вывод

Формулы Левенштейна и Дамерау — Левенштейна для расчета расстояния между строками задаются рекурсивно, а следовательно, алгоритмы могут быть реализованы рекурсивно или итерационно.

## 2 Конструкторская часть

### 2.1 Разработка алгоритмов

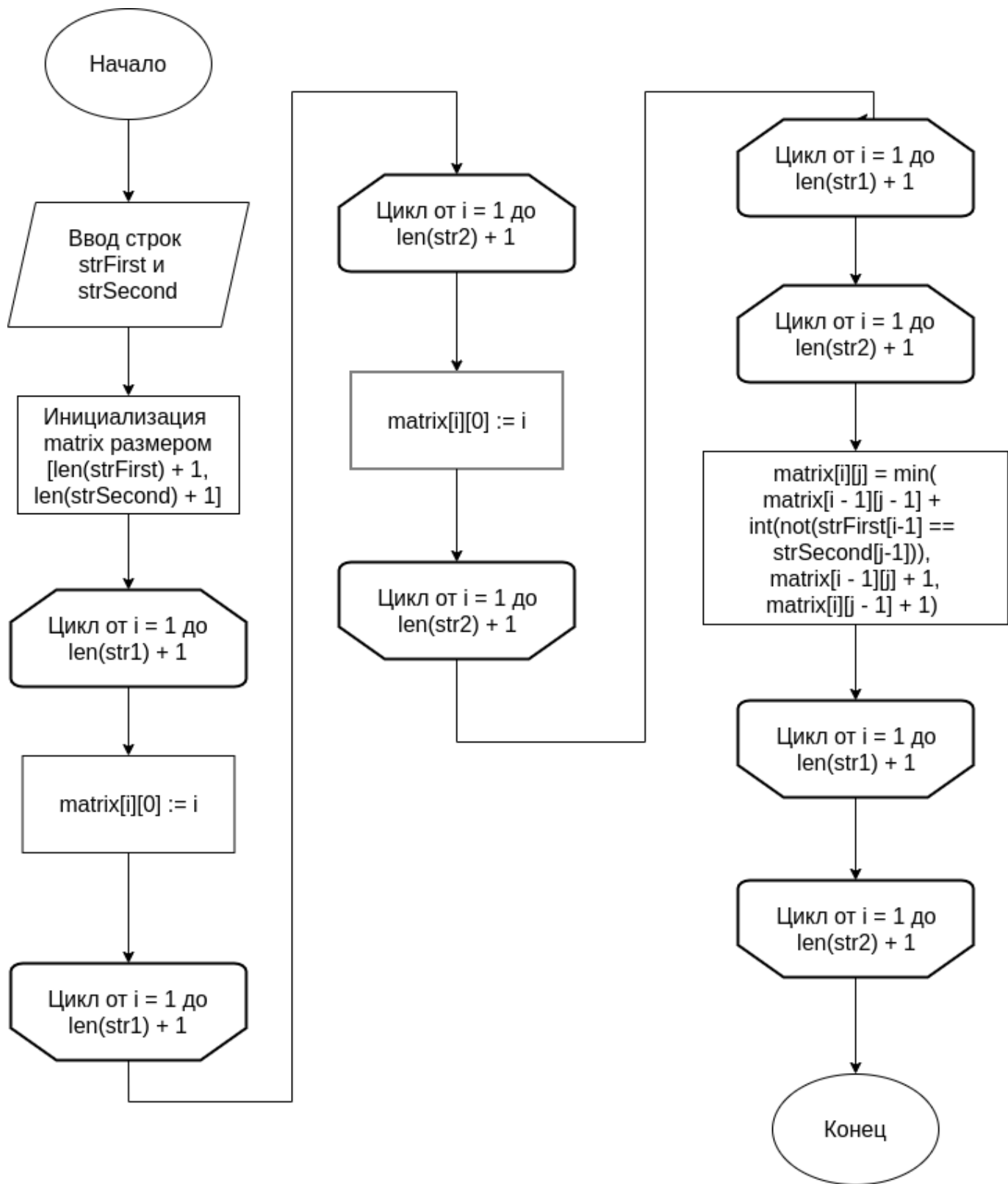


Рисунок 2.1 – Схема алгоритма поиска расстояния Левенштейна

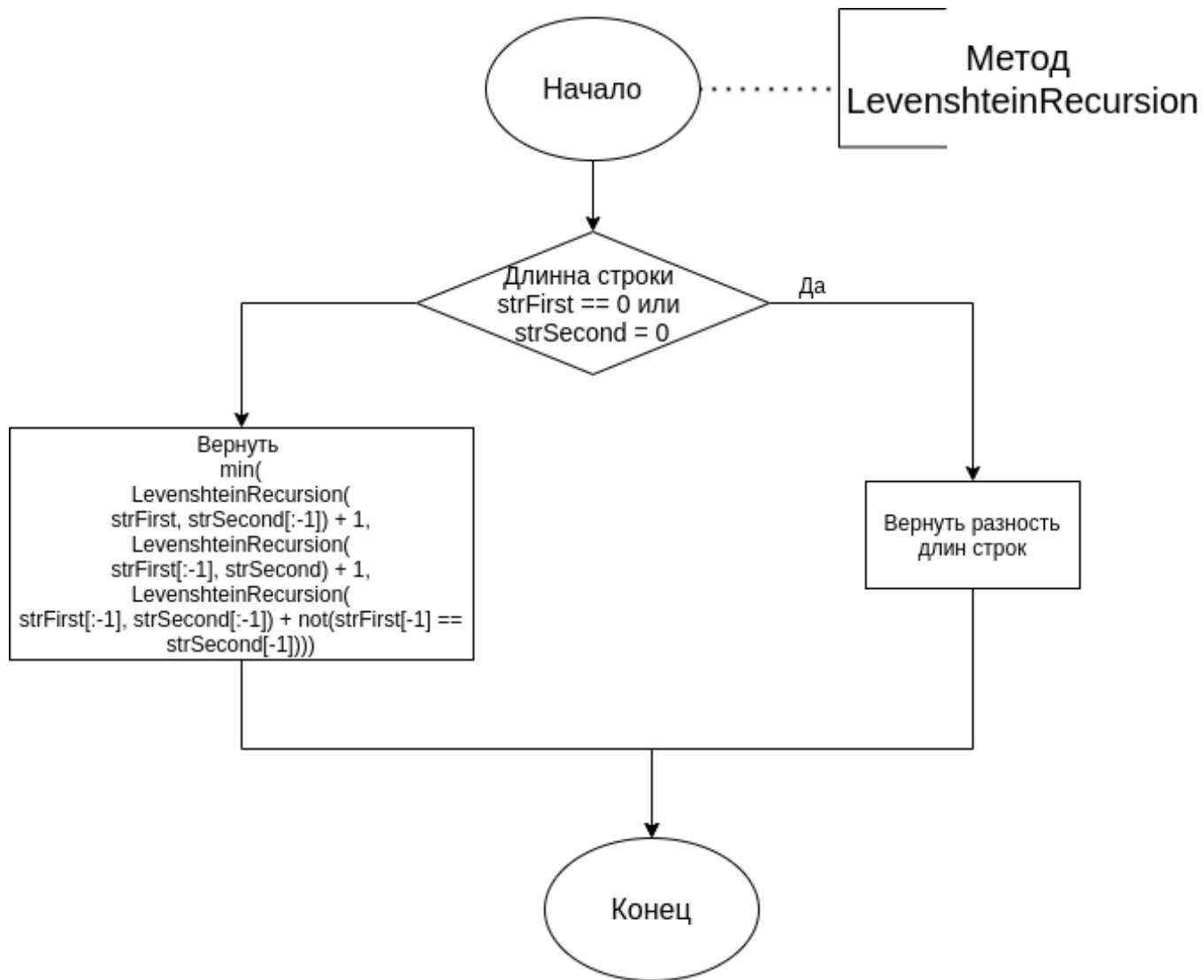


Рисунок 2.2 – Схема рекурсивного алгоритма нахождения расстояния Левенштейна



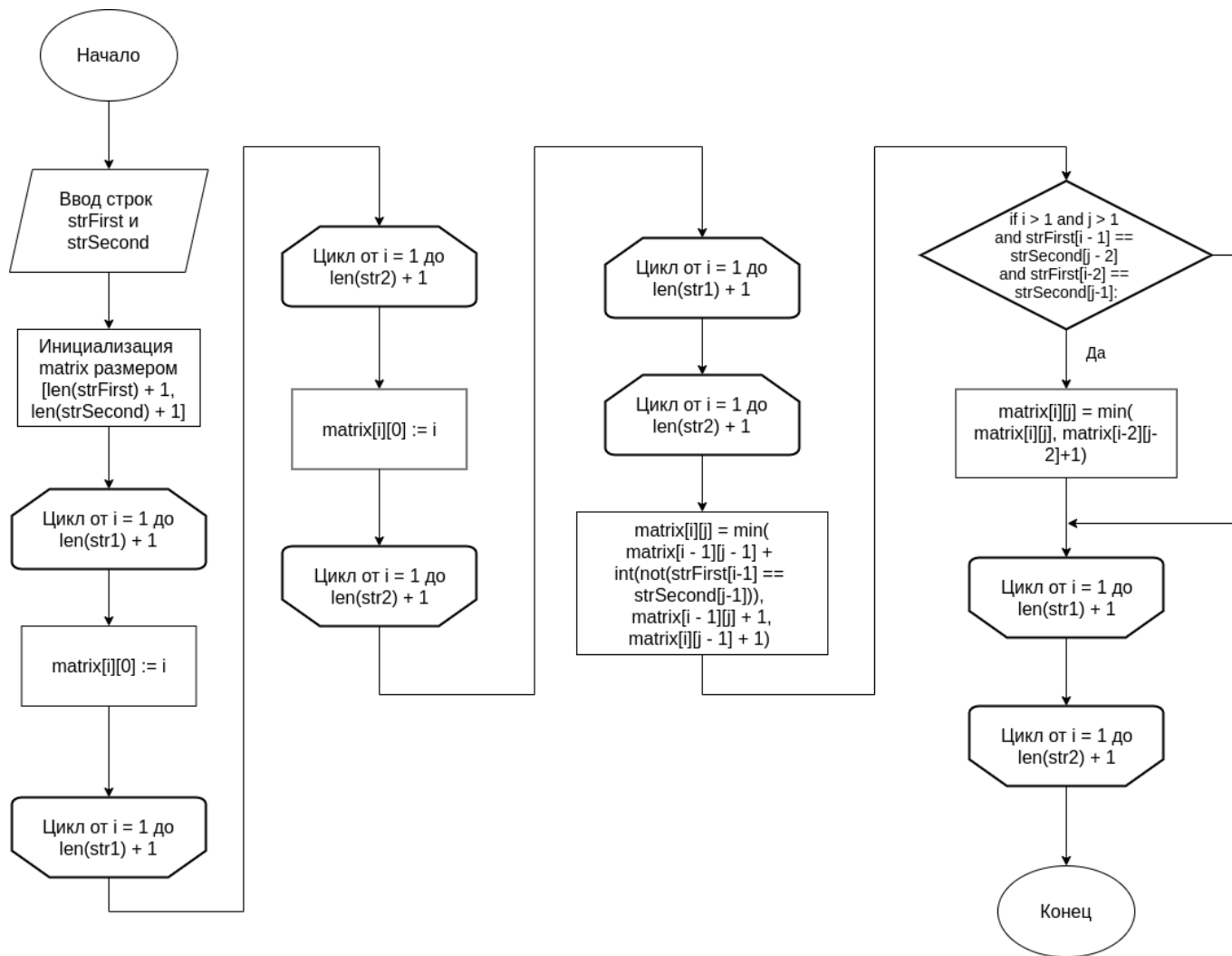


Рисунок 2.3 – Схема алгоритма нахождения расстояния  
Дамерау-Левенштейна

## 3 Технологическая часть

В данном разделе приведены требования к программному обеспечению, выбор ЯП и листинги кода.

### 3.1 Выбор ЯП

В данной лабораторной работе использовался язык программирования - Java. Данный язык простой и понятный, поэтому использовался мною. В качестве среды разработки была выбрана IntelliJ IDEA. IntelliJ IDEA подходит не только для Windows, но и для Linux.

### 3.2 Требования к ПО

К программе предъявляется ряд требований:

- на вход подаются две строки на русском или английском языке в любом регистре;
- на выходе — искомое расстояние для всех четырех методов и матрицы расстояний для всех методов, за исключением рекурсивного.
- замеры времени работы каждой реализации

### 3.3 Листинг кода

В листинге 3.1 приведена реализация алгоритмов нахождения расстояния Левенштейна и Дamerau — Левенштейна, а также вспомогательные функции.

Листинг 3.1 – Листинг с алгоритмами

```
1 import java.util.ArrayList;
2
3 public class Levenshtein {
4     public static int getIntBoolean(char first, char second) {
5         if (first == second) {
```

```

6         return 0;
7     }
8     return 1;
9 }
10
11 public static int Levenshtein(String first, String second) {
12     int n = first.length();
13     int m = second.length();
14
15     int[][] matrix = new int[n + 1][m + 1];
16     matrix[0][0] = 0;
17     for (int i = 1; i < n + 1; i++) {
18         matrix[i][0] = i;
19     }
20     for (int i = 1; i < m + 1; i++) {
21         matrix[0][i] = i;
22     }
23
24     for (int i = 1; i < n + 1; i++) {
25         for (int j = 1; j < m + 1; j++) {
26             matrix[i][j] = Math.min(
27                 Math.min(
28                     matrix[i - 1][j - 1] + getIntBoolean(first.charAt(i - 1),
29                         second.charAt(j - 1)), matrix[i - 1][j] + 1),
30                     matrix[i][j - 1] + 1
31                 );
32         }
33     }
34
35     return matrix[n][m];
36 }
37
38 public static int LevenshteinRecursion(String first, String second) {
39     if (first == "" || second == "") {
40         return Math.abs(first.length() - second.length());
41     }
42
43     int temp = (first.charAt(first.length() - 1) == second.charAt(second.length() -
44         1)) ? 0 : 1;
45     return Math.min(
46         Math.min(
47             LevenshteinRecursion(first, second.substring(0, second.length() -
48                 1)) + 1,
49             LevenshteinRecursion(first.substring(0, first.length() - 1),
50                 second) + 1),
51         LevenshteinRecursion(first.substring(0, first.length() - 1),
52             second.substring(0, second.length() - 1)) + temp
53     );

```

```

50     }
51
52
53     public static int DamerauLevenshteinRecursion(String first, String second) {
54         if (first == "" || second == "") {
55             return Math.abs(first.length() - second.length());
56         }
57         int temp = (first.charAt(first.length() - 1) == second.charAt(second.length() -
58             1)) ? 0 : 1;
59
60         int result = Math.min(
61             DamerauLevenshteinRecursion(first, second.substring(0, second.length() -
62                 1)) + 1,
63             Math.min(
64                 DamerauLevenshteinRecursion(first.substring(0, first.length() -
65                     1), second) + 1,
66                 DamerauLevenshteinRecursion(first.substring(0, first.length() -
67                     1), second.substring(0, second.length() - 1)) + temp
68             )
69         );
70
71         if (first.length() > 1 && second.length() > 1 &&
72             first.charAt(first.length() - 1) == second.charAt(second.length() - 2) &&
73             first.charAt(first.length() - 2) == second.charAt(second.length() - 1) ) {
74             result = Math.min(
75                 result,
76                 DamerauLevenshteinRecursion(first.substring(0, first.length() - 2),
77                     second.substring(0, second.length() - 2)) + 1);
78         }
79
80         return result;
81     }
82
83
84     public static int DamerauLevenshtein(String first, String second) {
85         int n = first.length();
86         int m = second.length();
87
88         int[][] matrix = new int[n + 1][m + 1];
89         matrix[0][0] = 0;
90         for (int i = 1; i < n + 1; i++) {
91             matrix[i][0] = i;
92         }
93         for (int i = 1; i < m + 1; i++) {
94             matrix[0][i] = i;
95         }
96
97         for (int i = 1; i < n + 1; i++) {
98             for (int j = 1; j < m + 1; j++) {
99                 int cost = (first.charAt(i - 1) == second.charAt(j - 1)) ? 0 : 1;
100                 matrix[i][j] = Math.min(
101                     matrix[i - 1][j] + 1,
102                     matrix[i][j - 1] + 1,
103                     matrix[i - 1][j - 1] + cost
104                 );
105             }
106         }
107         return matrix[n][m];
108     }

```

```

94     for (int j = 1; j < m + 1; j++) {
95         matrix[i][j] = Math.min(
96             Math.min(
97                 matrix[i - 1][j - 1] + getIntBoolean(first.charAt(i - 1),
98                     second.charAt(j - 1)),
99                 matrix[i - 1][j] + 1),
100             matrix[i][j - 1] + 1
101         );
102         if (i > 1 && j > 1 && first.charAt(i - 1) == second.charAt(j - 2) &&
103             first.charAt(i - 2) == second.charAt(j - 1)) {
104             matrix[i][j] = Math.min(matrix[i][j], matrix[i - 2][j - 2] + 1);
105         }
106     }
107
108     return matrix[n][m];
109 }
110 }

```

В таблице 3.1 приведены функциональные тесты для алгоритмов вычисления расстояния Левенштейна и Дамерау — Левенштейна. Все тесты пройдены успешно.

Таблица 3.1 – Функциональные тесты

Строка 1	Строка 2	Ожидаемый результат	
		Левенштейн	Дамерау — Левенштейн
cvd	dvc	2	2
mother	money	3	3
just	jest	1	1
turnover	turnovre	2	1
member	morning	6	6
death	health	2	2

## Вывод

В данном разделе были рассмотрены листинги кода, был выбран язык программирования и среда разработки, а также было произведено функциональное тестирование. Сравнивая листинги программ видно, что написание рекуррентных подпрограмм значительно проще, чем матричных.

## 4 Исследовательская часть

### 4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялось тестирование:

- Операционная система: Windows 10.
- Память: 16 GiB.
- Процессор: Intel(R) Core(TM) i7-4700HQ CPU @ 2.40GHz.

### 4.2 Временные характеристики

Для сравнения времени выполнения программ брались строки длиной [10, 20, 30, 50, 100, 200] символов. Для получения точных временных характеристик замеры времени прогонялись 500 раз.

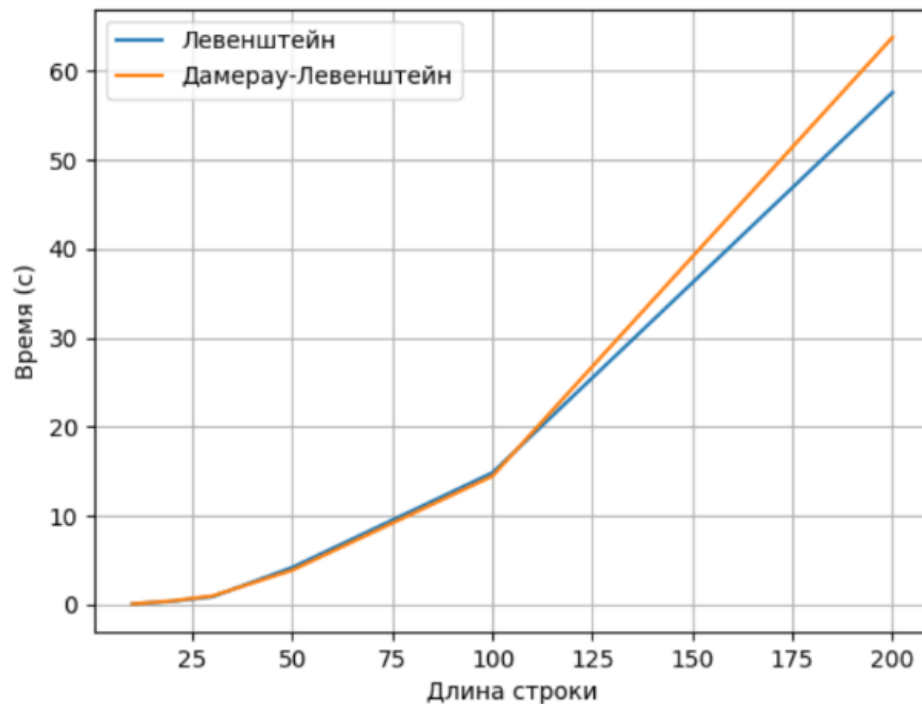


Рисунок 4.1 – Сравнение времени работы алгоритма поиска расстояния Левенштейна и Дамерау-Левенштейна

Из рисунка 4.1 видно, что при короткой длине разница по времени минимальна, при увеличении длины алгоритм поиска расстояния Левенштейна с небольшим опережением вырывается вперед.

### 4.3 Характеристики по памяти

Алгоритмы нахождения расстояний Левенштейна и Дамерау — Левенштейна не отличаются друг от друга с точки зрения использования памяти, следовательно, достаточно рассмотреть лишь разницу рекурсивной и матричной реализаций этих алгоритмов.

Максимальная глубина стека вызовов при рекурсивной реализации равна сумме длин входящих строк, при этом для каждого вызова рекурсии в моей реализации требуется:

- переменная типа `int`, в моем случае: 4 байта;
- 2 аргумента типа строка:  $2 \cdot 24 = 48$  байт;
- адрес возврата: 8 байт;
- место для записи возвращаемого функцией значения: 8 байт.

Таким образом получается, что при обычной рекурсии на один вызов требуется (4.1):

$$M_{percall} = 4 + 48 + 8 + 8 = 68 \quad (4.1)$$

Следовательно память, расходуемая в момент, когда стек вызовов максимален, равна (4.2):

$$M_{recursive} = 80 \cdot depth \quad (4.2)$$

где  $depth$  - максимальная глубина стека вызовов, которая равна (4.3):

$$depth = |S_1| + |S_2| \quad (4.3)$$

где  $S_1, S_2$  - строки.



Память, требуемая для при итеративной реализации, состоит из следующего:

- 2 локальные переменные типа `int`, в моем случае:  $2 \cdot 4 = 8$  байт;
- 2 аргумента типа строка:  $2 \cdot 24 = 48$  байт;
- адрес возврата: 8 байт;
- место для записи возвращаемого функцией значения: 8 байт;
- матрица:  $M_{Matrix}$  размером  $4 \cdot (n + 1) \cdot (m + 1)$ .

Таким образом общая расходуемая память итеративных алгоритмов (4.4):

$$M_{iter} = M_{Matrix} + 72 \quad (4.4)$$

## Вывод

В данном разделе было произведено сравнение количества затраченного времени и памяти вышеизложенных алгоритмов. Самым быстрым оказался матричный алгоритм нахождения расстояния Левенштейна.

# Заключение

Алгоритмы Левенштейна и Дамерау-Левенштейна являются самыми популярными алгоритмами, которые помогают найти редакционное расстояние.

В ходе выполнения лабораторной работы была проделана следующая работа:

- были теоретически изучены алгоритмы нахождения расстояний Левенштейна и Дамерау-Левенштейна;
- для некоторых реализаций были применены методы динамического программирования, что позволило сделать алгоритмы быстрее;
- были практически реализованы алгоритмы в 2 вариантах: рекурсивном и итеративном;
- на основе полученных в ходе экспериментов данных были сделаны выводы по поводу эффективности всех реализованных алгоритмов;
- был подготовлен отчет по ЛР.

# Литература

- [1] Двоичные коды с исправлением выпадений, вставок и замещений символов — 1965, V. 163. — P. 845–848.
- [2] Windows. - Microsoft, 1985. <https://www.microsoft.com/ru-ru/windows>.