

Отчет по лабораторной работе №1  
по курсу "Анализ алгоритмов"  
по теме "Расстояние Левенштейна"

Студент: Барсуков Н.М. ИУ7-56  
Преподаватель: Волкова Л.Л., Строганов Ю.В.

# Содержание

<b>1 Аналитическая часть</b>	<b>3</b>
1.0.1 Вывод . . . . .	4
<b>2 Конструкторский раздел</b>	<b>5</b>
2.1 Алгоритм . . . . .	5
2.2 Структура программы . . . . .	6
2.3 Вывод . . . . .	7
<b>3 Технологический раздел</b>	<b>8</b>
3.1 Требования . . . . .	8
3.2 Выбор языка и среды разработки . . . . .	8
3.3 Интерфейс . . . . .	8
3.4 Листинг . . . . .	9
3.5 Вывод . . . . .	12
<b>4 Исследовательский раздел</b>	<b>13</b>
4.0.1 Железо . . . . .	13
4.1 Сравнение . . . . .	13
4.1.1 Общие результаты замеров . . . . .	13
4.1.2 Детальное сравнение . . . . .	14
4.2 Вывод . . . . .	16
<b>Список литературы</b>	<b>18</b>

# Введение

Расстояние Левенштейна (также редакционное расстояние или дистанция редактирования) между двумя последовательностями символов. В теории информации и компьютерной лингвистике — это минимальное количество операций вставки одного символа, удаления одного символа и замены одного символа на другой, необходимых для превращения одной строки в другую.

Расстояние Левенштейна и его обобщения активно применяются [4]:

1. для исправления ошибок в слове (в поисковых системах, базы данных, при вводе текста, при автоматическом распознавании отсканированного текста и речи);
2. в биоинформатике для сравнения генов, хромосом, белков.

# 1 Аналитическая часть

Цель: Изучить, реализовать и сравнить алгоритмы вычисления редакционного расстояния между 2 строками

Для этого необходимо выполнить следующие пункты:

1. Исследовать алгоритмы вычисления редакционного расстояния;
2. Реализовать следующие алгоритмы:
  - (a) Итеративный;
  - (b) Рекурсивный;
  - (c) Модифицированный итеративный.
3. Сделать замеры выбранных алгоритмов;
4. Сравнить результаты замеров;
5. Сделать вывод.

В нашей лабораторной работе будет рассматриваться 2 способа реализации данного алгоритма.

1. Итеративный;
2. Рекурсивный.

Все они работают по следующему принципу. Допустим, что существует две строки S1 и S2 над некоторым алфавитом. Длина одной из них - M, второй - N. Для нахождения расстояния Левенштейна между ними D(S1, S2) можно применить следующую формулу (D(S1, S2) == D(M, N)) [4] (1 2):

$$D(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) == 0 \\ \min \begin{cases} D(i, j-1) + 1 \\ D(i-1, j) + 1 \\ D(i-1, j-1) + (S1[i] \neq S2[j]) \end{cases} & \end{cases} \quad (1)$$

$$D(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) == 0 \\ \min \begin{cases} D(i, j-1) + 1 \\ D(i-1, j) + 1 \\ D(i-1, j-1) + (S1[i] \neq S2[j]) \end{cases} & \text{if } i, j > 1 \\ \min \begin{cases} D(i, j-1) + 1 \\ D(i-1, j) + 1 \\ D(i-1, j-1) + (S1[i] \neq S2[j]) \end{cases} & \end{cases} \quad (2)$$

Так же вводится понятие стоимости операций. По умолчанию предполагается, что каждая операция, кроме операции сравнения символом, стоят по 1. Введем собственную стоимость для нашего случая:

1. Вставка в S1 (I) - 1;
2. Удаление из S1 (D) - 1;
3. Замена символа в S1 (R) - 1;
4. Обмен местами 2 соседних символов (S) - 1;
5. Совпадение символов в S1 и S2 - 0.

### **1.0.1 Вывод**

В данном разделе были рассмотрены алгоритмы поиска редакционного расстояния Левенштейн и Левенштейн - Дамерау, который является модифицированной версией первого и вводит дополнительную операцию перестановки двух близких символов

Мы будем реализовывать 2





Рис. 2: Блок схема рекурсивного алгоритма

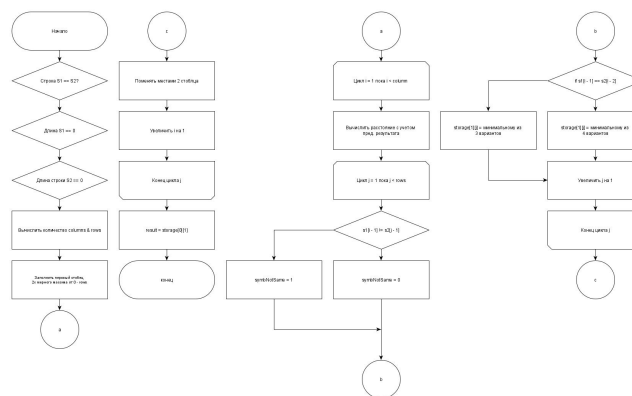


Рис. 3: Блок схема мод. алгоритма

## 2.2 Структура программы

Программа обладает следующей структурой:

1. приветственное окно в котором пользователь знакомят с программой и просят выбрать режим работы. Режимы следующие:
  - (a) поиск редакционного расстояния на 2 строках введенных пользователем;
  - (b) поиска редакционного расстояния на случайных строках длины указанной пользователем.

## 2.3 Вывод

В данном разделе были представлены схемы итеративного (матричного), рекурсивного и модифицированного Левенштейна и структура программы



## 3 Технологический раздел

В данном разделе будут описаны требования к программному продукту и инструментарий реализации. Приведен интерфейс программы и листинг кодов алгоритмов

### 3.1 Требования

Требования к данному программному продукту следующие:

1. программа должна корректно работать;
2. возвращать редакционное расстояние между 2 строк длиной до 10000 символов;
3. не должна обрабатывать некорректный ввод, ответственность ложиться на пользователя

### 3.2 Выбор языка и среды разработки

Для решения данной поставленной задачи, мной был выбран язык с++ стандарта c11 по причине использования структурного подхода и его мультиплатформенности.

Так же мной используются среда разработки под названием Visual Studio 2019 по причине удобства отладки и функциональности. И так же желая изучить данную среду по лучше.

### 3.3 Интерфейс

Интерфейс представляет из себя простую консоль в котором пользователь взаимодействует с помощью ввода команд. Такой тип взаимодействия выбран, по причине простоты разработки и удобства тестирования программы (4, 5).

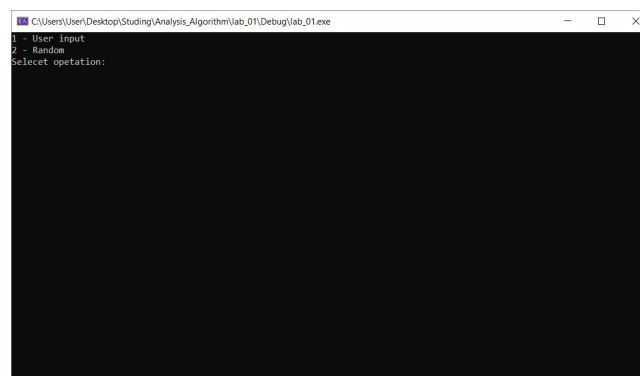


Рис. 4: Примера работы

```
Консоль отладки Microsoft Visual Studio
2 - Random
Select operation: 2
Input len of strings (1 - 10000): 10
Input count of runs: (1-1000): 2
Running test of LEV...

### LEV ###
lev: 10
Time: 1.792e-05

Running test of REC...

### LEV REC ###
pydhsuylk
strp7a39
lev: 10
Time: 0.648904

Running test of MOD...

### LEV MOD ###
lev: 5
Time: 1.8773e-05

C:\Users\User\Desktop\Studing\Analysis_Algorithm\lab_01\Debug\lab_01.exe (процесс 3212) завершает работу с кодом 0.
Чтобы автоматически закрывать консоль при остановке отладки, установите параметр "Сервис" -> "Параметры" -> "Отладка" ->
"Автоматически закрыть консоль при остановке отладки".
Чтобы закрыть это окно, нажмите любую клавишу...
```

Рис. 5: Примера работы

### 3.4 Листинг

Ниже будут представлены листинги итеративного (1), рекурсивного (2), модифицированного (3)

Листинг 1: Итеративный алгоритм

```
1 int levenstain3(char* str1, char* str2) {
2
3     if (strcmp(str1, str2) == 0) {
4         return 0;
5     }
6
7     if (strlen(str1) == 0) {
8         return strlen(str2);
9     }
10
11    if (strlen(str2) == 0) {
12        return strlen(str1);
13    }
14
15    unsigned columns = strlen(str1) + 1;
16    unsigned rows = strlen(str2) + 1;
17
18    int* storage[2];
19    storage[0] = new int[rows];
20    storage[1] = new int[rows];
21
22    for (unsigned i = 0; i < rows; i++) {
23        storage[0][i] = i;
24    }
25
26    for (unsigned i = 1; i < columns; i++) {
27        storage[1][0] = storage[0][0] + 1;
28        for (unsigned j = 1; j < rows; j++) {
29            bool symbNotSame = false;
```

```

30         if (str1[i - 1] != str2[j - 1]) {
31             symbNotSame = true;
32         }
33
34         storage[1][j] = min(min(
35             storage[0][j] + 1,
36             storage[1][j - 1] +
37                 1),
38             storage[0][j - 1] +
39                 symbNotSame);
40     }
41     swap(storage[0], storage[1]);
42 }
43
44 int result = storage[0][rows - 1];
45
46 delete storage[0];
47 delete storage[1];
48
49 return result;
50 }

```

Листинг 2: Рекурсивный алгоритм

```

1 int levenstainRec3(char* str1, char* str2) {
2
3     if (strcmp(str1, str2) == 0) {
4         return 0;
5     }
6
7     if (strlen(str1) == 0) {
8         return strlen(str2);
9     }
10
11    if (strlen(str2) == 0) {
12        return strlen(str1);
13    }
14
15
16    bool symbNotSame = false;
17    if (*str1 != *str2) {
18        symbNotSame = true;
19    }
20
21    int result = min(min(levenstainRec3(str1 + 1, str2)
22        + 1, levenstainRec3(str1, str2 + 1) + 1),
23        levenstainRec3(str1 + 1, str2 + 1) +
24            symbNotSame);

```

```

24     return result;
25 }

```

Листинг 3: Левенштейн модифицированный

```

1  int levenstain4(char* str1, char* str2) {
2
3      unsigned lenstr1 = strlen(str1);
4      unsigned lenstr2 = strlen(str2);
5
6      if (lenstr1 < 2 || lenstr2 < 2)
7          return levenstain3(str1, str2);
8
9      unsigned columns = lenstr1 + 1;
10     unsigned rows = lenstr2 + 1;
11
12     int* storage[3];
13     for (int i = 0; i < 3; i++)
14         storage[i] = new int[rows];
15
16     for (unsigned i = 0; i < rows; i++) {
17         storage[0][i] = i;
18     }
19
20     storage[1][0] = 1;
21     for (unsigned i = 1; i < rows; i++) {
22         bool symbNotSame = false;
23         if (str1[i - 1] != str2[0]) {
24             symbNotSame = true;
25         }
26         storage[1][i] = min(min(storage[1][i - 1] + 1,
27                                 storage[0][i] + 1), storage[0][i - 1] +
28                                 symbNotSame);
29     }
30
31     for (unsigned i = 2; i < columns; i++) {
32         storage[2][0] = storage[1][0] + 1;
33         bool symbNotSame = false;
34         if (str1[1] != str2[i - 1]) {
35             symbNotSame = true;
36         }
37         storage[2][1] = min(min(storage[2][0] + 1,
38                                 storage[1][1] + 1), storage[1][0] +
39                                 symbNotSame);
40     }
41
42     for (unsigned j = 2; j < rows; j++) {
43         symbNotSame = false;
44         if (str1[j] != str2[j - 1]) {
45             symbNotSame = true;
46         }
47         storage[2][j] = min(min(storage[2][j - 1] + 1,
48                                 storage[1][j] + 1), storage[1][j - 1] +
49                                 symbNotSame);
50     }
51
52     return storage[2][rows - 1];
53 }

```

```

41         if (str2[j - 1] != str1[i - 1]) {
42             symbNotSame = true;
43         }
44
45         int tmp = min(min(min(
46             storage[1][j] + 1, storage[2][j - 1] +
47                 1),
48             storage[1][j - 1] + symbNotSame),
49             storage[0][j - 2] + 1);
50
51         storage[2][j] = tmp;
52     }
53
54     swap(storage[0], storage[1]);
55     swap(storage[1], storage[2]);
56 }
57
58 int result = storage[1][rows - 1];
59
60 delete storage[0];
61 delete storage[1];
62 delete storage[2];
63
64 return result;
65 }

```

### 3.5 Вывод

В данном разделе вы ввели требования по программе, выбрали среду и язык разработки, так же определились с будущим интерфейсом нашей программы

## 4 Исследовательский раздел

В данном разделе представлена конфигурация компьютера на котором проводилось тестирование. Результаты замеров и сравнение алгоритмов

### 4.0.1 Железо

1. Компьютер:
  - (a) Тип компьютера Компьютер с ACPI на базе x64;
  - (b) Операционная система Microsoft Windows 10 Pro.
2. Системная плата:
  - (a) тип ЦП DualCore Intel Core i5-6200U, 2700 MHz (27 x 100);
  - (b) системная плата HP 8079;
  - (c) чипсет системной платы Intel Sunrise Point-LP, Intel Skylake-U;
  - (d) системная память 8072 МБ (DDR4 SDRAM).

## 4.1 Сравнение

### 4.1.1 Общие результаты замеров

Таблица 1(1) содержит в себе общие результаты замеров времени, необходимого для вычисления редакционного расстояния между 2 строками равной длины. Считается что слова полностью различны.

Условия замеров:

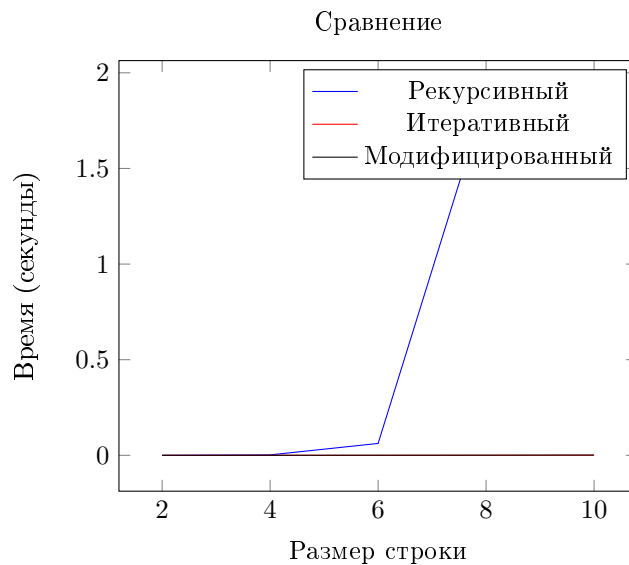
1. Время
  - (a) Время измерений: секунды;
  - (b) Берется суммарное время за все итерации
2. Количество итераций алгоритма: 100;
3. Символы в строках полностью различны

Таблица 1: Общие результаты замеров на 100 итераций

LEN	Рекурсивный	Итеративный	Модифицированный
2	0.000105387	0.000169814	0.000147626
4	0.00224298	0.000243199	0.0003136
6	0.0621811	0.000329813	0.000431786
8	1.87619	0.000616959	0.000716373
10	56.2996	0.000702292	0.00151893

Таблица 2: Таблица сравнения времени

Длина слов	Сравнение с Итеративным	Сравнение с Модифицированным
2	-6.442700000000001e-05	-4.2239000000000014e-05
4	0.001999781	0.0019293799
6	0.061851287	0.061749314
8	1.875573041	1.875473627
10	56.298897708	56.29808



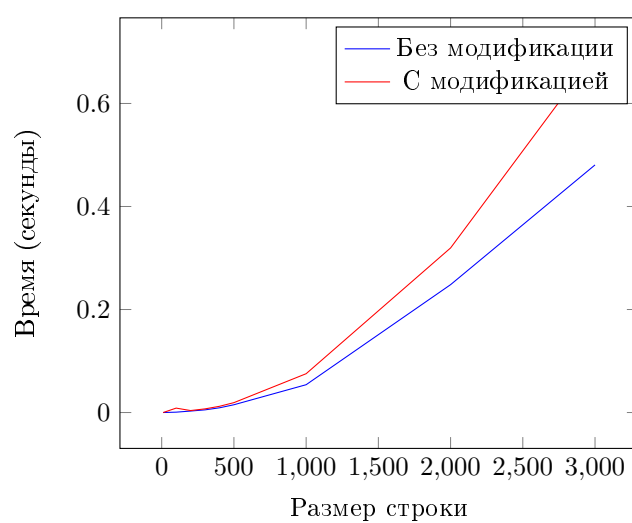
В таблице (2) находятся результаты сравнения рекурсивной с другими реализациями алгоритма. Как мы можем заметить на строках длины 2 рекурсивный справился быстрее на  $-6.442700000000001e-05$  итеративного и на  $-4.2239000000000014e-05$  быстрее модифицированного, но начиная с длины строк равным 4 и больше он начнет проигрывать по времени. В следствии этого дальнейшее детальное сравнение рекурсивного и другими проводится не будет.

#### 4.1.2 Детальное сравнение

Сравним детально время работы классического Левенштейна и модифицированного (с операцией перестановки). На каждый метод берется среднее время 1000 прогонов.

L	Итеративный	Модифицированный
10	6.322e-06	1.112e-05
100	0.000724276	0.00859501
200	0.00286479	0.00383968
300	0.00517626	0.00711481
400	0.0091726	0.0121067
500	0.0151016	0.0191977
1000	0.0540165	0.0753726
2000	0.248242	0.319514
3000	0.480696	0.696582

Сравнение





## 4.2 Вывод

В данном разделе были проведены замеры времени необходимого для вычисления редакционного расстояния трех различных алгоритмов зависимости от:

1. Длины слов;
2. От их редакционного расстояния

Рекурсивный алгоритм при более простой реализации работает чрезвычайно долго. На пример: для строк длины 8 при суммарном времени 100 прогонов, он работает 1.87 дольше итеративного и модифицированного, что делает его использование нецелесообразным. Итеративный алгоритм и его модификация значительно превосходит его по эффективности.

Не смотря схожую асимпотику Левенштейна и Дамерау - Левенштейн, второй работает в среднем медленнее из более сложной внутренней реализации

## Заключение

В ходе работы было проведено сравнение алгоритмов поиска расстояния Левенштейна (рекурсивной и итеративной реализации) и Дамерау-Левенштейна (итеративной реализации). Были исследованы зависимости времени выполнения программ, реализующих данные алгоритмы, от искомого расстояния и от размеров строк для случаев с одинаковой длиной исходных строк и случая, когда одна из строк значительно меньше другой. В ходе исследования были сделаны следующие выводы:

1. Рекурсивная реализация алгоритма Левенштейна выполняется за приемлемое время лишь в случаях, когда размер одной из строк до 10-15 символов. Среднее время составляет 0.5 секунды, что в разы дольше модифицированного и тем более итеративного.
2. 2) Итеративные реализации алгоритмов поиска расстояний Дамерау-Левенштейна и Левенштейна имеют схожую асимптотику, но алгоритм поиска расстояния Дамерау-Левенштейна из-за более сложной внутренней логики в среднем работает медленнее.

## Список литературы

- [1] Дж. Макконнелл. Анализ алгоритмов. Активный обучающий подход.- М.:Техносфера, 2009.
- [2] Методика идентификации пассажира по установленным данным В. М . Черненко , Ю. Е . Гапанюк [Электронный ресурс]. – Режим доступа: <http://www.engjournal.ru/articles/89/89.pdf>, свободный – (07.10.2019)
- [3] Библиография в LaTeX [Электронный ресурс]. – Режим доступа: <https://habr.com/ru/post/114997/>, свободный – (26.09.2019) <https://habr.com/ru/post/114997/>
- [4] Нечеткий поиск, расстояние левенштейна алгоритм [Электронный ресурс]. – Режим доступа: <https://steptosleep.ru/antananarivo-106/>, свободный – (01.10.2019)

□