

МГТУ им. БАУМАНА

ЛАБОРАТОРНАЯ РАБОТА №1

По курсу: "АНАЛИЗ АЛГОРИТМОВ"

## Расстояние Левенштейна

Работу выполнил: Ерохин Никита ИУ7-51Б

Преподаватели: Волкова Л.Л.

*Москва, 2020*

# Оглавление

<b>Введение</b>	<b>2</b>
<b>1 Аналитическая часть</b>	<b>3</b>
1.1 Описание задачи . . . . .	3
1.2 Цель работы . . . . .	4
<b>2 Конструкторская часть</b>	<b>6</b>
2.1 Схемы алгоритмов . . . . .	6
<b>3 Технологическая часть</b>	<b>11</b>
3.1 Выбор ЯП . . . . .	11
3.2 Сведения о модулях программы . . . . .	11
3.3 Тестирование программы по методу черного ящика . . . . .	15
<b>4 Исследовательская часть</b>	<b>19</b>
4.1 Сравнительный анализ на основе замеров времени работы алгоритмов . . . . .	19
4.2 Сравнительный анализ ёмкостной характеристики алгорит- мов . . . . .	23
<b>Заключение</b>	<b>27</b>
<b>Список использованных источников</b>	<b>28</b>

# Введение

**Расстояние Левенштейна** - минимальное количество операций вставки, удаления или замены символа, необходимых для превращения одной строки в другую.

Поиск расстояний Левенштейна применяется в теории информации и компьютерной лингвистике для следующих задач:

- сравнения текстовых файлов утилитой diff;
- исправления ошибок в слове;
- для сравнения генов, хромосом и белков в биоинформатике.

Впервые задачу формализовал советский математик Владимир Левенштейн в 1965 году при изучении последовательностей. Более общую задачу для произвольного алфавита связали с его именем. Большой вклад в изучение вопроса внёс Дэн Гасфилд.

**Расстояние Дameraу-Левенштейна** - Эта вариация вносит в определение расстояния Левенштейна еще одно правило — транспозиция (перестановка) двух соседних букв также учитывается как одна операция, наряду со вставками, удалениями и заменами.

# Глава 1

## Аналитическая часть

### 1.1 Описание задачи

Расстояние Левинштейна между строками  $S_1$  и  $S_2$  - минимальное количество операций над этими строками, необходимых для преобразования одной строки в другую.

При нахождении расстояния Дамерау — Левенштейна добавляется операция транспозиции (перестановки соседних символов).

**Операции обозначаются так:**

1. I — вставить. Штраф - 1;
2. D — удалить. Штраф - 1;
3. R — заменить. Штраф - 1;
4. M - совпадение. Штраф - 0.

Пусть  $S_1$  и  $S_2$  — две строки длиной M и N соответственно над некоторым алфавитом. Тогда расстояние Левенштейна можно получить из следующей рекуррентной формулы:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min( & j > 0, i > 0 \\ D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[j]) \\ ), \end{cases}$$

где  $m(a, b)$  равна нулю при  $a = b$ , и единице в противном случае;  $\min\{a, b, c\}$  возвращает наименьший из аргументов.

Расстояние Дамрау-Левенштейна можно получить из следующей рекуррентной формулы:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & i > 0, j = 0 \\ j, & i = 0, j > 0 \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[i]), \\ D(i - 2, j - 2) + m(S_1[i], S_2[i]), \end{cases} & \begin{array}{l} \text{, если } i, j > 0 \\ \text{и } S_1[i] = S_2[j - 1] \\ \text{и } S_1[i - 1] = S_2[j] \end{array} \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[i]), \end{cases} & \text{, иначе} \end{cases}$$

## 1.2 Цель работы

Целью лабораторной работы является изучение и применение алгоритмов поиска расстояния Левенштейна и Дамрау-Левенштейна, а также получение практических навыков реализации следующих алгоритмов:

1. поиск расстояния Левенштейна рекурсивным способом;
2. поиск расстояния Левенштейна матричным способом;
3. поиск расстояния Левенштейна матрично-рекурсивным способом;

4. поиск расстояния Дамерау-Левенштейна табличным способом.

## Глава 2

# Конструкторская часть

### Требования к программе:

- ввод двух строк;
- при введенных двух пустых строках - корректный вывод, программа не должна завершиться аварийно;
- вывод редакционного расстояния;
- осуществление замера процессорного времени работы алгоритмов.

## 2.1 Схемы алгоритмов

В данной части будут рассмотрены схемы алгоритмов поиска расстояния Левенштейна:

- схема матричного алгоритма нахождения расстояния Левенштейна (рис. 2.1);
- схема рекурсивного алгоритма нахождения расстояния Левенштейна (рис. 2.3);
- схема матричного-рекурсивного алгоритма нахождения расстояния Левенштейна (рис. ??);
- схема матричного алгоритма нахождения расстояния Дамерау-Левенштейна изображена (рис. 2.4).

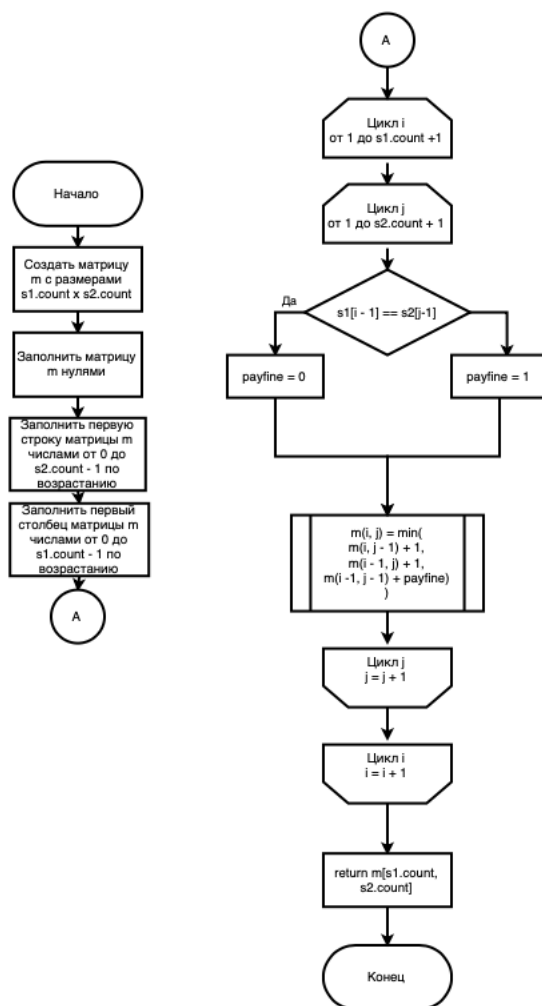


Рис. 2.1: Схема матричного алгоритма нахождения расстояния Левенштейна



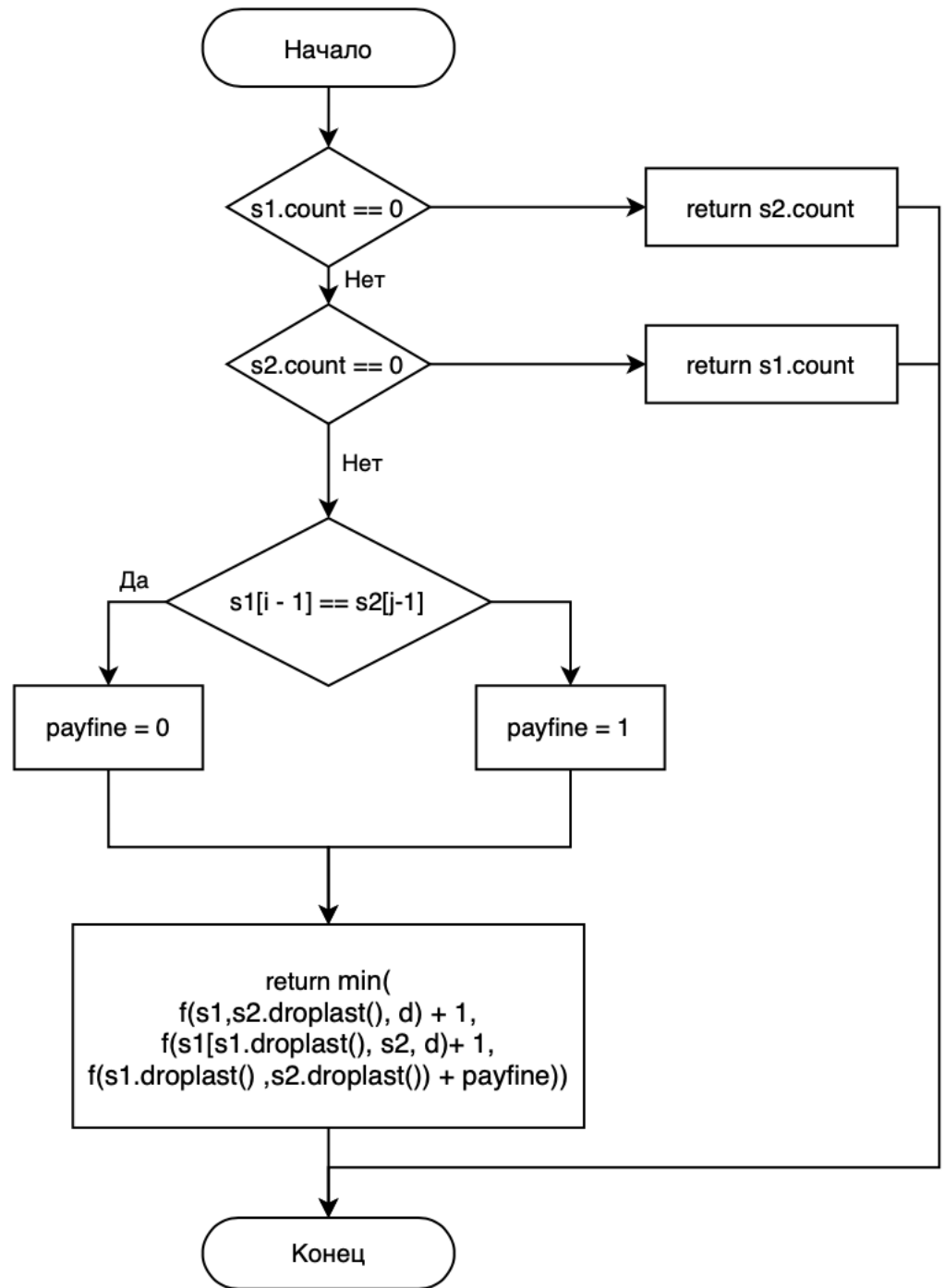


Рис. 2.2: Схема рекурсивного алгоритма нахождения расстояния Левенштейна

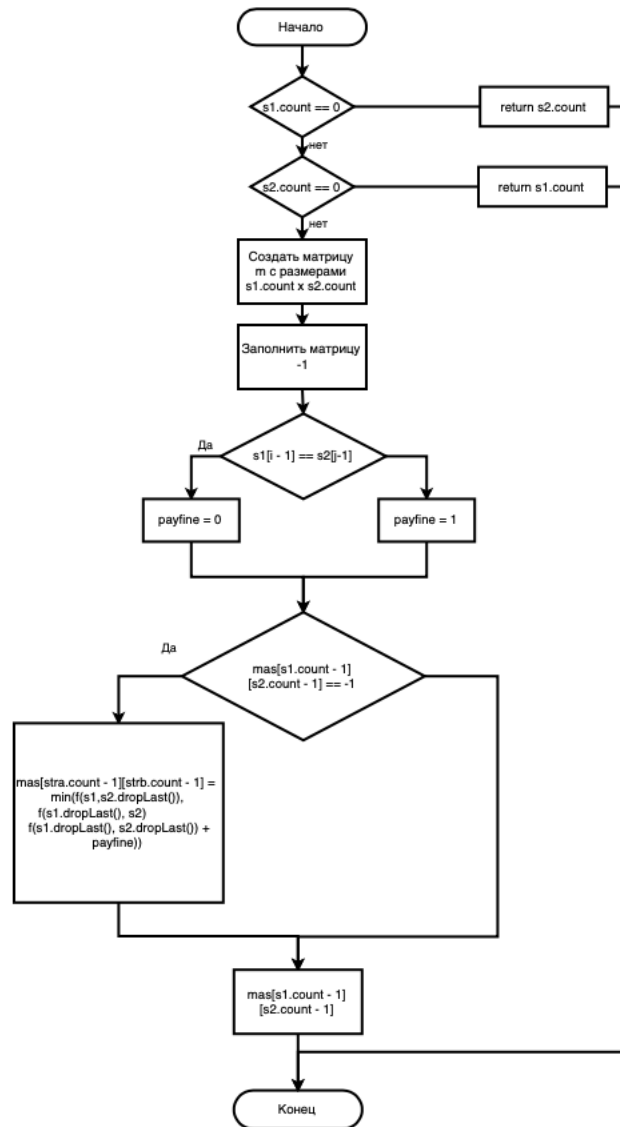


Рис. 2.3: Схема матричного-рекурсивного алгоритма нахождения расстояния Левенштейна

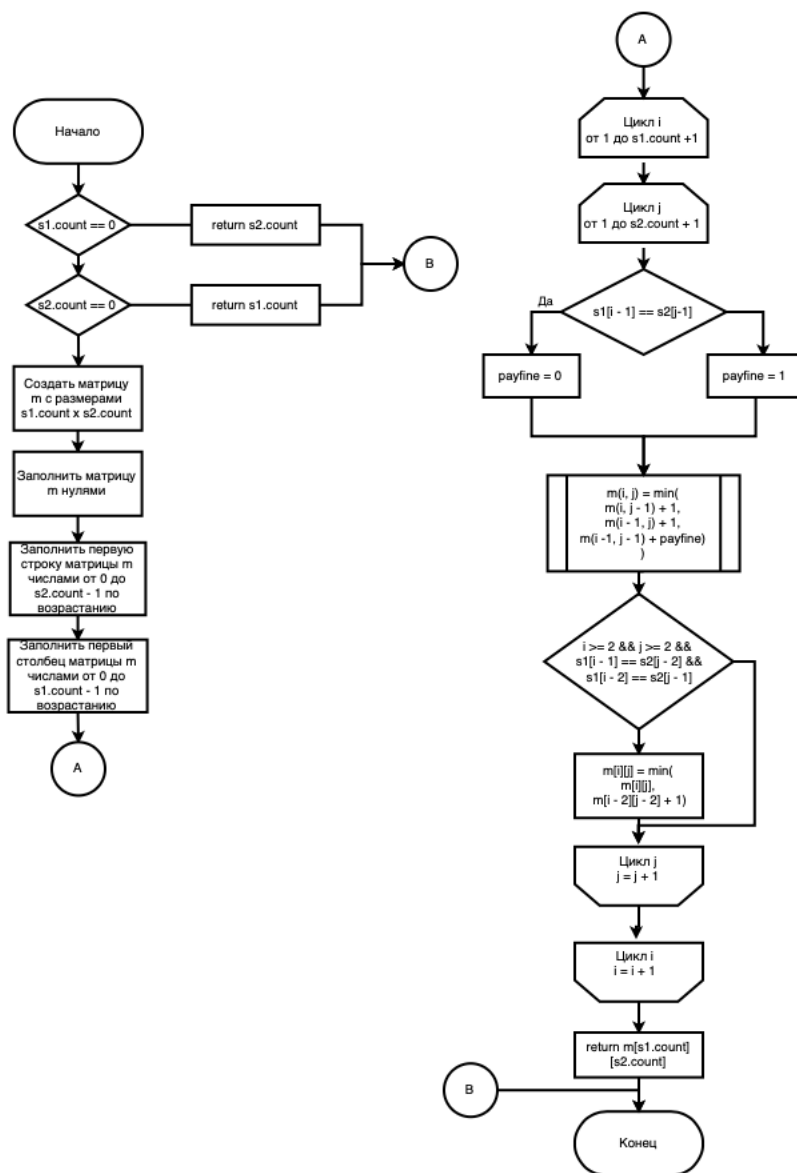


Рис. 2.4: Схема матричного алгоритма нахождения расстояния Дамерау-Левенштейна

## Глава 3

# Технологическая часть

### 3.1 Выбор ЯП

Для реализации программы был выбран язык программирования swift, тк данный язык отличается скоростью компиляции, что значительно сократло время разработки.

### 3.2 Сведения о модулях программы

Программа состоит из:

- main.swift - тесты
- levenstein.swift - реализация представленных алгоритмов

Рекурсивный алгоритм нахождения расстояния Левенштейна изображен на листинге 3.1.

Листинг 3.1: Функция нахождения расстояния Левенштейна рекурсивно

```
1 func levensteinRecursion(_ stra: String, _ strb: String)
  -> Int {
2   if (stra.count == 0) {
3     return strb.count
4   }
5   if (strb.count == 0) {
6     return stra.count
7   }
8   let payfine = stra.last != strb.last ? 1 : 0
9
10
11  return min(levensteinRecursion(stra, String(strb.
    dropLast(1))) + 1, levensteinRecursion(String(stra.
    dropLast(1)), strb) + 1,
12    levensteinRecursion(String(stra.dropLast(1)),
    String(strb.dropLast(1))) + payfine)
13 }
```

Матричный алгоритм нахождения расстояния Левенштейна изображен на листинге 3.2.

Листинг 3.2: Функция нахождения расстояния Левенштейна матрично

```
1 func levensteinMatrix(_ stra: String, _ strb: String) ->
  Int {
2   if stra.count == 0 {
3     return strb.count
4   }
5   if strb.count == 0 {
6     return stra.count
7   }
8
9   var mas = getStartMas(stra.count, strb.count)
10
11  for i in 1 ..< stra.count + 1 {
12    for j in 1..<strb.count + 1 {
13      let payfine = stra[i - 1] != strb[j - 1] ? 1 :
14        0
15
16      mas[i][j] = min(mas[i - 1][j] + 1,
17        mas[i][j - 1] + 1,
18        mas[i - 1][j - 1] + payfine)
19    }
20  }
21  return mas[stra.count][strb.count]
```

Матрично-рекурсивный алгоритм нахождения расстояния Левенштейна изображен на листинге 3.3.

Листинг 3.3: Функция нахождения расстояния Левенштейна матрично-рекурсивно

```
1 func levensteinRecursionMatrix(_ stra: String, _ strb:
  String, _ mas : inout [[Int]]) -> Int {
2   if stra.count == 0 {
3     return strb.count
4   }
5   if strb.count == 0 {
6     return stra.count
7   }
8
9   if mas.count == 0 {
10    mas = getStartMas(stra.count - 1, strb.count - 1,
      -1)
11  }
12
13
14  let payfine = stra[stra.count - 1] != strb[strb.count -
    1] ? 1 : 0
15
16  if (mas[stra.count - 1][strb.count - 1] == -1) {
17    mas[stra.count - 1][strb.count - 1]
18      = min(levensteinRecursionMatrix(stra, String(
        strb.dropLast(1)), &mas) + 1,
19            levensteinRecursionMatrix(String(stra.
        dropLast(1)), strb, &mas) + 1,
20            levensteinRecursionMatrix(String(stra.
        dropLast(1)), String(strb.dropLast(1))
        , &mas) + payfine)
21  }
22  return mas[stra.count - 1][strb.count - 1]
23 }
```

Матричный алгоритм нахождения расстояния Дамерау-Левенштейна изображен на листинге 3.4.

Листинг 3.4: Функция матричного нахождения расстояния Дамерау-Левенштейна

```
1 func damerauLevenstein(_ stra: String, _ strb: String) ->
  Int {
2
3   if stra.count == 0 {
4     return strb.count
5   }
6   if strb.count == 0 {
7     return stra.count
8   }
9
10
11  var mas = getStartMas(stra.count, strb.count)
12
13  for i in 1..<stra.count + 1 {
14    for j in 1..<strb.count + 1 {
15      let payfine = stra[i - 1] != strb[j - 1] ? 1 :
        0
16      mas[i][j] = min(mas[i - 1][j] + 1,
17                      mas[i][j - 1] + 1,
18                      mas[i - 1][j - 1] + payfine)
19      if (i >= 2 && j >= 2 && stra[i - 1] == strb[j -
        2] && stra[i - 2] == strb[j - 1]) {
20        mas[i][j] = min(mas[i][j], mas[i - 2][j -
        2] + 1)
21      }
22    }
23  }
24  return mas[stra.count][strb.count]
25
26 }
```

### 3.3 Тестирование программы по методу черного ящика

В данном разделе будут приведены результаты тестирования программы.



На Рис. 3.1 показано поведение программы при вводе пустых строк.

```
String one:  
String two:  
expect answer: 0  
Levenstein matrix result 0  
Levenstein recursion result 0  
Levenstein recursion-matrix result 0  
Passed
```

Рис. 3.1: Тест при строках размером 0

На Рис. 3.2 показано поведение программы одной пустой и одной непустой строке.

```
String one:  
String two: 12345  
expect answer: 5  
Levenstein matrix result 5  
Levenstein recursion result 5  
Levenstein recursion-matrix result 5  
Passed
```

Рис. 3.2: Тест при одной строке длины 0 и одной непустой строке

На Рис. 3.3 показано поведение программы строках "1234" и "1".

```
String one: 1234
String two: 1
expect answer: 3
Levenstein matrix result 3
Levenstein recursion result 3
Levenstein recursion-matrix result 3
Passed
```

Рис. 3.3: Тест при введенных строках "1234" и "1"

На Рис. 3.4 показано поведение программы строках "*kitten*" и "*mittens*".

```
String one: kitten
String two: mittens
expect answer: 2
Levenstein matrix result 2
Levenstein recursion result 2
Levenstein recursion-matrix result 2
Passed
```

Рис. 3.4: Тест при введенных строках "*kitten*" и "*mittens*"

На Рис. 3.5 показано поведение тестирования функции поиска расстояния Дамерау-Левенштейна для строк "*iphone*" и "*iphnoe*".

```
String one: iphone
String two: iphnoe
expect answer: 1
Damerau-Levenstein matrix result 1
Passed
```

Рис. 3.5: Тестирование функции поиска расстояния Дамерау-Левенштейна для строк "*iphone*" и "*iphnoe*"

На Рис. 3.6 показано поведение тестирования функции поиска расстояния Дамерау-Левенштейна для строк "*supercat*" и "*uppercta*".

```
String one: supercat
String two: uppercta
expect answer: 3
Damerau-Levenstein matrix result 3
Passed
```

Рис. 3.6: Тестирование функции поиска расстояния Дамерау-Левенштейна для строк "*supercat*" и "*uppercta*"

## Глава 4

### Исследовательская часть

#### 4.1 Сравнительный анализ на основе замеров времени работы алгоритмов

Был произведен замер времени работы каждой из реализаций алгоритмов (таблица 4.1).

Введем обозначения:

- `lmatrix` - матричный алгоритм поиска расстояния Левенштейна;
- `lrnomatrix` - рекурсивный алгоритм поиска расстояния Левенштейна;
- `lrmatrix` - матрично-рекурсивный алгоритм поиска расстояния Левенштейна;
- `damerauLevenstein` - матричный алгоритм поиска расстояния Дамерау-Левенштейна.

Таблица 4.1: Время работы алгоритмов (в секундах)

Количество символов	lmatrix	lnomatrix	lrmatrix	damerauLevenstein
3	0.000023	0.000051	0.000038	0.000018
4	0.000019	0.000209	0.000052	0.000021
5	0.000026	0.001001	0.000082	0.000030
6	0.000035	0.004965	0.000106	0.000034
7	0.000048	0.026186	0.000141	0.000045

На основе таблицы был построен график(Рис. 4.1).

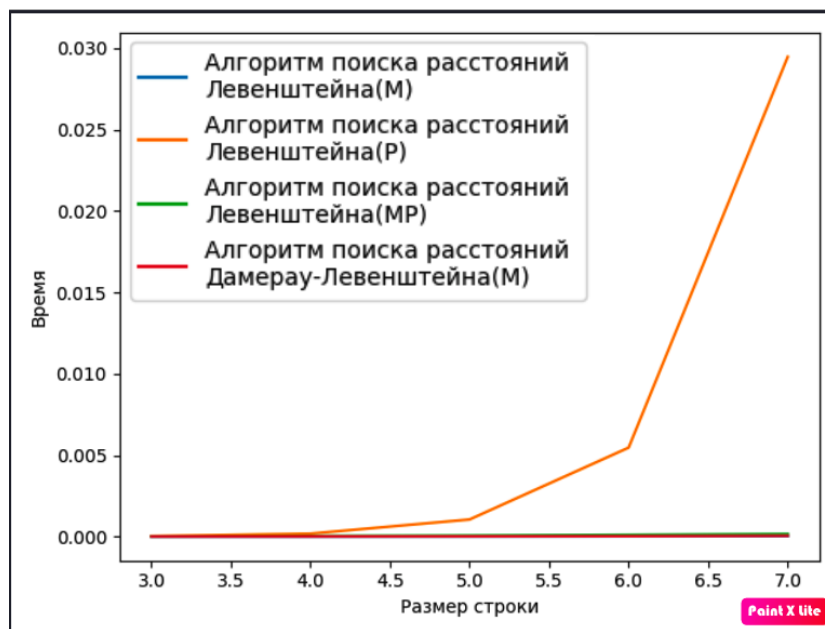


Рис. 4.1: Сравнение времени выполнения алгоритмов

Данный график наглядно демонстрирует недостатки рекурсивной реализации алгоритма поиска расстояния Левенштейна. Однако масштаб не позволяет сравнить более эффективные алгоритмы. Далее, на (рис. 4.2) будут рассмотрены 3 графика из 4.

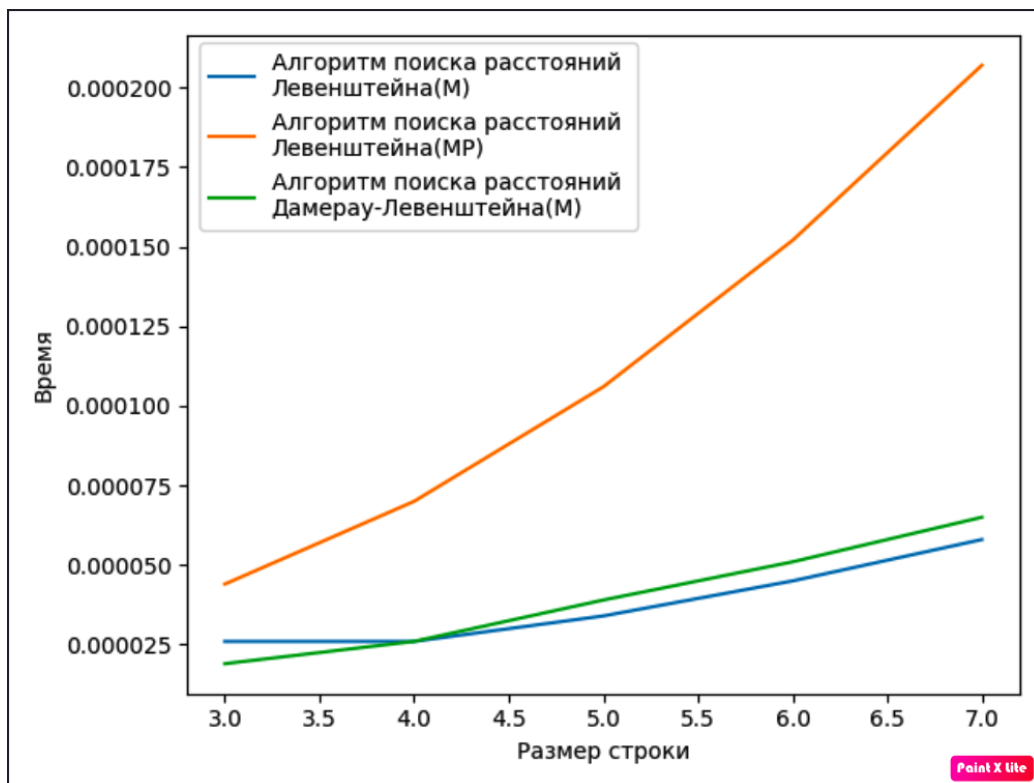


Рис. 4.2: Сравнение времени выполнения алгоритмов Левенштейна(Матричная реализация), Левенштейна(Матрично-рекурсивная реализация), Дамерау-Левенштейна(Матричная реализация)

Из полученных графиков можно сделать вывод, что рекурсивный алгоритм является самым медленным. Но можно уменьшить количество повторных вызовов в рекурсивном алгоритме Левенштейна, если использовать матрицу, в которой все поля изначально -1(бесконечность). Так как не будет повторного вызова. Самым быстрым алгоритмом является матричный алгоритм Левенштейна. Из полученных графиков можно сделать вывод, что рекурсивная реализация алгоритма поиска расстояния Левенштейна является самой медленной. Ее возможно улучшить, уменьшить количество повторных вызовов в рекурсивном алгоритме. Для этого необходимо использовать матрицу, в которой все поля изначально -1. Самым быстрым алгоритмом является матричная реализация алгоритма поиска расстояния Левенштейна.

## 4.2 Сравнительный анализ ёмкостной характеристики алгоритмов

Проведем анализ ёмкостной характеристики у разных структур данных:

Ёмкостная характеристика у различных структур данных приведена в таблице 4.2.

Таблица 4.2: Ёмкостная характеристика у различных структур данных(в байтах)

Структура данных	Память(байты)
Пустая строка	49
Целое число	28
Пустой список	56
Список, хранящий указатель	64



Емкостная характеристика матричного алгоритма Дамерау-Левенштейна приведена в таблице 4.3.

Таблица 4.3: Емкостная характеристика матричного алгоритма Левенштейна

Структура данных	Занимаемая память(байты)
Матрица	$56 + 28 * (n1 + 1) + (n1 + 1) * (56 + 28 * (n2 + 1))$
Счетчики	84
Передача параметров	$2 * [49 + n]$
Вспомогательные переменные	56

Емкостная характеристика матричного алгоритма поиска расстояния Дамерау-Левенштейна приведена в таблице 4.4.

Таблица 4.4: Емкостная характеристика матричного алгоритма Дамерау-Левенштейна

Структура данных	Занимаемая память(байты)
Матрица	$56 + 28 * (n1 + 1) + (n1 + 1) * (56 + 28 * (n2 + 1))$
Счетчики	56
Передача параметров	$2 * [49 + n]$
Вспомогательные переменные	112

Емкостная характеристика рекурсивного алгоритма поиска расстояния Левенштейна приведена в таблице 4.5.

Таблица 4.5: Емкостная характеристика рекурсивного алгоритма поиска расстояния Левенштейна

Структура данных	Занимаемая память(байты)
Передача параметров + глубина рекурсии	$2 * (49 + n) * (n + m)$

Емкостная характеристика рекурсивного алгоритма поиска расстояния Левенштейна приведена в таблице 4.6.

Таблица 4.6: Емкостная характеристика матрично-рекурсивного алгоритма поиска расстояния Левенштейна.

Структура данных	Занимаемая память(байты)
Матрица	$56 + 28 * (n1 + 1) + (n1 + 1) * (56 + 28 * (n2 + 1))$
Передача параметров	$2 * (49 + n) * (n + m)$

Исходя из данных, приведенных выше, можно сделать вывод, что матричная реализация алгоритма поиска расстояния Левенштейна и Дамерау-Левенштейна имеют одинаковые затраты по памяти. Самым неэффекти-

вынм по памяти является матрично-рекурсивный способ Левенштейна, так как содержит в себе матрицу.

## Заключение

В процессе выполнения данной лабораторной работы были реализованы и проанализированы различные реализации алгоритма поиска расстояний Левенштейна. Было проведено тестирование по методу "" и теоретический замер памяти.

Было подтверждено различие в эффективности различных реализаций (в частности, рекурсивной и нерекурсивной и матрично-рекурсивной) алгоритма определения расстояния Левенштейна между строками.

## Список использованных источников

1. Дж. Макконнелл. Анализ Алгоритмов. Активный обучающий подход.- М.: Техносфера, 2009.
3. Расстояние Дамерау–Левенштейна [Электронный ресурс] Режим доступа: <https://elementy.ru/problems/1068> (дата обращения 01.10.20).
4. Вычисление редакционного расстояния [Электронный ресурс]. Режим доступа: <https://habr.com/ru/post/117063/> (дата обращения 30.09.20).