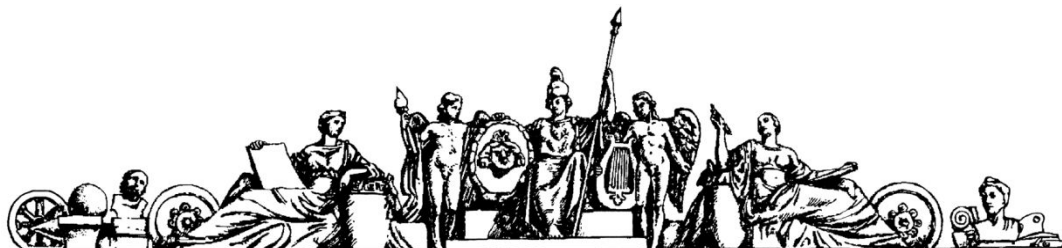


**Министерство образования Российской Федерации**  
**МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ**  
**им. Н.Э. БАУМАНА**



**ОТЧЁТ ПО АНАЛИЗУ АЛГОРИТМОВ**  
к лабораторной работе №1 на тему:

Расстояние Левенштейна. Расстояние Дameraу-Левенштейна

Студент: Зыкин Д.А. ИУ7-53(Б)

Москва 2018

## **Оглавление**

Введение.	3
1 Аналитическая часть	4
1.1 Описание алгоритмов	4
1.1.1 Расстояние Левенштейна	4
1.1.2 Расстояние Дамерау-Левенштейна	5
2 Конструкторский раздел	5
2.1 Разработка алгоритмов	5
2.1.1 Математическое описание	5
2.1.2 Блок-схемы алгоритмов	7
3 Технологический раздел	11
3.1 Требования к программному обеспечению	11
3.2 Средства реализации	11
3.3 Листинг кода	11
4 Исследовательский раздел	13
4.1 Пример работы	13
4.2 Постановка эксперимента	14
Заключение	17

## **Введение.**

Цель работы: изучение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна.

В данной работе требуется изучить и применить метод динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна, а также получить практические навыки реализации алгоритмов.

Данные алгоритмы активно используются при проверке правописания в редакторских программах и в поисковых системах.

В данной работе рассматриваются 3 алгоритма:

- 1) Алгоритм Левенштейна;
- 2) Алгоритм Дамерау-Левенштейна;
- 3) Рекурсивный алгоритм Левенштейна.

# 1 Аналитическая часть

1.1.1 Расстояние Дамерау-Левенштейна, как и метрика Левенштейна, является мерой "схожести" двух строк. Алгоритм его поиска находит применение в реализации нечеткого поиска, а также в биоинформатике (сравнение ДНК), несмотря на то, что изначально алгоритм разрабатывался для сравнения текстов, набранных человеком (Дамерау показал, что 80% человеческих ошибок при наборе текстов составляют перестановки соседних символов, пропуск символа, добавление нового символа, и ошибка в символе. Поэтому метрика Дамерау-Левенштейна часто используется в редакторских программах для проверки правописания).

Наиболее часто применяемой метрикой является расстояние Левенштейна, или расстояние редактирования, алгоритмы вычисления которого можно найти на каждом шагу.

Исходный вариант этого алгоритма имеет временную сложность  $O(mn)$  и потребляет  $O(mn)$  памяти, где  $m$  и  $n$  — длины сравниваемых строк.

Расстояние Левенштейна — минимальное количество действий, необходимых для преобразования одного слова в другое. При вычислении расстояния Левенштейна следует выбирать минимальное количество действий необходимые для полного совпадения слов.

При вычислении расстояния Левенштейна используется понятие цена(штраф).

Всего в алгоритме Левенштейна есть 3 операции.

- 1) замена;
- 2) вставка;
- 3) удаления.

### 1.1.2 Расстояние Дамерау-Левенштейна

Если поиск слова осуществляется в тексте, который набран с клавиатуры, то вместо расстояния Левенштейна используют усовершенствованное расстояние Дамерау – Левенштейна. Исследования Ф. Дамерау показали, что наиболее частая ошибка при наборе слова – перестановка двух соседних букв, транспозиция Т (transposition). В случае одной транспозиции расстояние Левенштейна равно 2. При использовании поправки Дамерау транспозиция принимается за единичное расстояние. При использовании расстояния Дамерау – Левенштейна за единичное расстояние принимают следующие действия: I (insert) – добавление символа; D (delete) – удаление символа; R (replace) – замена символа; Т (transposition) – перестановка двух соседних символов.

## 2 Конструкторский раздел

В этом разделе будут рассмотрены алгоритмы с точки зрения формул и блок-схем их работы.

### 2.1 Разработка алгоритмов

Создание и представление работы алгоритмов.

#### 2.1.1 Математическое описание

Пусть имеются две строки 1 и 2 длиной S1, S2.

Тогда расстояние Левенштейна (1, 2) можно подсчитать по рекуррентной формуле:

$$D(i,j) = \begin{cases} 0; i = 0, j = 0 \\ i; j = 0, i > 0 \\ j; i = 0, j > 0 \\ D(i-1, j-1); S1[i] = S2[j] \\ \min(D(i, j-1) + insertCost, \\ D(i-1, j) + deleteCost, \\ D(i-1, j-1) + replaceCost); j > 0, i > 0, S1[i] \neq S2[j] \end{cases} .$$

Классификация разрешенных операций и штрафы на выполнение операции:

- 1) замена символа = 1;
- 2) вставка символа = 1;
- 3) удаление символа = 1;
- 4) совпадение символа = 0.

В алгоритме Дамерау-Левенштейна добавляется еще одна операция - перестановка символа = 1.

В рекуррентную формулу добавляется еще один член:

$$da, b(i, j) \begin{cases} \max(i, j) \text{ if } \min(i, j) = 0 \\ \min(da, b(i-1, j) + 1, \text{ if } i, j > 1 \text{ and } ai = bj - 1 \text{ and } ai - 1 = bj \\ da, b(i, j-1) + 1, \\ da, b(i-1, j-1) + 1 (ai \neq bj) \\ da, b(i-2, j-2) + 1) \\ \min(da, b(i-1, j) + 1, \\ da, b(i, j-1) + 1, \\ da, b(i-1, j-1) + 1 (ai \neq bj)) \end{cases}$$

2.1.2 Блок-схемы алгоритмов

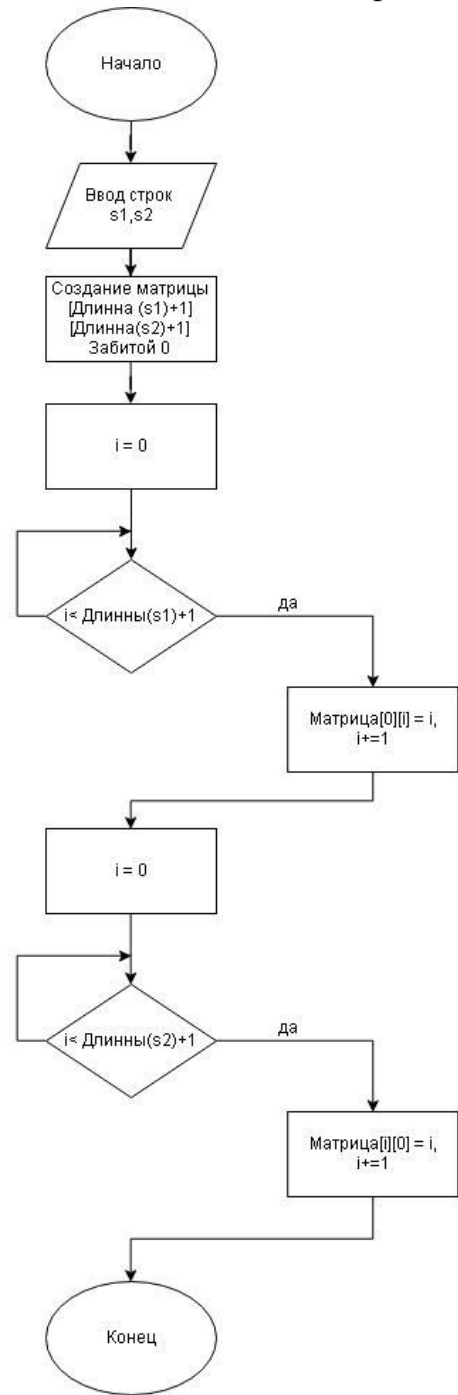


Рисунок 2.1 — процесс инициализации матрицы.

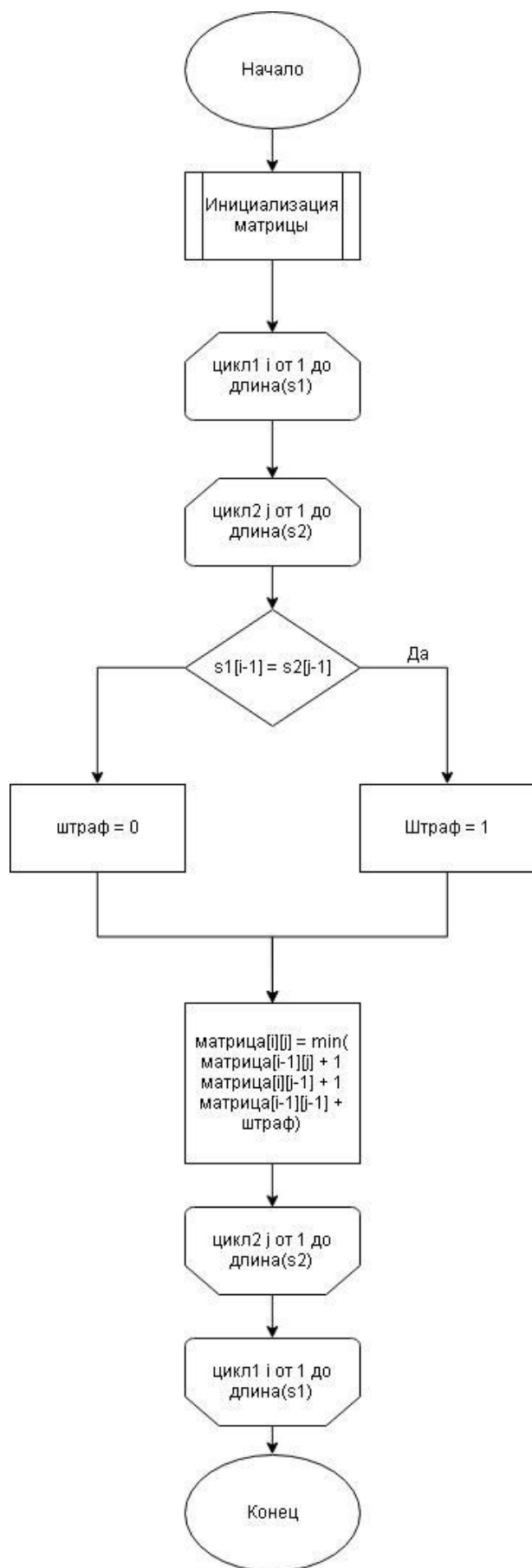


Рисунок 2.2 - алгоритм Левенштейна



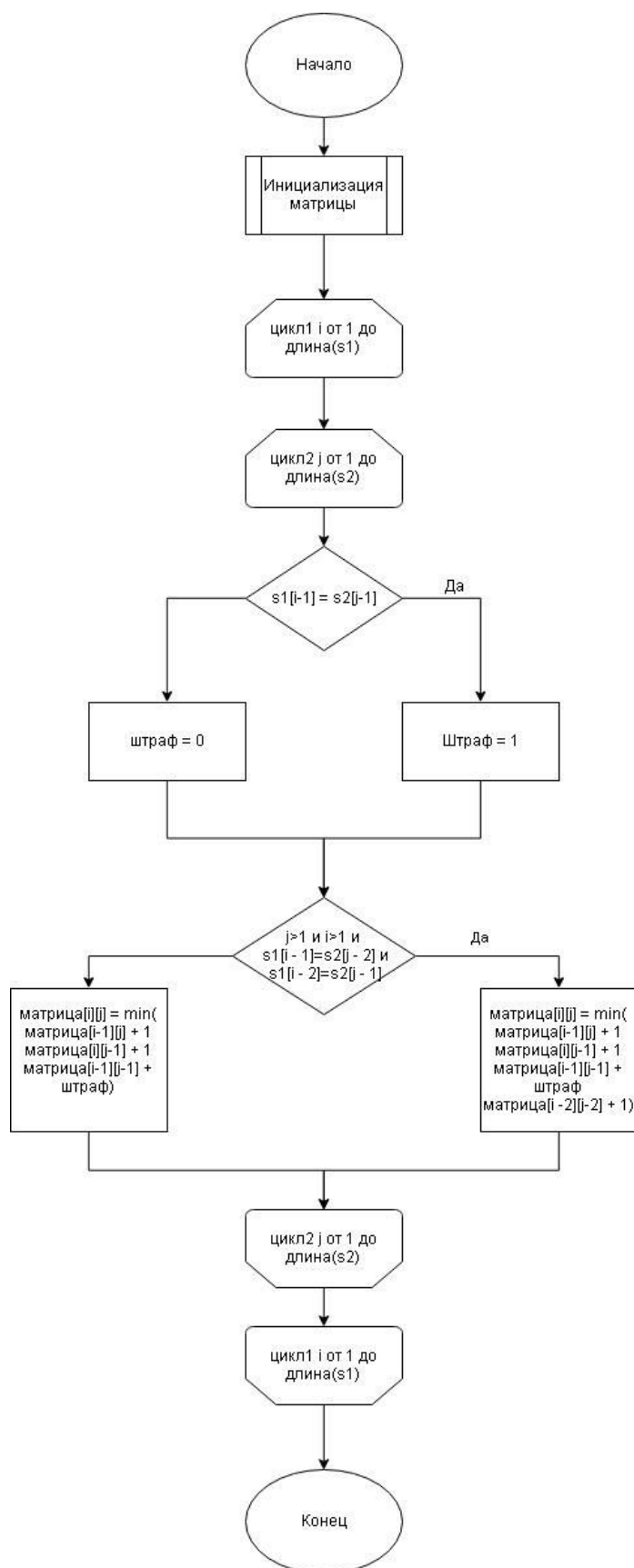


Рисунок 2.3 – алгоритм Дameraу-Левенштейна

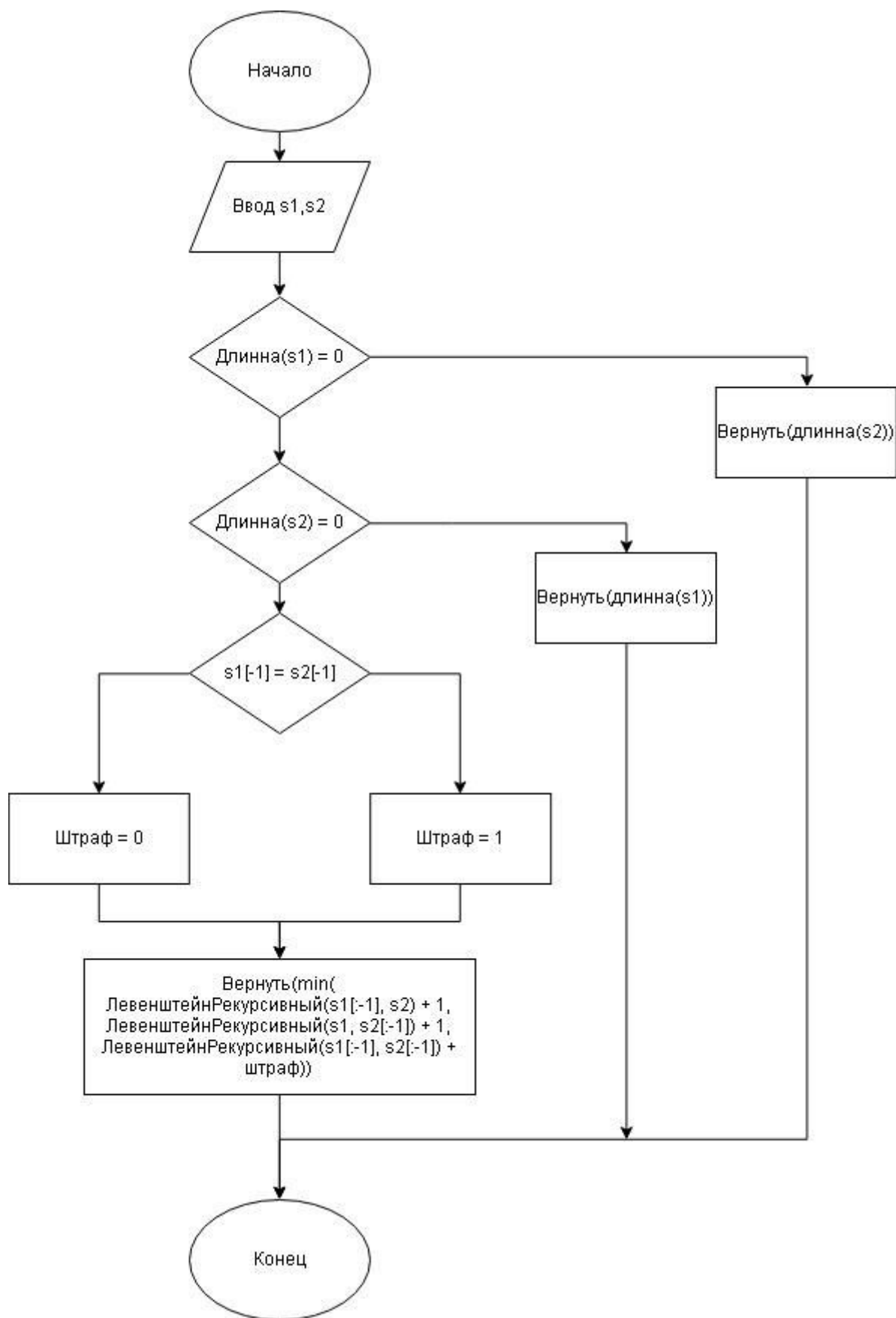


Рисунок 2.4 - алгоритм Левенштейна Рекурсивный

### 3 Технологический раздел

Создание алгоритмов на языке программирование и создание готового продукта.

#### 3.1 Требования к программному обеспечению

Средние требования:

- Python 3.4.3 для 64-битной системы;
- ОС: OS 64-bit Windows 7;
- Процессор: Intel CPU Core i5-2500K 3.3GHz;
- Оперативная память: 8 GB ОЗУ.

#### 3.2 Средства реализации

Данная программа разрабатывалась на языке Python, т.к. python высокоуровневый язык программирования общего назначения, ориентированный на повышение производительности разработчика и читаемости кода. Синтаксис ядра Python минималистичен. В то же время стандартная библиотека включает большой объём полезных функций.

Python поддерживает несколько парадигм программирования, в том числе структурное, объектно-ориентированное, функциональное, императивное и аспектно-ориентированное. Основные архитектурные черты — динамическая типизация, автоматическое управление памятью, полная интроспекция, механизм обработки исключений, поддержка многопоточных вычислений и удобные высокоуровневые структуры данных. Код в Python организовывается в функции и классы, которые могут объединяться в модули (они в свою очередь могут быть объединены в пакеты).

### 3.3 Листинг кода

Алгоритм Левенштейна:

```
def Livenshtein(s1, s2):  
    matrix = []  
    for i in range(len(s2) + 1):  
        matrix.append([0] * (len(s1) + 1))  
    for i in range(1, len(s2)+1 ):  
        matrix[i][0] = i  
    for i in range(1, len(s1)+1):  
        matrix[0][i] = i  
    for i in range (1, len(s1)+1):  
        for j in range (1, len(s2)+1):  
            cost = 0 if s2[j - 1] == s1[i - 1] else 1  
            matrix[j][i] = min(matrix[j-1][i]+1,  
                                matrix[j][i-1]+1,  
                                matrix[j-1][i-1] + cost)
```

Алгоритм Дамерау-Левенштейна:

```
def DamLivenshtein(s1, s2):  
    matrix = []  
    string = s1  
  
    for i in range(len(s2) + 1):  
        matrix.append([0] * (len(s1) + 1))  
    for i in range(1, len(s2)+1 ):  
        matrix[i][0] = i  
    for i in range(1, len(s1)+1):  
        matrix[0][i] = i  
    for i in range (1, len(s1)+1):
```

```

for j in range (1, len(s2)+1):
    cost = 0 if s2[j - 1] == s1[i - 1] else 1
    if j>1 and i>1 and s1[i - 1] == s2[j - 2] and s1[i - 2] == s2[j - 1]:
        matrix[j][i] = min(matrix[j - 1][i] +1,
                            matrix[j][i - 1] +1,
                            matrix[j - 1][i - 1] + cost,
                            matrix[j - 2][i - 2] +1)
    else:
        matrix[j][i] = min(matrix[j - 1][i] +1,
                            matrix[j][i - 1] +1,
                            matrix[j - 1][i - 1] + cost)

```

Рекурсивный алгоритм Левенштейна:

```

def RecLivenshtein(s1, s2):
    if s1 == "":
        return len(s2)
    if s2 == "":
        return len(s1)
    cost = 0 if s1[-1] == s2[-1] else 1

    result = min(RecLivenshtein(s1[:-1], s2) + 1,
                  RecLivenshtein(s1, s2[:-1]) + 1,
                  RecLivenshtein(s1[:-1], s2[:-1]) + cost)

```

## 4 Исследовательский раздел

### 4.1 Пример работы

Левенштейн

<p>п р о л е т</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>э [1, 1, 2, 3, 4, 5, 6]</p> <p>а [2, 2, 2, 3, 4, 5, 6]</p> <p>л [3, 3, 3, 3, 3, 4, 5]</p> <p>е [4, 4, 4, 4, 4, 3, 4]</p> <p>т [5, 5, 5, 5, 5, 4, 3]</p>	<p>к о р о в а</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>к [1, 0, 1, 2, 3, 4, 5]</p> <p>о [2, 1, 0, 1, 2, 3, 4]</p> <p>о [3, 2, 1, 1, 1, 2, 3]</p> <p>р [4, 3, 2, 1, 2, 2, 3]</p> <p>в [5, 4, 3, 2, 2, 2, 3]</p> <p>а [6, 5, 4, 3, 3, 3, 2]</p>	<p>л а м а</p> <p>[0, 1, 2, 3, 4]</p> <p>д [1, 1, 2, 3, 4]</p> <p>а [2, 2, 1, 2, 3]</p> <p>л [3, 2, 2, 2, 3]</p> <p>а [4, 3, 2, 3, 2]</p> <p>й [5, 4, 3, 3, 3]</p> <p>л [6, 5, 4, 4, 4]</p> <p>а [7, 6, 5, 5, 4]</p> <p>м [8, 7, 6, 5, 5]</p> <p>а [9, 8, 7, 6, 5]</p>
--	---	--

<p>1 д в а 3 4</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>1 [1, 0, 1, 2, 3, 4, 5]</p> <p>2 [2, 1, 1, 2, 3, 4, 5]</p> <p>3 [3, 2, 2, 2, 3, 3, 4]</p> <p>д [4, 3, 2, 3, 3, 4, 4]</p> <p>в [5, 4, 3, 2, 3, 4, 5]</p> <p>а [6, 5, 4, 3, 2, 3, 4]</p> <p>4 [7, 6, 5, 4, 3, 3, 3]</p> <p>5 [8, 7, 6, 5, 4, 4, 4]</p>	
---	--

## Дамерау-Левенштейн

<p>п р о л е т</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>э [1, 1, 2, 3, 4, 5, 6]</p> <p>а [2, 2, 2, 3, 4, 5, 6]</p> <p>л [3, 3, 3, 3, 3, 4, 5]</p> <p>е [4, 4, 4, 4, 4, 3, 4]</p> <p>т [5, 5, 5, 5, 5, 4, 3]</p>	<p>к о р о в а</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>к [1, 0, 1, 2, 3, 4, 5]</p> <p>о [2, 1, 0, 1, 2, 3, 4]</p> <p>о [3, 2, 1, 1, 1, 2, 3]</p> <p>р [4, 3, 2, 1, 1, 2, 3]</p> <p>в [5, 4, 3, 2, 2, 1, 2]</p> <p>а [6, 5, 4, 3, 3, 2, 1]</p>	<p>л а м а</p> <p>[0, 1, 2, 3, 4]</p> <p>д [1, 1, 2, 3, 4]</p> <p>а [2, 2, 1, 2, 3]</p> <p>л [3, 2, 2, 2, 3]</p> <p>а [4, 3, 2, 3, 2]</p> <p>й [5, 4, 3, 3, 3]</p> <p>л [6, 5, 4, 4, 4]</p> <p>а [7, 6, 5, 5, 4]</p> <p>м [8, 7, 6, 5, 5]</p> <p>а [9, 8, 7, 6, 5]</p>
--	---	--

<p>1 д в а 3 4</p> <p>[0, 1, 2, 3, 4, 5, 6]</p> <p>1 [1, 0, 1, 2, 3, 4, 5]</p> <p>2 [2, 1, 1, 2, 3, 4, 5]</p> <p>3 [3, 2, 2, 2, 3, 3, 4]</p> <p>д [4, 3, 2, 3, 3, 4, 4]</p> <p>в [5, 4, 3, 2, 3, 4, 5]</p> <p>а [6, 5, 4, 3, 2, 3, 4]</p> <p>4 [7, 6, 5, 4, 3, 3, 3]</p> <p>5 [8, 7, 6, 5, 4, 4, 4]</p>	<p>1 3 2 4</p> <p>[0, 1, 2, 3, 4]</p> <p>1 [1, 0, 1, 2, 3]</p> <p>2 [2, 1, 1, 1, 2]</p> <p>3 [3, 2, 1, 1, 2]</p> <p>4 [4, 3, 2, 2, 1]</p> <p>5 [5, 4, 3, 3, 2]</p>
---	--

## 4.2 Постановка эксперимента

Исследование скорости работы алгоритма. Для временного сравнительного анализа рекурсивной реализации алгоритма Левенштейна, нерекурсивной реализации алгоритма Левенштейна и алгоритма Дамерау-Левенштейна использовалось несколько запусков. Проводились вычисления по тикам для слов одинаковой длины (5, 10, 12). Можно сделать вывод, что рекурсивную версию алгоритма Левенштейна не стоит применять в крупных базах данных, так как он существенно замедляет работу.

Левинштейн:  
4.346780077378723e-05  
Дамерал-Левинштейн:  
5.5844049605212505e-05  
Левинштейн Рекурсивный:  
0.0012421527790564893

Левинштейн:  
0.00012406434804185105  
Дамерал-Левинштейн:  
0.00016813586827082918  
Левинштейн Рекурсивный:  
5.965388461774144

///

Левинштейн:  
0.00016481541126727654  
Дамерал-Левинштейн:  
0.0002327338499763191  
Левинштейн Рекурсивный:  
185.55000261108663

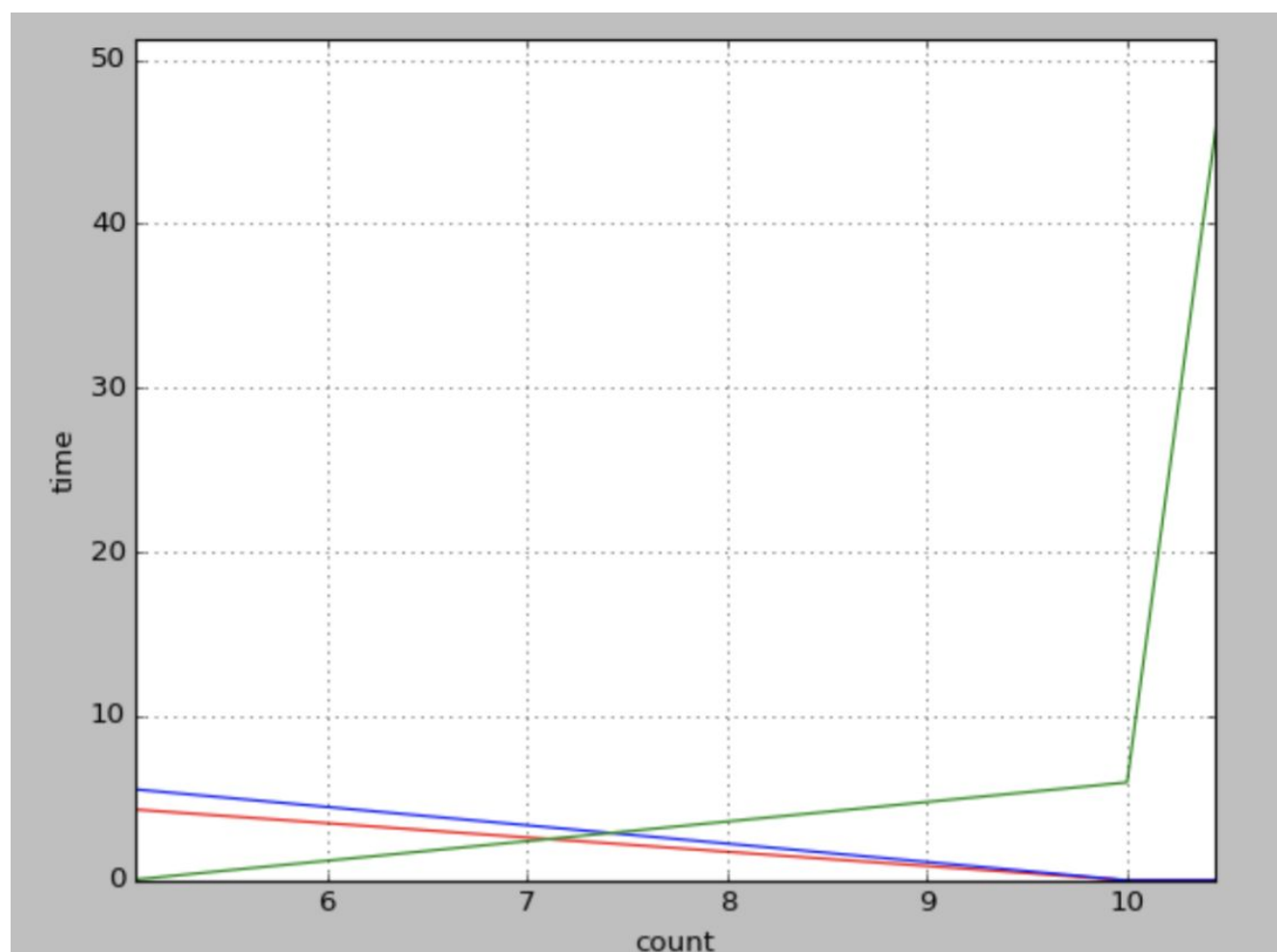


Рисунок 4.2.1 – Пример работы алгоритмов по скорости со словами одной длины

Эксперимент с разными длинами слов, где первое слово (10), второе (9, 8, 7)

Левинштейн:  
0.0001465059864402284  
Дамерал-Левинштейн:  
0.00021384733459345676  
Левинштейн Рекурсивный:  
3.3649366451568463

Левинштейн:  
0.00013416864021368284  
Дамерал-Левинштейн:  
0.0001927710347897748  
Левинштейн Рекурсивный:  
1.4498592257389942

Левинштейн:  
0.00012234535007991002  
Дамерал-Левинштейн:  
0.0002919838606949103  
Левинштейн Рекурсивный:  
0.5419934992466507

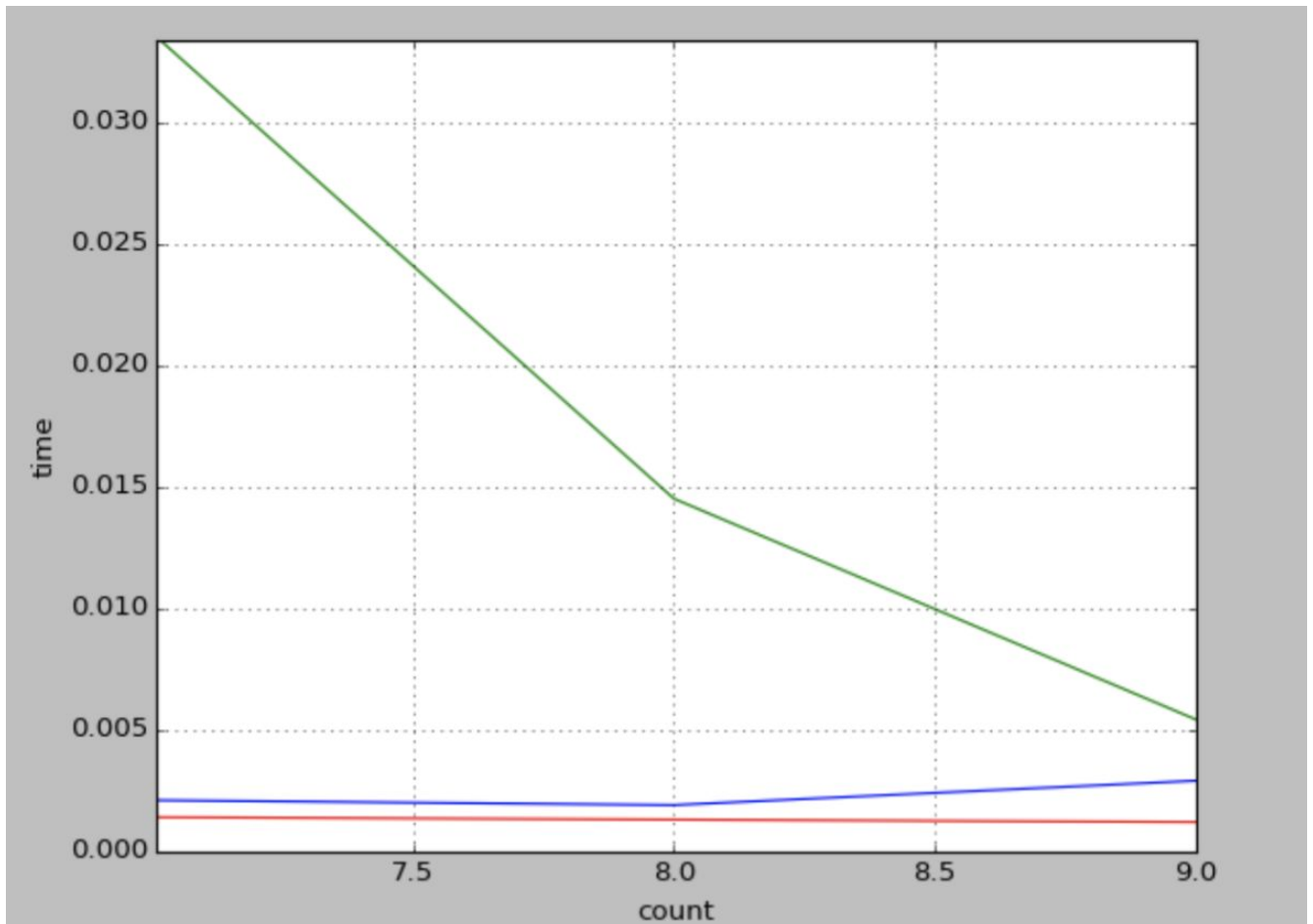


Рисунок 4.2.2 – Пример работы алгоритмов по скорости с словами разной длины



## **Заключение**

1. Изучены теоретические понятия в алгоритмах для нахождения расстояния Левенштейна и Дамерау-Левенштейна;
2. Проведен аналитический вывод формул для заполнения матриц расстояний;
3. Проведено сравнение трех реализаций заданного алгоритма
4. В рамках данной работы было сделано заключение, что рекурсивный алгоритм сильно проигрывает по скорости двум другим реализациям. Скоростные отличия между алгоритмом Дамерау-Левенштейна и Левенштейна в рамках данной работы найдены не были;
5. Определение расстояния по Левенштейну имеет недостатки:
  - 1) при перестановке местами слов или частей слов получаются большие расстояния;
  - 2) расстояния между совершенно разными короткими словами будут меньше в отличии от расстояния между двумя длинными и очень похожими, но длинными словами.