

Оглавление

Введение	2
1 Аналитическая часть	3
2 Конструкторская часть	5
2.1 Расстояние Левенштейна, матричный метод	5
2.2 Расстояние Дамерау-Левенштейна, матричный метод .	5
2.3 Расстояние Левенштейна, рекурсивный метод	6
2.4 Расстояние Левенштейна, рекурсивный метод с запол- нением матрицы	6
2.5 Требования к программному обеспечению	7
2.6 Заготовки тестов	7
3 Технологическая часть	12
3.1 Выбор языка программирования	12
3.2 Листинг кода	12
3.3 Результаты тестирования	15
3.4 Оценка памяти	17
3.5 Оценка памяти	18
Исследовательская часть	20
Заключение	20
Результат экспериментов	20
Сравнительный анализ	21
Заключение	22

Введение

Расстояние Левенштейна – минимальное количество редакционных операций, которые необходимы для превращения одной строки в другую. Существуют следующие редакционные операции:

- вставка символа;
- удаление символа;
- замена символа;

Расстояние Дамерау-Левенштейна также учитывает и операцию транспозиции – перестановки двух соседних символов местами.

Данные расстояния имеют большое количество применений. Они используются для автокоррекции при выполнении поисковых запросов и печати на клавиатуре, а также в биоинформатике для сравнения генов, представленных в строковом формате.

1. Аналитическая часть

Целью лабораторной работы является реализация и сравнение алгоритмов поиска расстояний Левенштейна и Дamerau-Левенштейна.

Выделены следующие задачи лабораторной работы:

- математическое описание расстояний Левенштейна и Дamerau-Левенштейна;
- описание и реализация алгоритмов поиска расстояний;
- проведение замеров процессорного времени работы алгоритмов при различных размерах строк;
- оценка наибольшей используемой каждым алгоритмом памяти;
- проведение сравнительного анализа алгоритмов на основании экспериментов;

Задача по поиску расстояний заключается в нахождении такой последовательности операций, применение которых даст минимальный суммарный штраф. Штрафы операций:

- вставка (I) – 1;
- замена (R) – 1;
- удаление (D) – 1;
- совпадение (M) – 0;
- транспозиция (T) – 1;

Для решения данной проблемы используется рекуррентная формула вычисления расстояний. Пусть $D(s1[1..i], s2[1..j])$ – расстояние Левенштейна для подстроки $s1$ длиной i и $s2$ длиной j . Формула

вычисления D:

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

Аналогично рекуррентно представляется формула расстояния Дамерау-Левенштейна:

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

2. Конструкторская часть

Рассмотрим алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна для строк $s1$ и $s2$ с длинами $n1$ и $n2$ соответственно.

2.1. Расстояние Левенштейна, матричный метод

Алгоритм матричного поиска расстояния Левенштейна основывается на вышеописанной рекуррентной формуле. Создаётся целочисленная матрица размерами $(n1+1) \times (n2+1)$. В каждой клетке $[i][j]$ этой матрицы будет записано значение $D(s1[1..i-1], s2[1..j-1])$. В случае, когда $i=1$ или $j=1$ вместо строк $s1$ и $s2$ соответственно будут выступать пустые строки. Искомым расстоянием Левенштейна будет значение ячейки $[n1+1][n2+1]$.

Нахождение расстояний алгоритм начинает с заполнения первого столбца и первой строки, так как они являются базой для рекуррентной формулы. После этого производится построчное заполнение остальной части матрицы.

Схема алгоритма приведена на рисунке 2.1.

2.2. Расстояние Дамерау-Левенштейна, матричный метод

Алгоритм является модификацией вышеописанного способа нахождения расстояния Левенштейна. Дополнительно для ячейки $[i][j]$ ($i > 2, j > 2$) рассматривается вариант перехода из клетки $[i-2][j-2]$, при условии, что $s1[i] = s2[j-1]$ и $s1[i-1] = s2[j]$. Искомым расстоянием Дамерау-Левенштейна также является значение ячейки $[n1+1][n2+1]$.

Схема алгоритма приведена на рисунке 2.2.

2.3. Расстояние Левенштейна, рекурсивный метод

Данный алгоритм использует только рекурсивную формулу нахождения $D(s1[1..i], s2[1..j])$. Для этого используется рекурсивная функция, принимающая в себя строки $s1$, $s2$ и длины подстрок i , j . Функция вызывает функции для тех же строк, и длин: $(i-1, j-1)$, $(i-1, j)$ и $(i, j-1)$, после чего возвращает минимальный из них.

Схема алгоритма приведена на рисунке 2.3.

2.4. Расстояние Левенштейна, рекурсивный метод с заполнением матрицы

В данном случае, в качестве основы используется алгоритм Дейкстры поиска расстояний до вершин в графе. Создаётся матрица размерами $(n1+1) \times (n2+1)$, все ячейки которой изначально заполнены значением $+\infty$. В каждой клетке $[i][j]$ этой матрицы будет записано значение $D(s1[1..i-1], s2[1..j-1])$.

Рекурсивная функция получает матрицу, индексы i , j положения в ней и две строки. Алгоритм начинает свою работу с ячейки $[1][1]$, которая заполняется значением 0. Из положения $[i][j]$ рассматривается переход в соседние ячейки $[i+1][j+1]$, $[i+1][j]$, $[i][j+1]$. В случае, если соседняя ячейка расположена в пределах матрицы, и расстояние R при переходе из данной ячейки меньше ныне хранимого в ней, то значение соседней ячейки меняется на R , после чего функция запускается уже для соседней ячейки. После завершения работы всех функций, расстояние Левенштейна расположено в ячейке $[n1+1][n2+1]$.

Схема алгоритма приведена на рисунке 2.4.

2.5. Требования к программному обеспечению

Для полноценной проверки и оценки алгоритмов необходимо выполнить следующее.

1. Обеспечить возможность консольного ввода двух строк и выбора алгоритма для поиска расстояния. Программа должна вывести вычисленное редакционное расстояние, а также вывести матрицу поиска, в случае использования её в выбранном алгоритме.
2. Реализовать функцию замера процессорного времени, затраченного функцией. Для этого также создать возможность ввода длины строк, на которых будет выполнен замер.

2.6. Заготовки тестов

При проверке алгоритмов необходимо будет использовать следующие классы тестов:

- две пустые строки;
- одна из строк пустая;
- одинаковые строки;
- применение транспозиции даёт минимальное расстояние (для Дамерау-Левенштайна);

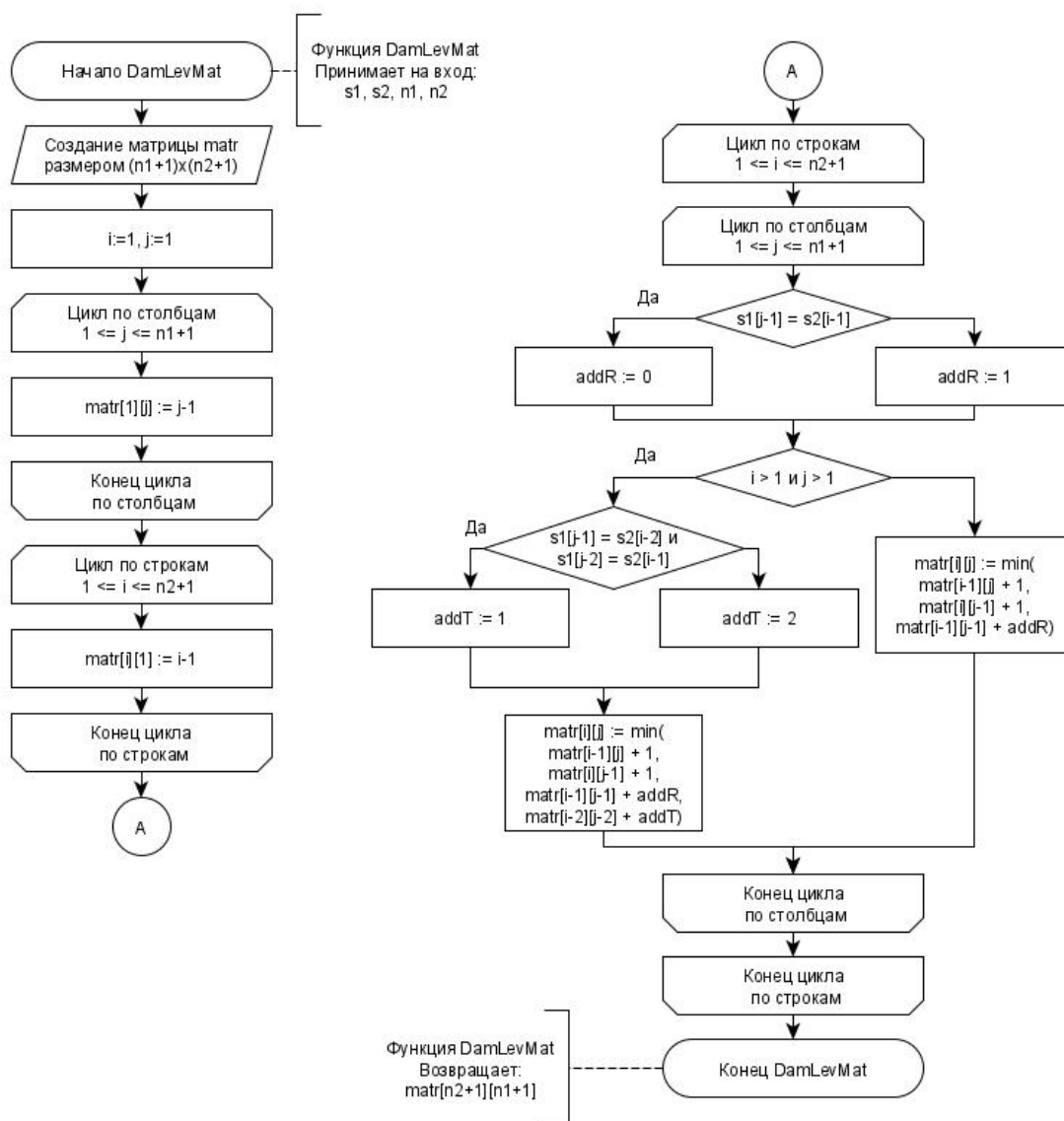


Рис. 2.2: Алгоритм нахождения расстояния Дameraу-Левенштейна, матричный метод

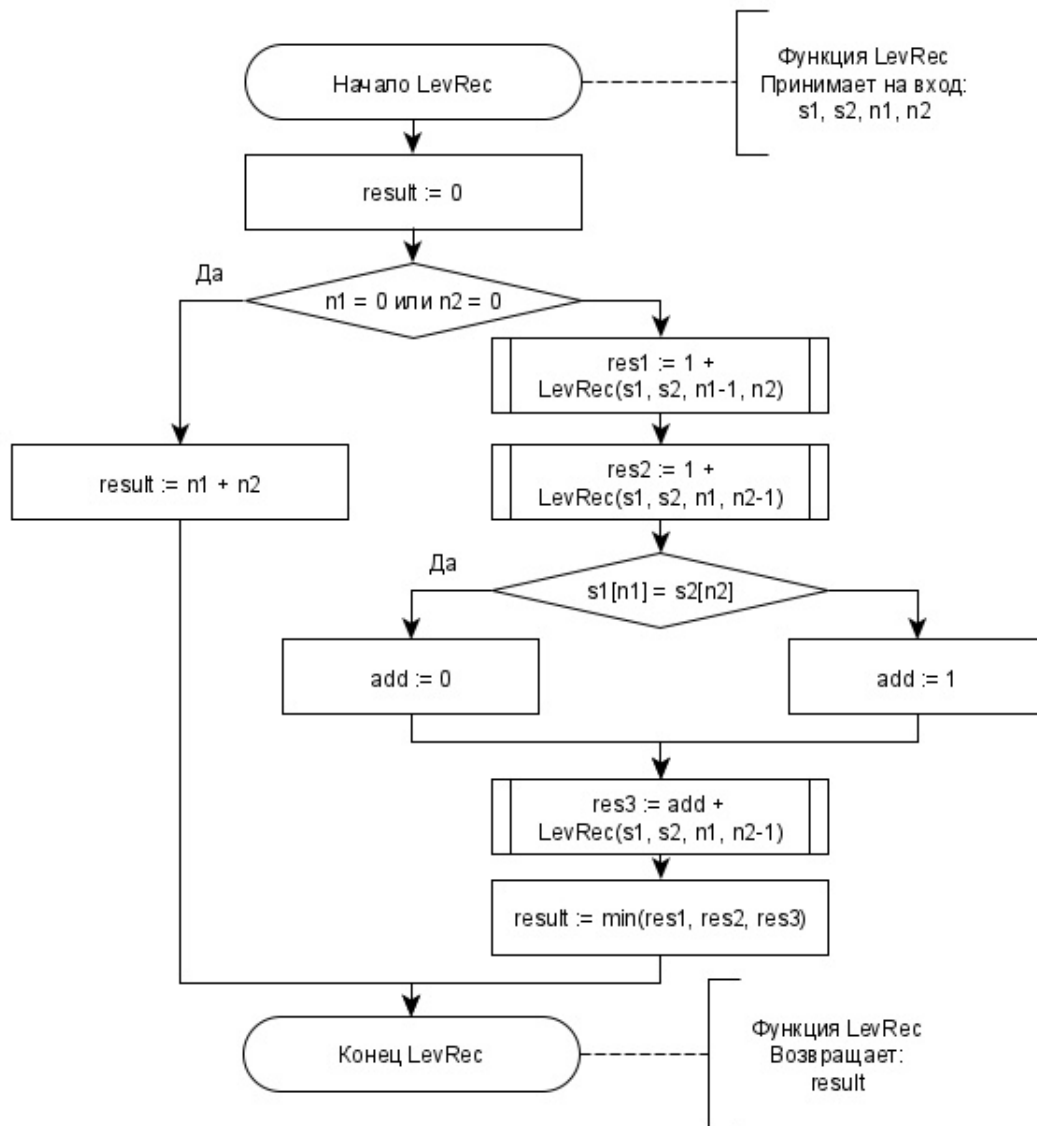


Рис. 2.3: Алгоритм нахождения расстояния Левенштейна, рекурсивный метод

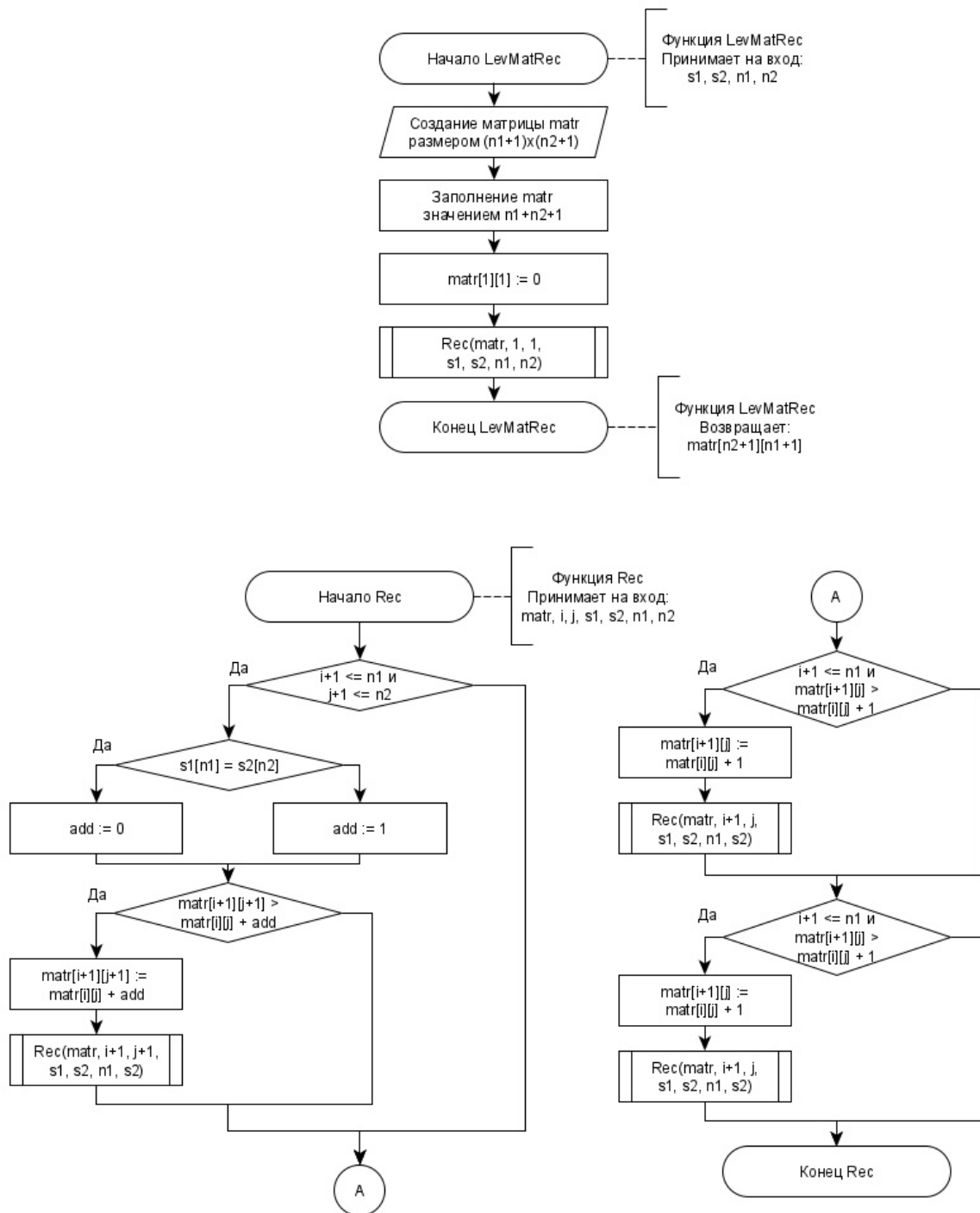


Рис. 2.4: Алгоритм нахождения расстояния Левенштейна, рекурсивный метод с заполнением матрицы

3. Технологическая часть

3.1. Выбор языка программирования

В качестве языка программирования был выбран Python 3, так как имеется опыт работы с ним, и с библиотеками, позволяющими провести исследование и тестирование программы.

3.2. Листинг кода

Реализация алгоритмов поиска расстояний представлена на листингах 3.1-3.4.

Листинг 3.1: Функция нахождения расстояния Левенштейна матричным методом.

```
1 def lev_matrix(s1, s2, is_print=False):
2     matr = [[0] * (len(s1)+1) for i in range(len(s2)+1)]
3
4     for j in range(len(s1)+1):
5         matr[0][j] = j
6     for i in range(len(s2)+1):
7         matr[i][0] = i
8
9     for i in range(1, len(s2)+1):
10        for j in range(1, len(s1)+1):
11            add = 0 if s1[j-1] == s2[i-1] else 1
12            matr[i][j] = min(matr[i-1][j]+1, matr[i][j-1]+1, matr[i-1][j-1]+add)
13
14    if is_print:
15        print("Расстояние: ", matr[i][j])
16        print_matrix(matr)
17    return matr[i][j]
18
```

Листинг 3.2: Функции нахождения расстояния Левенштейна рекур-

СИВНЫМ МЕТОДОМ.

```
1 def _lev_rec(s1, s2, len1, len2):
2     if len1 == 0: return len2
3     elif len2 == 0: return len1
4     else:
5         return min(_lev_rec(s1, s2, len1, len2-1)+1,
6                     _lev_rec(s1, s2, len1-1, len2)+1,
7                     _lev_rec(s1, s2, len1-1, len2-1) +
8                     (0 if s1[len1-1] == s2[len2-1]
9                      else 1))
10 def lev_recursion(s1, s2, is_print=False):
11     res = _lev_rec(s1, s2, len(s1), len(s2))
12     if is_print:
13         print("Расстояние: ", res)
14     return res
15
```

Листинг 3.3: Функции нахождения расстояния Левенштейна рекурсивным методом с заполнением матрицы.

```
1 def _lev_mr(matr, i, j, s1, s2):
2     if i+1 < len(matr) and j+1 < len(matr[0]):
3         add = 0 if s1[j] == s2[i] else 1
4         if matr[i+1][j+1] > matr[i][j] + add:
5             matr[i+1][j+1] = matr[i][j] + add
6             _lev_mr(matr, i+1, j+1, s1, s2)
7     if j+1 < len(matr[0]) and (matr[i][j+1] > matr[i][j] + 1):
8         matr[i][j+1] = matr[i][j] + 1
9         _lev_mr(matr, i, j+1, s1, s2)
10    if i+1 < len(matr) and (matr[i+1][j] > matr[i][j] + 1):
11        matr[i+1][j] = matr[i][j] + 1
12        _lev_mr(matr, i+1, j, s1, s2)
13
14 def lev_matrix_recursion(s1, s2, is_print=False):
```

```

15 max_len = max(len(s1), len(s2)) + 1
16 matr = [[max_len] * (len(s1)+1) for i in range(len(s2)+1)]
17 matr[0][0] = 0
18 _lev_mr(matr, 0, 0, s1, s2)
19
20 if is_print:
21     print("Расстояние: ", matr[-1][-1])
22     print_matrix(matr)
23 return matr[-1][-1]
24

```

Листинг 3.4: Функция нахождения расстояния Дамерау-Левенштейна матричным методом.

```

1 def dem_lev_matrix(s1, s2, is_print=False):
2     if len(s1) == 0: return len(s2)
3     elif len(s2) == 0: return len(s1)
4     matr = [[0] * (len(s1) + 1) for i in range(len(s2) + 1)]
5     for j in range(len(s1)+1):
6         matr[0][j] = j
7     for i in range(len(s2)+1):
8         matr[i][0] = i
9
10    for i in range(1, len(s2) + 1):
11        addM = 0 if s1[0] == s2[i-1] else 1
12        matr[i][1] = min(matr[i-1][1] + 1, matr[i][0] + 1,
13            matr[i-1][0] + addM)
14    for j in range(2, len(s1) + 1):
15        addM = 0 if s1[j-1] == s2[0] else 1
16        matr[1][j] = min(matr[0][j] + 1, matr[1][j-1] + 1,
17            matr[0][j-1] + addM)
18
19    for i in range(2, len(s2)+1):
20        for j in range(2, len(s1)+1):

```

```

21         addM = 0 if s1[j-1] == s2[i-1] else 1
22         addT = 1 if (s1[j-2] == s2[i-1] and s1[j-1] == s2[i-2])
           else 2
23         matr[i][j] = min(matr[i-1][j]+1, matr[i][j-1]+1,
24         matr[i-1][j-1]+addM, matr[i-2][j-2]+addT)
25
26     if is_print:
27         print("Расстояние: ", matr[i][j])
28         print_matrix(matr)
29     return matr[i][j]
30

```

3.3. Результаты тестирования

Для тестирования написанных функций была использована библиотека unittest. Тестирование функций проводилось за счёт сравнения результата, возвращённого функцией и ожидаемого расстояния для разных наборов строк.

Состав тестов приведён в листинге 3.5.

Листинг 3.5: Модульные тесты

```

1 import unittest
2 import main
3
4 # Общий набор тестов для всех алгоритмов
5 class GeneralTest(unittest.TestCase):
6     # Данный класс является абстрактным, поэтому для него тесты
       пропускаются
7     @unittest.skip("Skip GeneralTest")
8     def setUp(self):
9         self.function = None
10
11 # Проверка пустыми строками

```

```

12 def test_empty(self):
13     self.assertEqual(self.function("", ""), 0)
14     self.assertEqual(self.function("a", ""), 1)
15     self.assertEqual(self.function("", "b"), 1)
16
17 # Проверка нахождения совпадений
18 def test_match(self):
19     self.assertEqual(self.function("abc", "abc"), 0)
20     self.assertEqual(self.function("a", "a"), 0)
21     self.assertEqual(self.function("A", "a"), 1)
22
23 # Прочие общие тесты
24 def test_other(self):
25     self.assertEqual(self.function("q", "w"), 1)
26     self.assertEqual(self.function("aq", "aw"), 1)
27     self.assertEqual(self.function("a", "aw"), 1)
28     self.assertEqual(self.function("aw", "a"), 1)
29
30
31 # Набор тестов для алгоритмов поиска расстояния Левенштейна
32 class LevTest(GeneralTest):
33     def test_lev(self):
34         self.assertEqual(self.function("stolb", "telo"), 3)
35         self.assertEqual(self.function("kult_tela", "tela_kult"),
36                             6)
37         self.assertEqual(self.function("развлечение", "увлечение"),
38                             3)
39
40 # Набор тестов для алгоритма поиска расстояния ДамерауЛевенштейна—
41 class DemLevMatrixTest(GeneralTest):
42     def setUp(self):
43         self.function = main.dem_lev_matrix

```



```

44 def dem_lev_test(self):
45     self.assertEqual(self.function("aba", "aab"), 1)
46     self.assertEqual(self.function("ab", "ba"), 1)
47     self.assertEqual(self.function("abb", "bab"), 1)
48
49
50 # Алгоритмы поиска расстояния Левенштейна проходят одинаковые тесты из
    класса LevTest
51 # Алгоритм поиска расстояния Левенштейна, матричный метод
52 class LevMatrixTest(LevTest):
53     def setUp(self):
54         self.function = main.lev_matrix
55 # Алгоритм поиска расстояния Левенштейна, рекурсивный метод
56 class LevRecursionTest(LevTest):
57     def setUp(self):
58         self.function = main.lev_recursion
59 # Алгоритм поиска расстояния Левенштейна, рекурсивный метод с
    заполнением матрицы
60 class LevMatRecTest(LevTest):
61     def setUp(self):
62         self.function = main.lev_matrix_recursion
63
64 # Точка входа, запуск тестов
65 if __name__ == "__main__":
66     unittest.main()

```

3.4. Оценка памяти

Произведём оценку наибольшей затрачиваемой алгоритмом памяти M_{max} при поиске расстояний для строк $s1$ и $s2$. Для удобства оценки примем длину обеих строк за n .

Расстояние Левенштейна, матричный метод. Память затрачивается на матрицу и две строки.

$$M_{max} = (n+1) * (n+1) * sizeof(int) + (n+n) * sizeof(char) = (n+1) * (n+1) * 16 + (n+n) * 17 = 16 * n^2 + 2 * 17n + 16 \text{ байт}$$

Расстояние Дамерау-Левенштейна, матричный метод. Аналогично.

$$M_{max} = 16 * n^2 + 2 * 17n + 16 \text{ байт}$$

Расстояние Левенштейна, рекурсивный метод. Память используется при каждом вызове функции. Одна функция принимает в качестве аргумента 2 строки по значению, 2 размера строк. Максимальная глубина рекурсии = $n+n$.

$$M_{max} = (n+n) * (2n * sizeof(char) + 2 * sizeof(int)) = 2n * (2n + 32) = 4n^2 + 64n \text{ байт}$$

Расстояние Левенштейна, рекурсивный метод с матрицей. Память используется для матрицы и при каждом вызове функции. Максимальная глубина рекурсии = $n+n$.

$$M_{max} = (n+1) * (n+1) * sizeof(int) + (n+n) * (2n * sizeof(char) + 2 * sizeof(int)) = (n^2 + 2n + 1) * 16 + 2n * (2n + 32) = 20n^2 + 96n + 16 \text{ байт}$$

3.5. Оценка памяти

Для замера процессорного времени исполнения функции используется библиотека `time`. Проведение измерений производится в функции, приведённой в листинге 3.6. Также в листинге приведена функция `random_str` для создания строки заданной длины из случайной последовательности символов, с использованием библиотеки `random`.

Листинг 3.6: Функция замера процессорного времени работы функции

```
1 def random_str(length):
```

```

2  a = []
3  for i in range(length):
4      a.append(random.choice("qwerty"))
5  return "".join(a)
6
7  def test_time(func):
8      length = int(input("Введите длину строки: "))
9      s1 = random_str(length)
10     s2 = random_str(length)
11     print("Строка 1:", s1)
12     print("Строка 2:", s2)
13
14     begin_t = time.process_time()
15     count = 0
16     while time.process_time() - begin_t < 1.0:
17         func(s1, s2)
18         count += 1
19
20     t = time.process_time() - begin_t
21     print("Выполнено {:} операций за {:} секунд".format(count, t))
22     print("Время: {:.4} секунд".format(t / count))
23

```

Исследовательская часть

Заключение

Измерения процессорного времени проводятся при равных длинах строк $s1$ и $s2$. Содержание строк сгенерировано случайным образом. Изучается время работы при длинах: 1, 3, 10, 20, 100, 1000. Для повышения точности, каждый замер производится пять раз, за результат берётся среднее арифметическое.

Результат экспериментов

По результатам измерений процессорного времени можно составить таблицу 4.1.

Таблица 4.1: Результат измерений процессорного времени (в секундах)

	1	3	10	20	100	1000
Лев., матрица	$7 * 10^{-6}$	$1.9 * 10^{-5}$	$1.3 * 10^{-4}$	$4.7 * 10^{-4}$	0.013	1.405
Лев., рекурсия	$3 * 10^{-6}$	$4.7 * 10^{-5}$	6.984	—	—	—
Лев., рекурсия с матрицей	$1 * 10^{-5}$	$4.1 * 10^{-5}$	$4.1 * 10^{-4}$	$2.5 * 10^{-3}$	0.38	—
Д-Л, матрица	$8 * 10^{-6}$	$2.8 * 10^{-5}$	$1.7 * 10^{-4}$	$6.1 * 10^{-4}$	0.016	2.031

В алгоритме нахождения расстояния Левенштейна с помощью рекурсии замеры на длине строк более 10 не проводились, так как время выполнения было слишком велико (более 10 минут). В алгоритме рекурсии с заполнением матрицы не удалось провести измерения при длине 1000, так как была превышена максимальная глубина рекурсии.

Сравнительный анализ

По результатам эксперимента можно заключить следующее.

- Наиболее быстрое действие алгоритмом поиска расстояния Левенштейна является алгоритм, использующий матрицу.
- Рекурсивный алгоритм с использованием матрицы показывает значительно более низкую скорость роста времени по сравнению с рекурсивным алгоритмом.
- Алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна с помощью матрицы показывают схожую скорость роста времени, однако первый алгоритм несколько быстрее.

Заключение

В ходе лабораторной работы достигнута поставленная цель: реализация и сравнение алгоритмов поиска расстояний Левенштейна и Дamerau-Левенштейна. Решены все задачи работы.

Были изучены и описаны понятия расстояний Левенштейна и Дamerau-Левенштейна. Также были описаны и реализованы алгоритмы поиска расстояний. Проведены замеры процессорного времени работы каждого алгоритма при различных строках, оценена наибольшая занимаемая память. На основании оценок и экспериментов проведён сравнительный анализ.