



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»

Лабораторная работа № 1

Дисциплина Анализ алгоритмов

Тема Расстояния Левенштейна и Дamerau-Левенштейна

Студент Куликов Д. А.

Группа ИУ7-52Б

Оценка (баллы) _____

Преподаватель Волкова Л. Л.

Москва — 2020 г.

Оглавление

Введение	3
1. Аналитическая часть	4
1.1 Описание расстояний	4
2. Конструкторская часть	6
2.1 Схемы алгоритмов	6
3. Технологическая часть	11
3.1 Требования к программному обеспечению	11
3.2 Средства реализации	11
3.3 Листинг кода	12
3.4 Сравнительный анализ матричной и рекурсивной реализаций	14
3.5 Описание тестирования	15
4. Экспериментальная часть	16
4.1 Примеры работы	16
4.2 Результаты тестирования	18
4.3 Постановка эксперимента по замеру времени	20
4.4 Сравнительный анализ на материале экспериментальных данных	20
Заключение	22
Список литературы	23

Введение

Расстояние Левенштейна — метрика, измеряющая разность между двумя последовательностями символов. Она определяется как минимальное количество односимвольных операций (а именно вставки, удаления, замены), необходимых для превращения одной последовательности символов в другую. В общем случае, операциям, используемым в этом преобразовании, можно назначить разные цены [1].

Расстояние Левенштейна широко используется в теории информации и компьютерной лингвистике. Примеры использования:

- исправление ошибок в слове;
- в биоинформатике для сравнения генов, хромосом и белков;
- сравнение текстовых файлов утилитой diff и ей подобными.

Цель работы: изучение и применение на практике алгоритмов нахождения расстояния Левенштейна и Дamerau-Левенштейна, а также получение практических навыков реализации этих алгоритмов и сравнении их между собой.

Задачи работы:

- 1) изучение алгоритмов Левенштейна и Дamerau-Левенштейна нахождения расстояния между строками;
- 2) получение практических навыков реализации указанных алгоритмов: двух алгоритмов в матричной версии и одного из алгоритмов в рекурсивной версии;
- 3) сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
- 4) экспериментальное подтверждение различий во временной эффективности рекурсивной и нерекурсивной реализаций выбранного алгоритма определения расстояния между строками при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализации на варьирующихся длинах строк;
- 5) описание и обоснование полученных результатов в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

1 | Аналитическая часть

Задача по нахождению расстояния Левенштейна заключается в поиске минимального количества операций вставки/удаления/замены для превращения одной строки в другую.

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

Операции обозначаются следующим образом:

- 1) D (англ. delete) — удалить;
- 2) I (англ. insert) — вставить;
- 3) R (англ. replace) — заменить;
- 4) M (англ. match) — совпадение;
- 5) T (англ. transpose) — перестановка.

Операция M имеет штраф 0, другие операции — штраф 1.

1.1. Описание расстояний

Пусть S_1 и S_2 — две строки (длиной M и N соответственно) над некоторым алфавитом, тогда расстояние Левенштейна можно подсчитать по следующей рекуррентной формуле:

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

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

Расстояние Дамерау-Левенштейна вычисляется по следующей рекуррентной формуле:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min(& \text{если } i, j > 1 \text{ and } a_i = b_{j-1}, a_{i-1} = b_j \\ D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[j]) \\ D(i - 2, j - 2) + 1, \\) \\ \min(& \text{иначе} \\ D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[j]) \\) \end{cases}$$

2 | Конструкторская часть

В этом разделе содержатся схемы алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна и сравнительный анализ матричной и рекурсивной реализаций.

2.1. Схемы алгоритмов

На рис. 2.1-2.4 приведены схемы следующих алгоритмов:

- 1) матричный алгоритм нахождения расстояния Левенштейна;
- 2) рекурсивный алгоритм нахождения расстояния Левенштейна;
- 3) рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы;
- 4) матричный алгоритм нахождения расстояния Дамерау-Левенштейна.

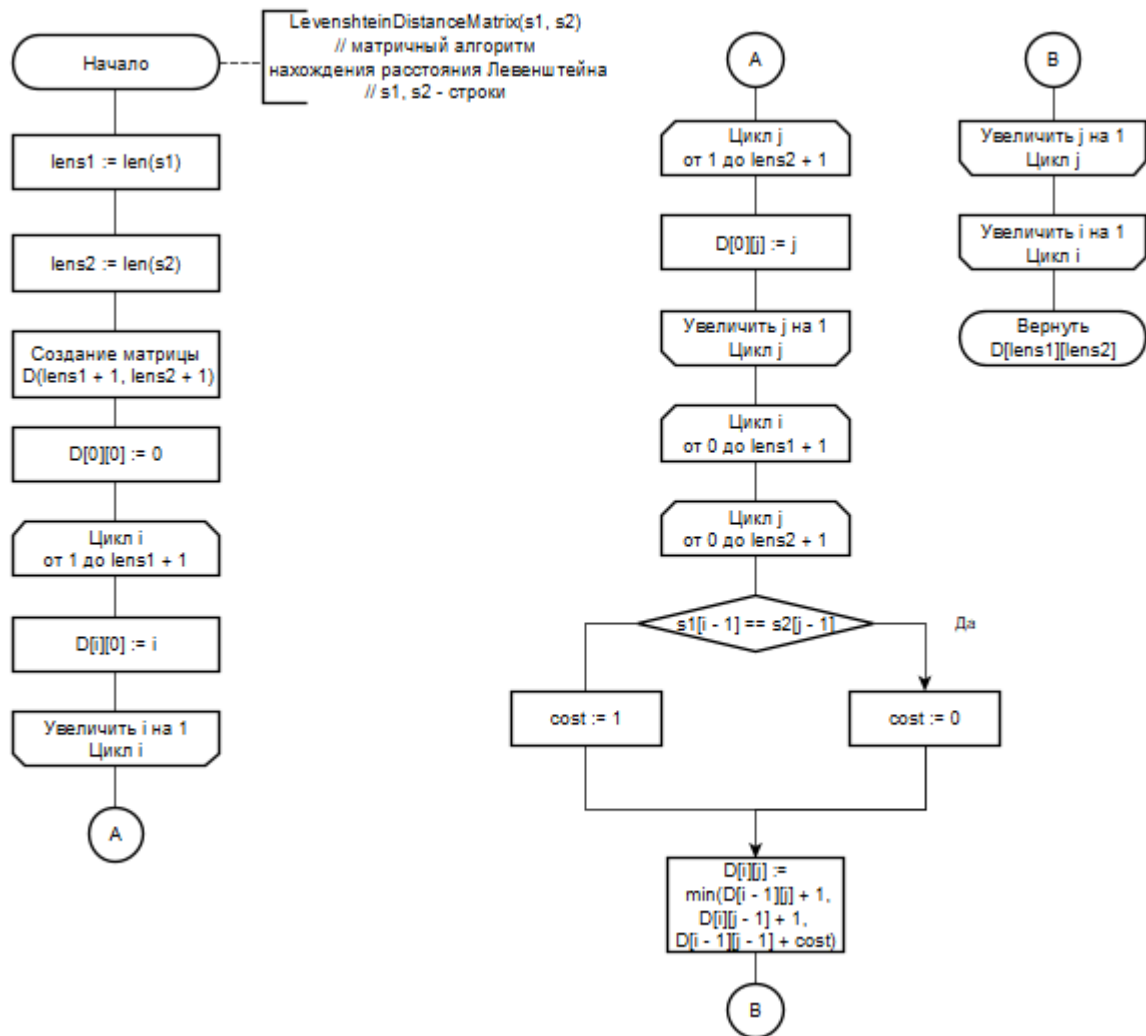


Рис. 2.1. Матричный алгоритм нахождения расстояния Левенштейна

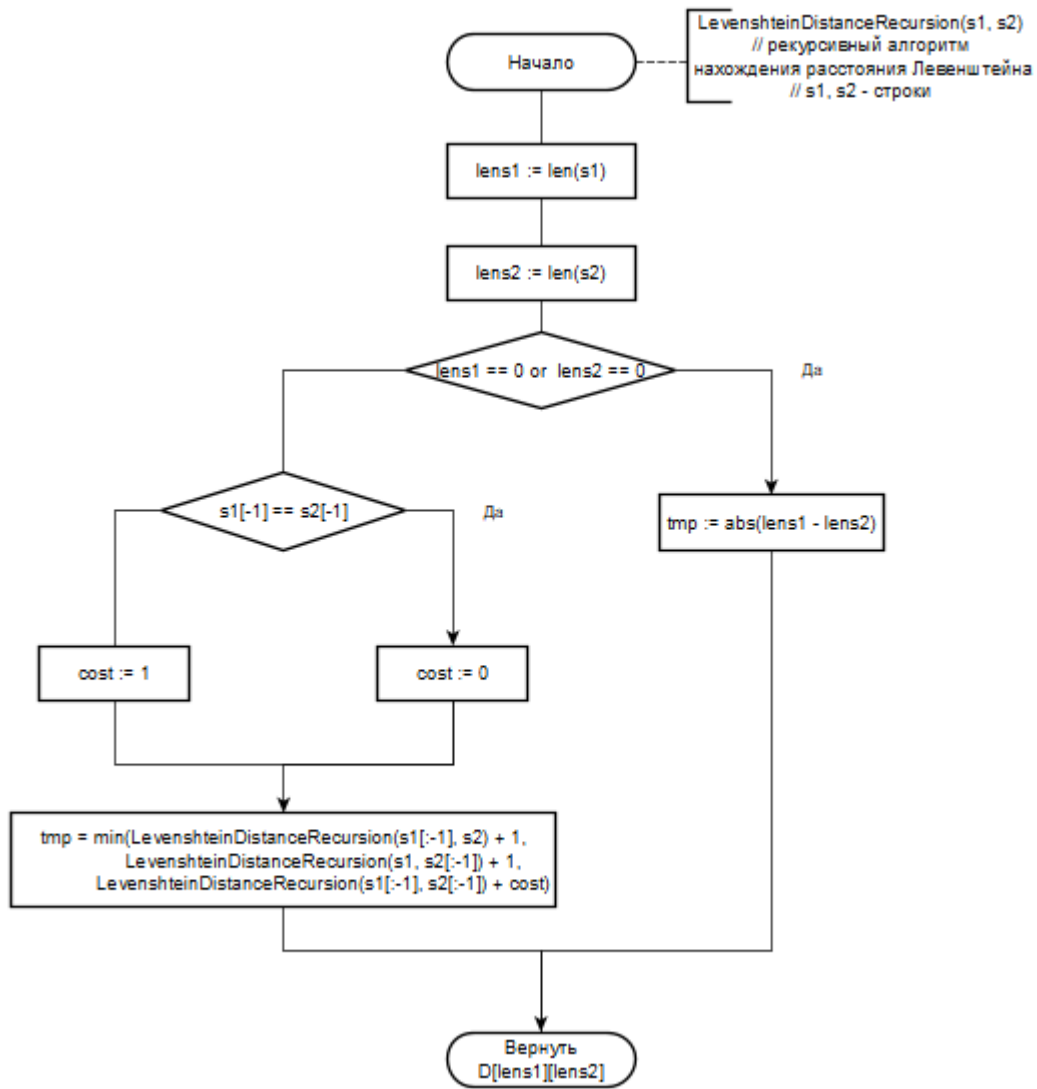


Рис. 2.2. Рекурсивный алгоритм нахождения расстояния Левенштейна

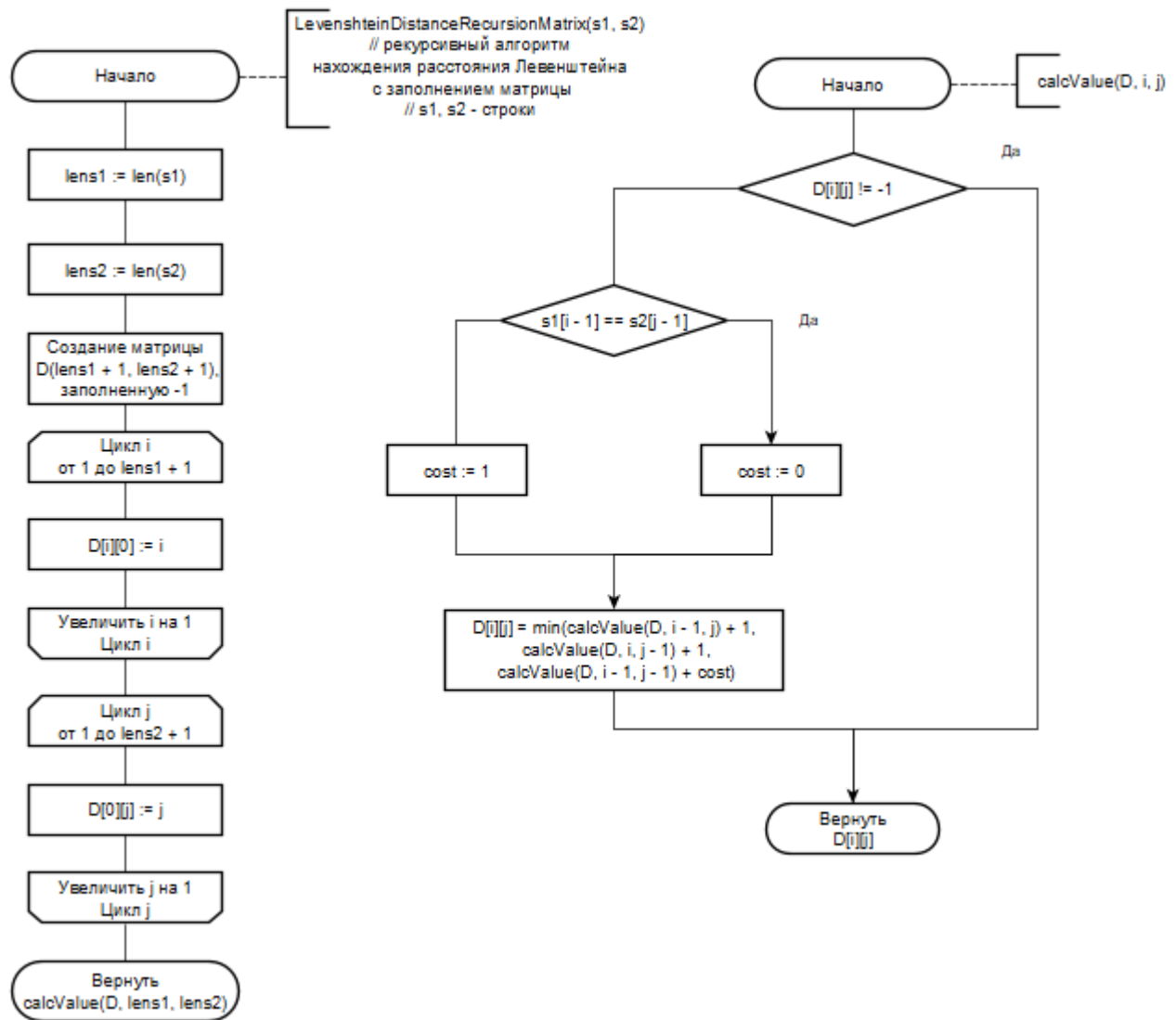


Рис. 2.3. Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы

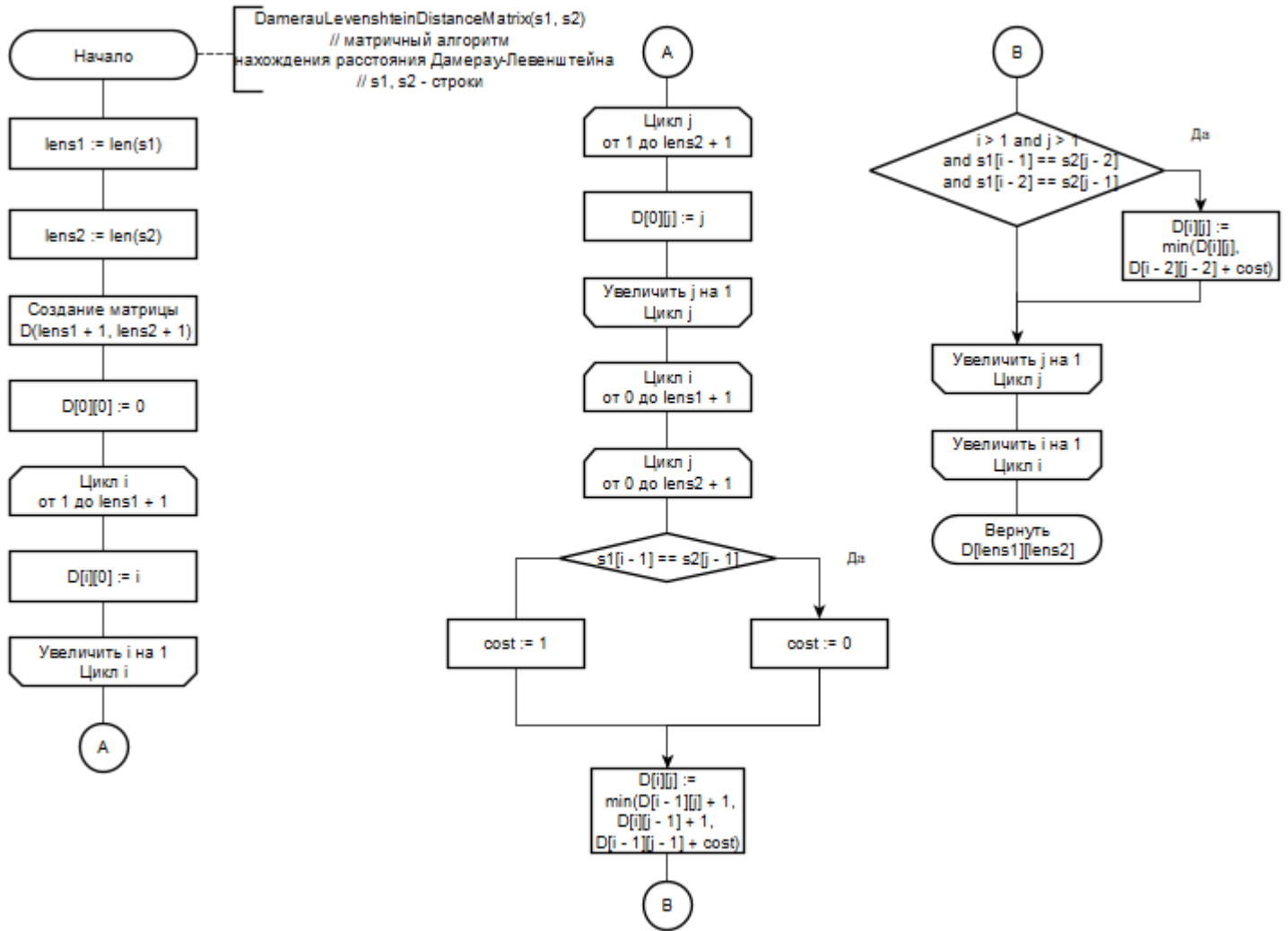


Рис. 2.4. Рекурсивный алгоритм нахождения расстояния Левенштейна

3 | Технологическая часть

В данном разделе будут рассмотрены требования к программному обеспечению, средства реализации, представлен листинг кода, сравнительный анализ матричной и рекурсивной реализации, описание тестирования.

3.1. Требования к программному обеспечению

Требования к вводу:

- 1) на вход подаются две строки;
- 2) uppercase и lowercase буквы считаются разными.

Требования к программе:

- 1) две пустые строки — корректный ввод, программа не должна аварийно завершаться;
- 2) на выходе необходимо получить число, являющиеся результатом работы алгоритма;
- 3) для матричных реализаций требуется вывести матрицу решений;
- 4) требуется замерить время работы каждого из алгоритмов.

3.2. Средства реализации

В качестве языка программирования был выбран Python т.к. я знаком с данным языком, он простой и лаконичный, имеющий немногословный и понятный синтаксис, похожий на псевдокод, обладающий сильной динамической типизацией, которая способствует быстрому написанию кода.

Среда разработки — PyCharm, которая предоставляет умную проверку кода, быстрое выявление ошибок и оперативное исправление, вкупе с автоматическим рефакторингом кода, и богатыми возможностями в навигации.

Матрица создается с помощью функции `eye(m, n)` из библиотеки `numpy` [2].

Время работы алгоритмов было замерено с помощью функции `process_time()` из библиотеки `time` [3].

3.3. Листинг кода

В листингах 3.1-3.4 представлена реализация алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна

Листинг 3.1. Матричный алгоритм нахождения расстояния Левенштейна

```
1 def LevenshteinDistanceMatrix(s1, s2, printMatrix = False):
2     lens1 = len(s1)
3     lens2 = len(s2)
4     D = np.eye(lens1 + 1, lens2 + 1)
5     D[0][0] = 0
6     for i in range(1, lens1 + 1):
7         D[i][0] = i
8         for j in range(1, lens2 + 1):
9             D[0][j] = j
10            for i in range(1, lens1 + 1):
11                for j in range(1, lens2 + 1):
12                    cost = 0 if s1[i - 1] == s2[j - 1] else 1
13                    D[i][j] = min(D[i - 1][j] + 1,
14                                D[i][j - 1] + 1, D[i - 1][j - 1] + cost)
15
16 if printMatrix:
17     print('Matrix:')
18     outputMatrix(s1, s2, D)
19     operations(s1, s2, D)
20
21 return D[lens1][lens2]
```

Листинг 3.2. Рекурсивный алгоритм нахождения расстояния Левенштейна

```
1 def LevenshteinDistanceRecursion(s1, s2):
2     lens1 = len(s1)
3     lens2 = len(s2)
4     if lens1 == 0 or lens2 == 0:
5         tmp = abs(lens1 - lens2)
6     else:
7         cost = 0 if s1[-1] == s2[-1] else 1
8         tmp = min(LevenshteinDistanceRecursion(s1[:-1], s2) + 1,
9                 LevenshteinDistanceRecursion(s1, s2[:-1]) + 1,
10                LevenshteinDistanceRecursion(s1[:-1], s2[:-1]) + cost)
11 return tmp
```

Листинг 3.3. Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы

```
1 def LevenshteinDistanceRecursionMatrix(s1, s2):
2     def calcValue(D, i, j):
3         if D[i][j] != -1:
4             return D[i][j]
5         else:
6             cost = 0 if s1[i - 1] == s2[j - 1] else 1
7             D[i][j] = min(calcValue(D, i - 1, j) + 1,
8                           calcValue(D, i, j - 1) + 1,
9                           calcValue(D, i - 1, j - 1) + cost)
10            return D[i][j]
11
12     lens1 = len(s1)
13     lens2 = len(s2)
14     D = np.full((lens1 + 1, lens2 + 1), -1)
15     for i in range(lens1 + 1):
16         D[i][0] = i
17     for j in range(lens2 + 1):
18         D[0][j] = j
19
20     return calcValue(D, lens1, lens2)
```

Листинг 3.4. Матричный алгоритм нахождения расстояния Дамерау-Левенштейна

```
1 def DamerauLevenshteinDistanceMatrix(s1, s2, printMatrix = False):
2     lens1 = len(s1)
3     lens2 = len(s2)
4     D = np.eye(lens1 + 1, lens2 + 1)
5     D[0][0] = 0
6     for i in range(1, lens1 + 1):
7         D[i][0] = i
8     for j in range(1, lens2 + 1):
9         D[0][j] = j
10    for i in range(1, lens1 + 1):
11        for j in range(1, lens2 + 1):
12            cost = 0 if s1[i - 1] == s2[j - 1] else 1
13            D[i][j] = min(D[i - 1][j] + 1,
14                          D[i][j - 1] + 1, D[i - 1][j - 1] + cost)
15            if i > 1 and j > 1 and s1[i - 1] == s2[j - 2] and
16                s1[i - 2] == s2[j - 1]:
17                D[i][j] = min(D[i][j], D[i - 2][j - 2] + cost)
18
19    if printMatrix:
20        print('Matrix:')
21        outputMatrix(s1, s2, D)
22        operations(s1, s2, D)
23    return D[lens1][lens2]
```

3.4. Сравнительный анализ матричной и рекурсивной реализаций

Рекурсивная реализация работает медленнее по сравнению с матричной из-за повторных вычислений, возникающих в ходе работы рекурсивного алгоритма. Это наглядно видно на Рис. 3.1, иллюстрирующем дерево рекурсивных вызовов. При каждом рекурсивном вызове необходимо передавать в функцию подстроки исходных строк, что затратно по памяти. Эту проблему возможно избежать, если занести строки в глобальные переменные, что является плохой практикой.

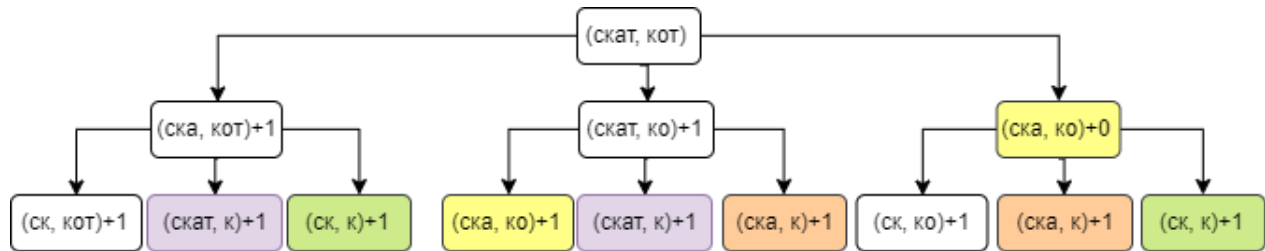


Рис. 3.1. Дерево рекурсивных вызовов

Для сравнения, вызов функции, которая реализует матричный алгоритм, происходит один раз и в функцию только один раз передаются обе строки. В этом алгоритме память будет затрачена на хранение матрицы, а время на вложенные циклы, однако затраты будут существенно меньше, чем при многократных вызовах рекурсивной функции.

Пусть длина строки $S1$ - n , длина строки $S2$ - m , тогда затраты памяти на приведенные выше алгоритмы будут следующими. Затраты памяти матричного алгоритма нахождения расстояния Левенштейна:

- строки $S1, S2$ - $(m + n) * \text{sizeof}(\text{char})$
- матрица - $((m + 1) * (n + 1)) * \text{sizeof}(\text{int})$
- длины строк - $2 * \text{sizeof}(\text{int})$
- вспомогательная переменная - $\text{sizeof}(\text{int})$

Затраты памяти рекурсивного алгоритма нахождения расстояния Левенштейна (для каждого вызова):

- строки $S1, S2$ - $(m + n) * \text{sizeof}(\text{char})$
- длины строк - $2 * \text{sizeof}(\text{int})$
- вспомогательные переменные - $2 * \text{sizeof}(\text{int})$
- адрес возврата

Чтобы получить итоговую оценку затрачиваемой памяти необходимо затрачиваемую память для одного вызова умножить на максимальную глубину рекурсии, которая равна сложению длин входных строк.

Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы:

- матрица - $((m + 1) * (n + 1)) * \text{sizeof}(\text{int})$
- строки S1, S2 - $(m + n) * \text{sizeof}(\text{char})$

Для каждого рекурсивного вызова:

- передача данных - $2 * \text{sizeof}(\text{int}^*) + \text{sizeof}(\text{int}^{**})$
- вспомогательная переменная - $\text{sizeof}(\text{int})$
- адрес возврата

Чтобы получить итоговую оценку затрачиваемой памяти на рекурсивные вызовы необходимо затрачиваемую память для одного рекурсивного вызова умножить на максимальную глубину рекурсии, которая равна сложению длин входных строк.

Затраты памяти матричного алгоритма нахождения расстояния Дамерау-Левенштейна:

- строки S1, S2 - $(m + n) * \text{sizeof}(\text{char})$
- матрица - $((m + 1) * (n + 1)) * \text{sizeof}(\text{int})$
- длины строк - $2 * \text{sizeof}(\text{int})$
- вспомогательная переменная - $\text{sizeof}(\text{int})$

3.5. Описание тестирования

Реализовано модульное тестирование отдельным файлом test.py с помощью библиотеки unittest [4]. Полученные результаты функций сравниваются с контрольными значениями.

Тестирование происходит по следующим данным:

- 1) проверка работы с пустыми строками;
- 2) проверка работы с идентичными строками;
- 3) проверка работы со строками, имеющие совпадающие символы;
- 4) проверка работы с полностью несовпадающими строками.

4 | Экспериментальная часть

В данном разделе приведены примеры работы программы и сравнительный анализ алгоритмов на основе экспериментальных данных.

4.1. Примеры работы

На рис. 4.1 представлено главное меню программы.

```
Меню:  
1) Расстояние Левенштейна матричный алгоритм  
2) Расстояние Левенштейна рекурсивный алгоритм по формуле  
3) Расстояние Левенштейна рекурсивный алгоритм, заполняющий матрицу  
4) Расстояние Дamerau-Левенштейна матричный алгоритм  
5) Сравнение по времени расстояние Левенштейна  
  
Выберите пункт меню:
```

Рис. 4.1. Главное меню программы

На рис. 4.2-4.5 приведены примеры работы программы при вводе строк «тело» и «столб» при выборе пунктов меню 1-4.

```
Выберите пункт меню: 1  
Введите первую строку: тело  
Введите вторую строку: столб  
Matrix:  
  с т о л б  
0 1 2 3 4 5  
т 1 1 1 2 3 4  
е 2 2 2 2 3 4  
л 3 3 3 3 2 3  
о 4 4 4 3 3 3  
Operation:  
I M R M R  
3.0
```

Рис. 4.2. Результат работы матричного алгоритма нахождения расстояния Левенштейна


```

Выберите пункт меню: 2
Введите первую строку: тело
Введите вторую строку: столб
3

```

Рис. 4.3. Результат работы рекурсивного алгоритма нахождения расстояния Левенштейна

```

Выберите пункт меню: 3
Введите первую строку: тело
Введите вторую строку: столб
3

```

Рис. 4.4. Результат работы рекурсивного алгоритма нахождения расстояния Левенштейна с заполнением матрицы

```

Выберите пункт меню: 4
Введите первую строку: тело
Введите вторую строку: столб
Matrix:
      с  т  о  л  б
    0  1  2  3  4  5
т  1  1  1  2  3  4
е  2  2  2  2  3  4
л  3  3  3  3  2  3
о  4  4  4  3  3  3
Operation:
I M R M R
3.0

```

Рис. 4.5. Результат работы матричного алгоритма нахождения расстояния Дameraу-Левенштейна

4.2. Результаты тестирования

Было проведено тестирование программы, результаты которого занесены в Таблицу 4.1, 1 столбец которой - номер тестового случая; 2 и 3 столбцы - строки, поступающие на вход; 4 и 5 столбец - ожидаемый и полученный результат, в которых приняты следующие обозначения:

- 1-ая цифра – результат работы матричного алгоритма нахождения расстояния Левенштейна
- 2-ая цифра – результат работы рекурсивного алгоритма нахождения расстояния Левенштейна
- 3-ая цифра – результат работы рекурсивного алгоритма нахождения расстояния Левенштейна с заполнением матрицы;
- 4-ая цифра – результат работы матричного алгоритма нахождения расстояния Дamerau-Левенштейна.

Таблица 4.1. Результаты тестирования

№	S1	S2	Ожидаемый результат	Полученный результат
1	пустая строка	пустая строка	0, 0, 0, 0	0, 0, 0, 0
2	кот	пустая строка	3, 3, 3, 3	3, 3, 3, 3
3	пустая строка	кот	3, 3, 3, 3	3, 3, 3, 3
4	кот	кот	0, 0, 0, 0	0, 0, 0, 0
5	кот	кит	1, 1, 1, 1	1, 1, 1, 1
6	от	кот	1, 1, 1, 1	1, 1, 1, 1
7	кот	кота	1, 1, 1, 1	1, 1, 1, 1
8	кот	кто	2, 2, 2, 1	2, 2, 2, 1
9	кот	рабы	4, 4, 4, 4	4, 4, 4, 4

Программа успешно прошла все тестовые случаи, все полученные результаты совпали с ожидаемыми.

4.3. Постановка эксперимента по замеру времени

Для произведения замеров времени выполнения реализации алгоритмов будет использована формула:

$$t = \frac{T}{N} \quad (4.1)$$

где t — среднее время выполнения алгоритма, N — количество замеров, T — время выполнения N замеров. Неоднократное измерение времени необходимо для получения более точного результата.

Количество замеров взято равным 100. Эксперимент проводится на рандомных строках одинаковой длины.

4.4. Сравнительный анализ на материале экспериментальных данных

Для сравнения скорости работы необходимо провести замеры процессорного времени для строк одинаковой длины.

Были проведены замеры времени работы каждого из алгоритмов, результаты которых занесены в Таблицу 4.2, столбцы которой обозначают следующее:

- 1 столбец — длина строк;
- 2 столбец — время работы матричного алгоритма нахождения расстояния Левенштейна;
- 3 столбец — время работы рекурсивного алгоритма нахождения расстояния Левенштейна;
- 4 столбец — время работы рекурсивного алгоритма нахождения расстояния Левенштейна с заполнением матрицы;
- 5 столбец — время работы матричного алгоритма нахождения расстояния Дамерау-Левенштейна.

Таблица 4.2. Сравнение алгоритмов по времени(в секундах)

len	LevMatr	LevRec	LevRecMatr	DamLevMatr
3	0.00000750	0.00003239	0.00002969	0.00000837
5	0.00008750	0.00123438	0.00012969	0.00008906
6	0.00009050	0.00155438	0.00014069	0.00009596
8	0.00024553	2.00000438	0.00035378	0.00026596
10	0.00029531	5.73437500	0.00046875	0.00032188

Вывод: У матричных алгоритмов нахождения расстояния Левенштейна время работы пропорционально квадрату длины строк (при увеличении строки в два раза, время увеличивается в четыре раза). Рекурсивный алгоритм по нахождению расстояние Левенштейна показывает наихудшее время, что связано с повторным вычислением одних и тех же значений (пример рис 2.5). Матричный алгоритм по нахождению расстояния Левенштейна и Дамерау-Левенштейна показывают практически одинаковое значение, однако необходимо учитывать то, что Дамерау-Левенштейн решает другую задачу.

Заключение

В ходе работы были изучены алгоритмы нахождения расстояний Левенштейна и Дамерау–Левенштейна. Реализованы 4 алгоритма поиска этих расстояний, приведен программный код реализации алгоритмов нахождения расстояний.

Было выполнено сравнение разных алгоритмов нахождения расстояния Левенштейна по затраченным ресурсам. Было установлено, что рекурсивный алгоритм занимает гораздо меньше памяти при работе со строками большой длины, чем матричные алгоритмы. Однако матричные алгоритмы отмечаются своим быстродействием.

Цель работы достигнута. Алгоритмы нахождения расстояния Левенштейна и Дамерау–Левенштейна применить на практике, получены практические навыки реализации этих алгоритмов.

Список литературы

1. Дж. Макконнелл. Анализ алгоритмов. Активный обучающий подход. – М.: Техносфера, 2017. – 267 с.
2. Библиотека numpy Python, документация [электронный ресурс] – Режим доступа: <https://pythonworld.ru/numpy.html>, свободный – (Дата обращения: 16.09.20)
3. Официальный сайт Python, документация [электронный ресурс]. Режим доступа: <https://docs.python.org/3/library/time.html>, свободный (Дата обращения: 16.09.20)
4. Официальный сайт Python, документация [электронный ресурс] – Режим доступа: <https://docs.python.org/3/library/unittest.html>, свободный – (Дата обращения: 16.09.20)