



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»

Отчет по лабораторной работе № 1 по курсу "Анализ алгоритмов"

Тема Расстояния Левенштейна и Дameraу-Левенштейна

Студент Хамзина Р. Р.

Группа ИУ7-53Б

Оценка (баллы) _____

Преподаватель Волкова Л. Л.

Содержание

| | |
|---|-----------|
| Введение | 3 |
| 1 Аналитическая часть | 5 |
| 1.1 Рекурсивный алгоритм нахождения расстояния Левенштейна | 5 |
| 1.2 Матричный алгоритм нахождения расстояния Левенштейна | 6 |
| 1.3 Рекурсивный алгоритм нахождения расстояния Левенштейна с использованием матрицы | 7 |
| 1.4 Алгоритм нахождения расстояния Дамерау-Левенштейна . . | 8 |
| 1.5 Вывод | 9 |
| 2 Конструкторская часть | 10 |
| 2.1 Разработка алгоритмов | 10 |
| 2.2 Описание используемых типов и структур данных | 16 |
| 2.3 Использование памяти | 16 |
| 2.4 Структура разрабатываемого ПО | 17 |
| 2.5 Классы эквивалентности при тестировании | 18 |
| 2.6 Вывод | 18 |
| 3 Технологическая часть | 19 |
| 3.1 Средства реализации | 19 |
| 3.2 Листинги кода | 19 |
| 3.3 Функциональные тесты | 22 |
| 3.4 Вывод | 23 |
| 4 Исследовательская часть | 24 |
| 4.1 Технические характеристики | 24 |
| 4.2 Демонстрация работы программы | 25 |
| 4.3 Время выполнения алгоритмов | 25 |
| 4.4 Вывод | 29 |
| Заключение | 30 |
| Список литературы | 31 |

Введение

Задача поиска и сравнения последовательностей возникает при обработке больших объемов информации, при работе с неструктурированными данными или при поисковых запросах. Алгоритмы для решения этой задачи разделяются на три группы: точный поиск подстрок, неточный поиск и поиск наибольшей общей подпоследовательности. Алгоритмы нечеткого поиска применяются в компьютерной лингвистике (коррекция ошибок, поиск с учетом формы одного и того же слова) и в биоинформатике (сравнение ДНК). В основе методов лежит определение расстояния между строками, которое задается в конкретном алгоритме.

Данная задача впервые была поставлена советским математиком Владимиром Левенштейном для двоичных кодов. Расстояние между строками для произвольного алфавита называли расстоянием Левенштейна и определили как минимальное количество операций вставки, удаления и замены одного символа, необходимых для превращения одной строки в другую [1].

Модификацией алгоритма Левенштейна является алгоритм Дамерау-Левенштейна. В данном алгоритме к операциям, необходимым для перевода одной строки в другую, добавляется транспозиция символов (перестановка двух соседних символов).

Данные алгоритмы реализуются методом динамического программирования. Существуют их рекурсивная и нерекурсивная реализации, которые различаются временной эффективностью.

Целью данной лабораторной работы является изучение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна. Для достижения поставленной цели требуется выполнить следующие задачи:

- изучить алгоритмы Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками;
- привести схемы изучаемых алгоритмов;
- описать используемые типы и структуры данных;
- оценить объем памяти для хранения данных;

- описать структуру разрабатываемого программного обеспечения;
- определить средства программной реализации выбранных алгоритмов;
- реализовать разработанные алгоритмы;
- провести функциональное тестирование программного обеспечения;
- провести сравнительный анализ по времени реализованных алгоритмов;
- подготовить отчет о выполненной лабораторной работе.

1 Аналитическая часть

В данном разделе будут описаны алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна.

1.1 Рекурсивный алгоритм нахождения расстояния Левенштейна

Рассмотрим перевод строки $A' = A[: i]$ в строку $B' = B[: j]$, последние символы которых соответственно a_{i-1} и b_{j-1} , при помощи следующих операций, каждая из которых имеет свою стоимость (x, y - символы, λ - строка):

- вставка символа в произвольное место - I , стоимостью $w(x, \lambda) = 1$;
- удаление символа с произвольной позиции - D , $w(\lambda, y) = 1$;
- замена символа на другой - R , $w(x, y) = 1, x \neq y$.

Расстояние Левенштейна - минимальное количество операций I, D, R для перевода $A' = A[: i]$ в $B' = B[: j]$ [1]. Его можно подсчитать по следующей рекуррентной формуле:

$$D(i, j) = \begin{cases} 0 & i = 0, j = 0 \\ i & j = 0, i > 0 \\ j & i = 0, j > 0 \\ \min\{ \\ \quad D(i, j - 1) + 1 \\ \quad D(i - 1, j) + 1 & i > 0, j > 0 \\ \quad D(i - 1, j - 1) + m(a_{i-1}, b_{j-1}) \quad (2.1) \\ \} \end{cases} \quad (1.1)$$

Функция 2.1 определена как:

$$m(x, y) = \begin{cases} 0 & \text{если } x = y, \\ 1 & \text{иначе} \end{cases}. \quad (1.2)$$

Рекурсия используется для реализации формулы 1.1. Пусть $D(i, j)$ - есть расстояние Левенштейна между строками $A' = A[: i]$ и $B' = B[: j]$. $B' = B[: j]$ может быть получен из $A' = A[: i]$ следующими способами:

- если символ a_{j-1} был добавлен при редактировании, то необходимо $A[: i]$ превратить в $B[: j - 1]$, на что необходимо $D(i, j - 1)$ операций редактирования. В этом случае $D(i, j) = D(i, j - 1) + 1$;
- если символ a_{i-1} был удален при редактировании, тогда необходимо $A[: i - 1]$ превратить в $B[: j]$, на что необходимо $D(i - 1, j)$ операций редактирования. В этом случае $D(i, j) = D(i - 1, j) + 1$;
- если последние символы префиксов совпадают, т. е. $a_{i-1} = b_{j-1}$, то в этом случае можно не менять эти последние символы. Тогда $D(i, j) = D(i - 1, j - 1)$;
- если $a_{i-1} \neq b_{j-1}$, то тогда можно потратить 1 операцию на замену символа a_{i-1} на b_{j-1} и также потратить операцию $D(i - 1, j - 1)$ на превращение $A[: i - 1]$ в $B[: j - 1]$. Тогда $D(i, j) = D(i - 1, j - 1) + 1$.

Далее при вычислении необходимо взять минимум из всех перечисленных возможностей (при этом из случаев 3 или 4 рассматривается только один, в зависимости от условия $a_{i-1} = b_{j-1}$).

Начальные значения: $D(i, 0) = i$, $D(0, j) = j$.

1.2 Матричный алгоритм нахождения расстояния Левенштейна

Промежуточные значения $D(i, j)$ вычисляются несколько раз, поэтому при больших значениях i, j временная эффективность алгоритма сни-

жается. Для решения этой проблемы используют матрицу для хранения значений расстояний.

Матрицу M размером $N+1$ и $L+1$, где N и L - длины строк $A' = A[:i]$ и $B' = B[:j]$, заполняют по следующим правилам:

- $M[0][0] = 0$;
- $M[i][j]$ вычисляется по формуле 1.3.

$$M[i][j] = \min \begin{cases} M[i][j-1] + 1 \\ M[i-1][j] + 1 \\ M[i-1][j-1] + m(A[i-1], B[j-1]) \end{cases} \quad (1.3)$$

$$(2.2)$$

Функция 2.2 определена как:

$$m(x, y) = \begin{cases} 0 & \text{если } x = y, \\ 1 & \text{иначе} \end{cases}. \quad (1.4)$$

$M[N][L]$ содержит значение расстояния Левенштейна.

1.3 Рекурсивный алгоритм нахождения расстояния Левенштейна с использованием матрицы

Для оптимизации рекурсивного алгоритма можно использовать матрицу как промежуточный буфер, в который при выполнении рекурсии заносятся значения расстояния Левенштейна для еще не рассмотренных символов. Рассмотренные ранее символы в матрицу не добавляются.

1.4 Алгоритм нахождения расстояния Дамерау-Левенштейна

Как было сказано ранее, алгоритм Дамерау-Левенштейна отличается от алгоритма Левенштейна добавлением операции транспозиции символов при переводе одной строки в другую. Расстояние Дамерау-Левенштейна может быть найдено по формуле:

$$d(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ \\ \quad d(i, j - 1) + 1, \\ \quad d(i - 1, j) + 1, \\ \quad d(i - 1, j - 1) + m(a_{i-1}, b_{j-1}), & \text{иначе} \\ \quad \left[\begin{array}{ll} d(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a_i = b_{j-1}; \\ & b_j = a_{i-1} \\ & \infty, & \text{иначе} \end{array} \right. \\ \} \end{cases}, \quad (1.5)$$

Пусть $d(i, j)$ - есть расстояние Дамерау-Левенштейна между строками $A' = A[: i]$ и $B' = B[: j]$. $B' = B[: j]$ может быть получен из $A' = A[: i]$ следующими способами:

- если символ a_{i-1} был удален при редактировании, тогда необходимо $A[: i - 1]$ превратить в $B[: j]$, на что необходимо $d(i - 1, j)$ операций редактирования. В этом случае $d(i, j) = d(i - 1, j) + 1$;
- если символ a_{j-1} был добавлен при редактировании, то необходимо $A[: i]$ превратить в $B[: j - 1]$, на что необходимо $d(i, j - 1)$ операций редактирования. В этом случае $d(i, j) = d(i, j - 1) + 1$;
- если последние символы префиксов совпадают, т. е. $a_{i-1} = b_{j-1}$, то в

этом случае можно не менять эти последние символы. Тогда $d(i, j) = d(i - 1, j - 1)$;

- если $a_{i-1} \neq b_{j-1}$, то тогда можно потратить 1 операцию на замену символа a_{i-1} на b_{j-1} и также потратить операцию $D(i - 1, j - 1)$ на превращение $A[: i - 1]$ в $B[: j - 1]$. Тогда $d(i, j) = d(i - 1, j - 1) + 1$;
- если последние два символа a_i и b_j были переставлены, то $d(i, j)$ может быть равно $d(i - 2, j - 2) + 1$.

Далее при вычислении необходимо взять минимум из всех перечисленных возможностей (при этом из случаев 3 или 4 рассматривается только один, в зависимости от условия $a_{i-1} = b_{j-1}$).

Начальные значения: $D(i, 0) = i$, $D(0, j) = j$.

1.5 Вывод

Были изучены способы нахождения расстояния перевода одной строки в другую при помощи рекурсивной и матричной версий алгоритма Левенштейна, алгоритма Левенштейна с использованием матрицы и рекурсивной версии алгоритма Дамерау-Левенштейна.

Программе, реализующей данные алгоритмы, на вход будет подаваться две строки. Выходными данными такой программы должно быть число - редакционное расстояние. Программа должна работать в рамках следующих ограничений: строки, поданные на вход, могут состоять из букв русского и английского языков. Кроме того, должен быть выдан корректный результат при вводе пустых строк.

Пользователь должен иметь возможность выбора алгоритма нахождения редакционного расстояния и вывода его значения на экран. Также должны быть реализованы сравнение алгоритмов по времени работы с выводом результатов на экран и получение графического представления результатов сравнения. Данные действия пользователь должен выполнять при помощи меню.

2 Конструкторская часть

В данном разделе будут представлены схемы алгоритмов вычисления расстояний Левенштейна и Дamerau-Левенштейна. Будут описаны типы и структуры данных, используемые для реализации, а также структура разрабатываемого программного обеспечения, оценено использование памяти. Кроме того будут выделены классы эквивалентности для тестирования.

2.1 Разработка алгоритмов

Схемы рекурсивного и матричного алгоритмов поиска расстояния Левенштейна представлены на рисунках 2.1-2.2 соответственно. Для алгоритма с использованием матрицы на рисунке 2.4 приведена схема рекурсивной части, схема всего алгоритма - на рисунке 2.3. Схема алгоритма Дamerau-Левенштейна приведена на рисунке 2.5.

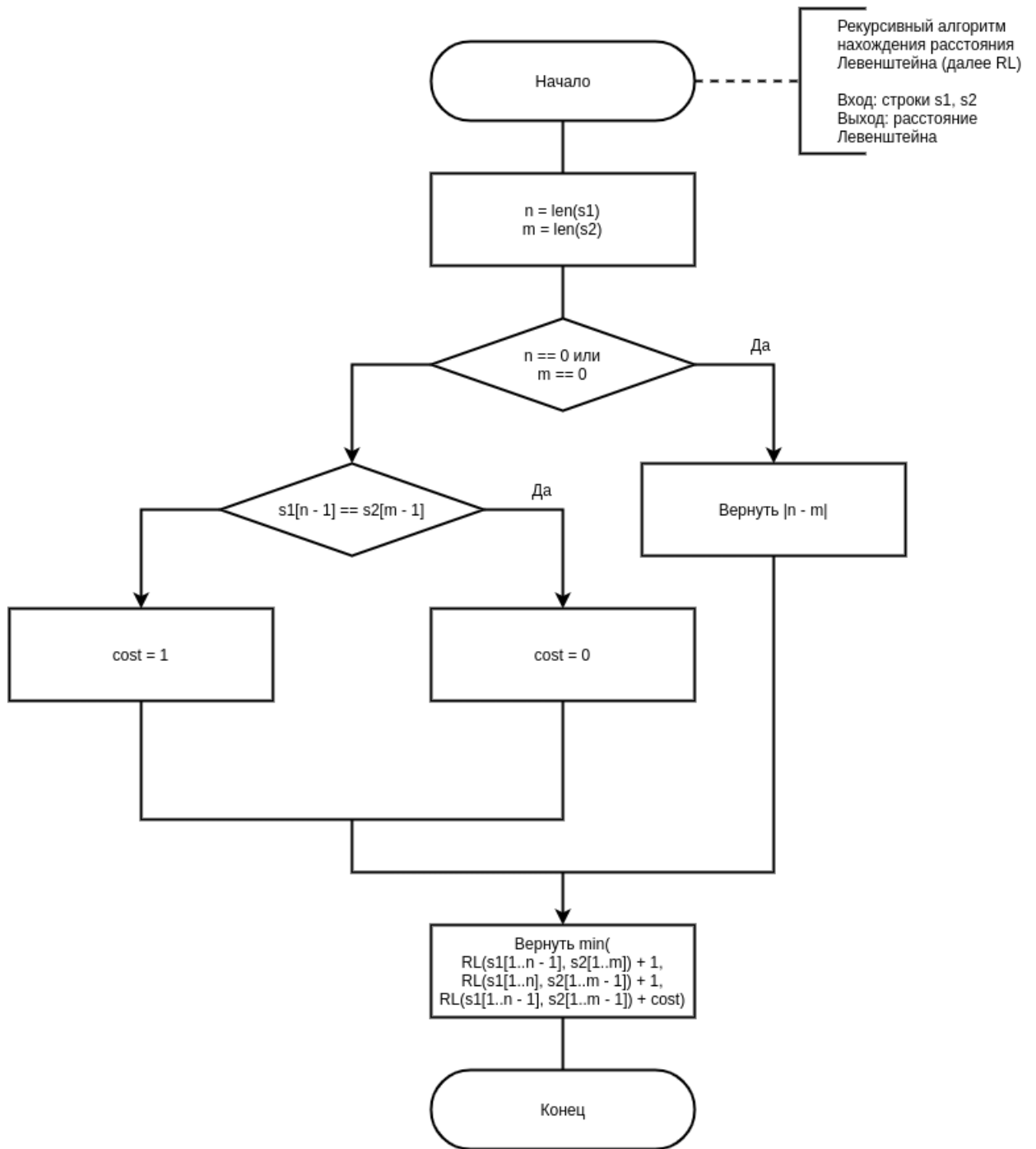


Рисунок 2.1 – Схема рекурсивного алгоритма нахождения расстояния Левенштейна

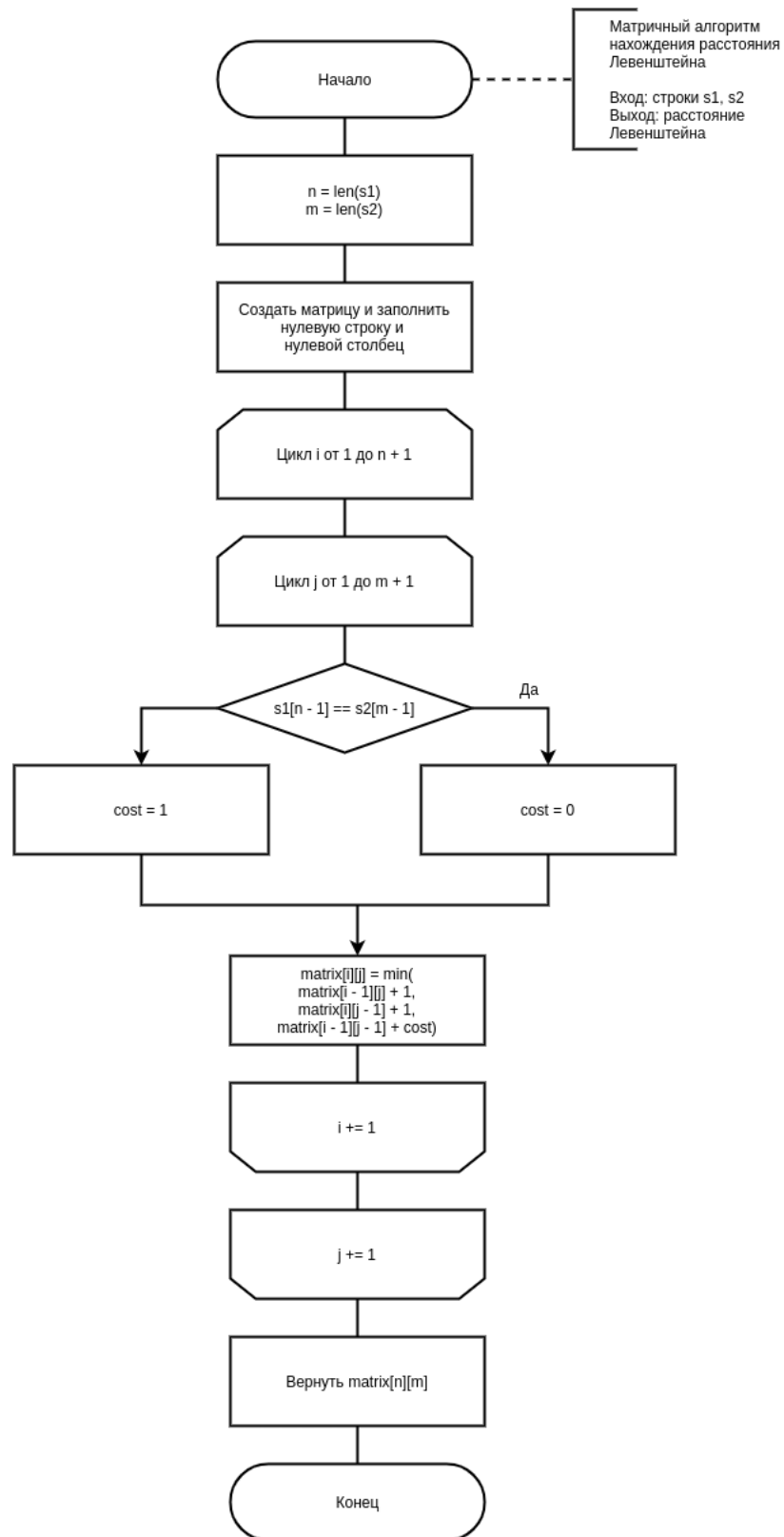


Рисунок 2.2 – Схема матричного алгоритма нахождения расстояния Левенштейна

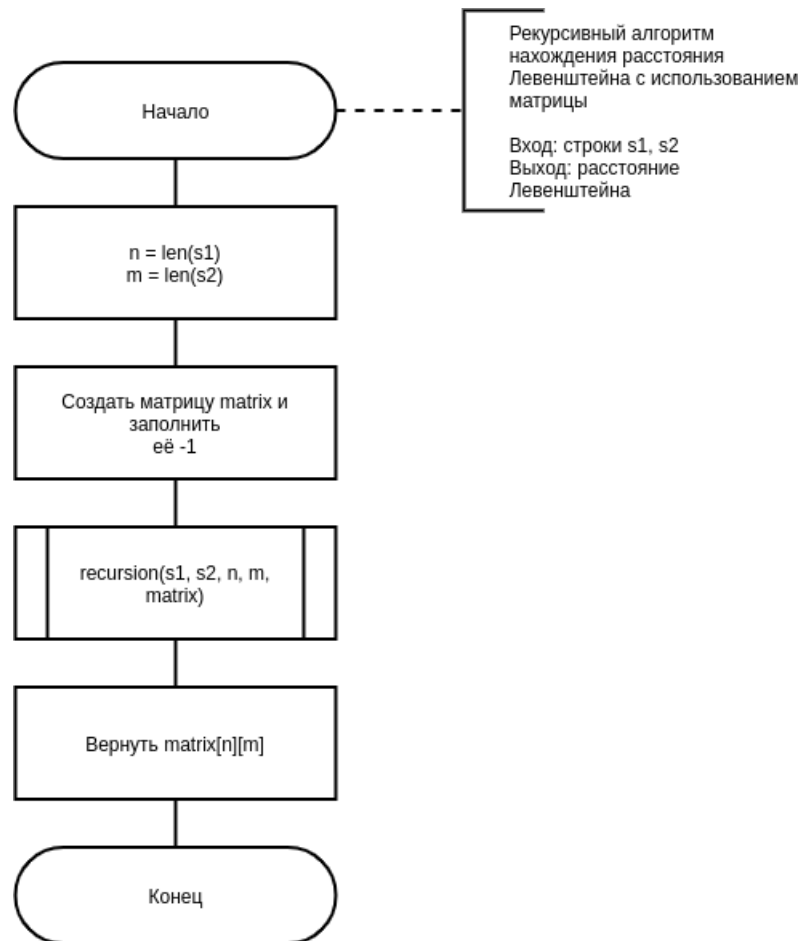


Рисунок 2.3 – Схема алгоритма нахождения расстояния Левенштейна с использованием матрицы

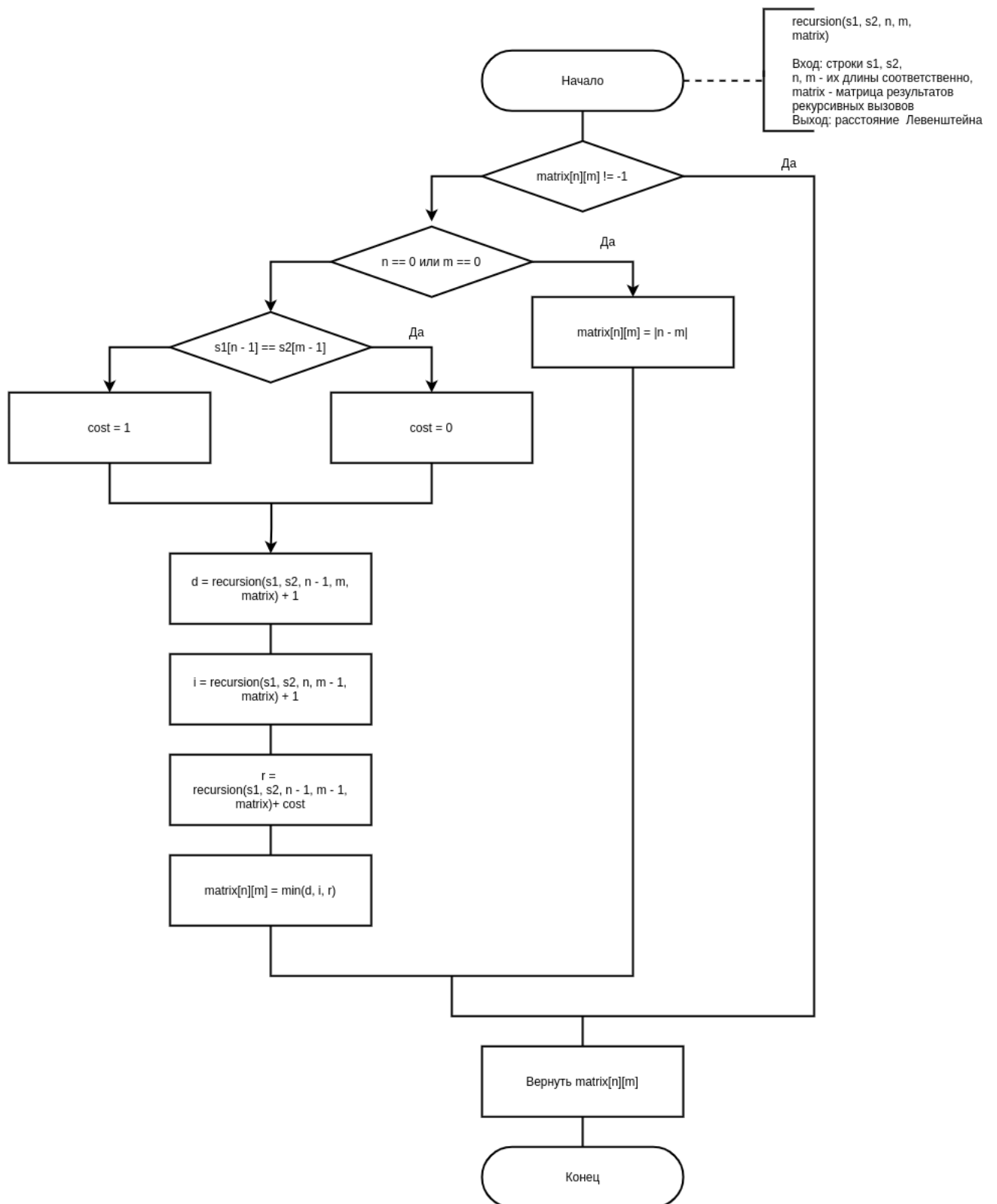


Рисунок 2.4 – Схема рекурсивной части алгоритма нахождения расстояния Левенштейна с использованием матрицы

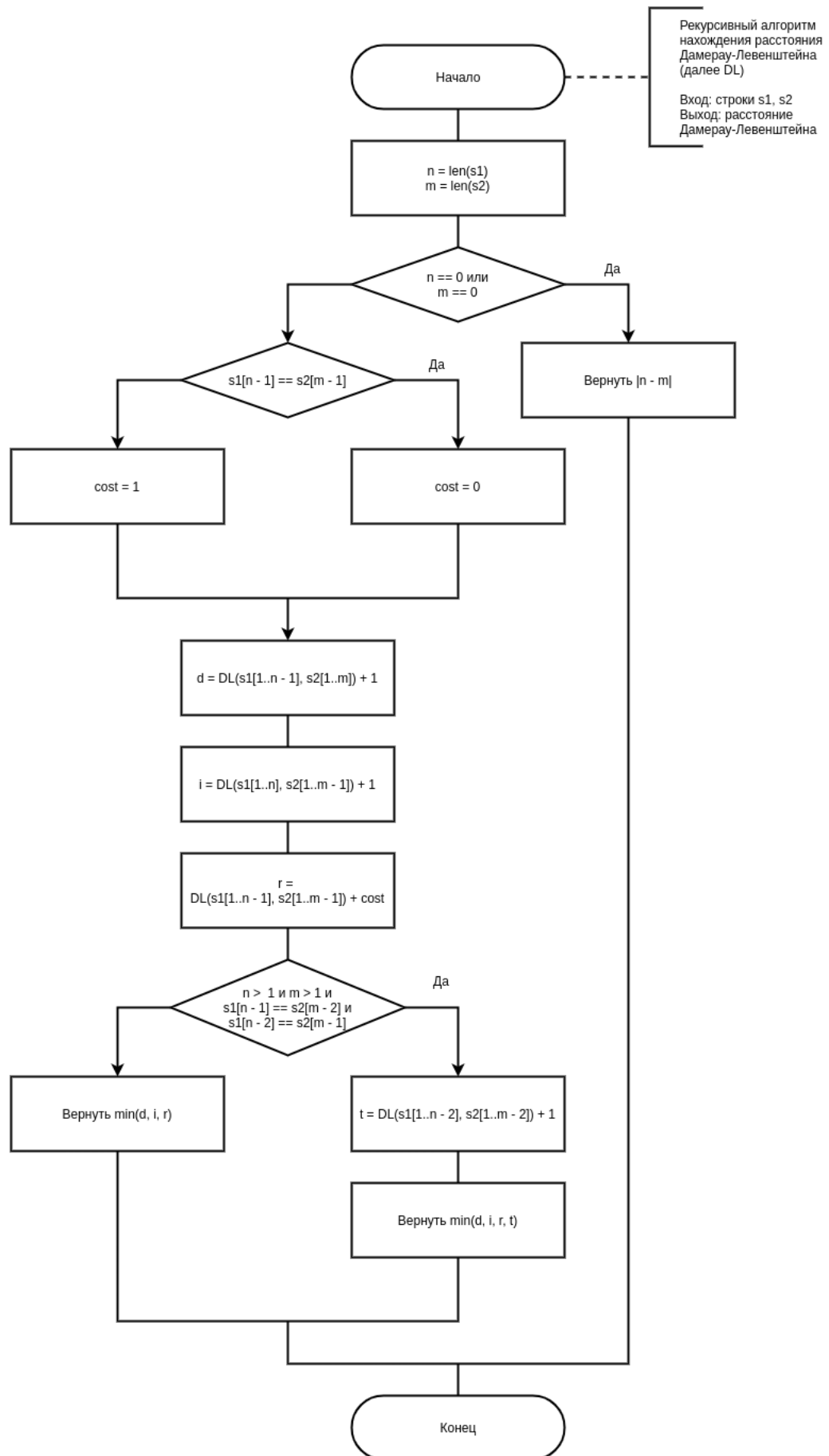


Рисунок 2.5 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна

2.2 Описание используемых типов и структур данных

Для реализации рассмотренных алгоритмов будут использованы следующие типы данных:

- *str* - для двух строк, поданных на вход;
- *int* - для длин каждой из строк, поданных на вход, и для возвращаемого значения искомого расстояния.

Для реализации матричной версии алгоритма нахождения расстояния Левенштейна и версии с использованием матрицы используется двумерный список значений типа *int*.

2.3 Использование памяти

Алгоритм нахождения расстояния Дамерау-Левенштейна отличается от алгоритма нахождения расстояния Левенштейна добавлением операции транспозиции, что не влияет на объем используемой памяти. Поэтому проведем анализ используемой памяти для рекурсивного и нерекурсивного версий алгоритма Левенштейна.

Для рекурсивной версии при каждом вызове происходит выделение памяти под:

- 2 строки типа *str*;
- 2 длины строк типа *int*;
- возвращаемое значение типа *int*;
- адрес возврата типа *int*;
- локальная переменная типа *int*.

При этом глубина рекурсии равна сумме длин двух строк. Таким образом, используемая память рекурсивной версии равна:

$$M_r = (5 \cdot \text{sizeof}(\text{int}) + 2 \cdot \text{sizeof}(\text{str})) \cdot (|s1| + |s2|) \quad (2.1)$$

Для нерекурсивной версии происходит выделение памяти под:

- 2 строки типа *str*;
- 2 длины строк типа *int*;
- матрицу с размерами, равными длинам строк, увеличенным на 1;
- возвращаемое значение типа *int*;
- адрес возврата типа *int*;
- 5 локальных переменных типа *int*.

Таким образом, используемая память нерекурсивной версии равна:

$$M_m = (9 + (|s1| + 1) \cdot (|s2| + 1)) \cdot \text{sizeof}(\text{int}) + 2 \cdot \text{sizeof}(\text{str}) \quad (2.2)$$

2.4 Структура разрабатываемого ПО

Программа состоит из следующих модулей:

- `main.py` - главный файл программы, предоставляющий пользователю меню для выполнения основных функций;
- `distances.py` - файл, содержащий алгоритмы нахождения редакционного расстояния;
- `str.py` - файл, содержащий функции работы со строками;
- `matrix.py` - файл, содержащий функции работы с матрицами;
- `time_test.py` - файл, содержащий функции замеров времени работы указанных алгоритмов;

- graph_result.py - файл, содержащий функции визуализации временных характеристик описанных алгоритмов.

2.5 Классы эквивалентности при тестировании

Для тестирования разрабатываемой программы будут выделены следующие классы эквивалентности:

- две пустые строки;
- одна пустая строка, одна нет;
- удаление/добавление символа в конец;
- удаление/добавление символа внутри строки;
- замена символов;
- транспозиция символов.

2.6 Вывод

Были представлены схемы алгоритмов вычисления расстояний Левенштейна и Дамерау-Левенштейна. Были указаны типы и структуры данных, используемые для реализации, и описана структура разрабатываемого программного обеспечения. Был проанализирован объем используемой памяти для рекурсивной и нерекурсивной версий алгоритмов. Также были выделены классы эквивалентности для тестирования ПО.

3 Технологическая часть

В данном разделе будут указаны средства реализации, будут представлены листинги кода, а также функциональные тесты.

3.1 Средства реализации

Реализация данной лабораторной работы выполнялась при помощи языка программирования Python [2]. Выбор ЯП обусловлен простотой синтаксиса, большим числом библиотек и эффективностью визуализации данных.

Замеры времени проводились при помощи функции `process_time` из библиотеки `time` [3].

3.2 Листинги кода

Реализации алгоритмов нахождения расстояния Левенштейна и расстояния Дамерау-Левенштейна представлены в листингах 3.1-3.4.

Листинг 3.1 – Рекурсивный алгоритм нахождения расстояния Левенштейна

```
1 def recursive(str1, str2):
2     n = len(str1)
3     m = len(str2)
4
5     if n == 0 or m == 0:
6         return abs(n - m)
7
8     cost_replacement = Cost.replacement if str1[-1] != str2[-1]
9         else 0
10
11     deletion = recursive(str1[:-1], str2) + Cost.deletion
12     insertion = recursive(str1, str2[:-1]) + Cost.insertion
13     replacement = recursive(str1[:-1], str2[:-1]) +
14         cost_replacement
15
16     min_distance = min(deletion, insertion, replacement)
17
18     return min_distance
```

Листинг 3.2 – Матричный алгоритм нахождения расстояния Левенштейна

```
1 def matrix_(str1, str2):
2     n, m = len(str1), len(str2)
3     matrix = create_matrix(n + 1, m + 1)
4
5     for i in range(1, n + 1):
6         for j in range(1, m + 1):
7
8             cost_replacement = Cost.replacement if str1[-1] !=
9                 str2[-1] else 0
10             matrix[i][j] = min(matrix[i - 1][j] + Cost.deletion,
11                 matrix[i][j - 1] + Cost.insertion,
12                 matrix[i - 1][j - 1] +
13                     cost_replacement)
14
15     return matrix[n][m]
```

Листинг 3.3 – Алгоритм нахождения расстояния Левенштейна с использованием матрицы

```
1 def recursive_(str1, str2, n, m, matrix):
2     if matrix[n][m] != -1:
3         return matrix[n][m]
4
5     if n == 0:
6         matrix[n][m] = m
7         return matrix[n][m]
8
9     if n > 0 and m == 0:
10        matrix[n][m] = n
11        return matrix[n][m]
12
13    cost_replacement = Cost.replacement if str1[n - 1] != str2[m
14        - 1] else 0
15
16    deletion = recursive_(str1, str2, n - 1, m, matrix) +
17        Cost.deletion
18
19    insertion = recursive_(str1, str2, n, m - 1, matrix) +
20        Cost.insertion
21
22    replacement = recursive_(str1, str2, n - 1, m - 1, matrix) +
23        cost_replacement
24
25    matrix[n][m] = min(deletion,
26        insertion,
27        replacement)
28
29    return matrix[n][m]
30
31 def recursive_with_cache(str1, str2):
32     n, m = len(str1), len(str2)
33     matrix = [[-1] * (m + 1) for _ in range(n + 1)]
34
35     recursive_(str1, str2, n, m, matrix)
36
37     return matrix[n][m]
```

Листинг 3.4 – Рекурсивный алгоритм нахождения расстояния
Дамерау-Левенштейна

```
1 def damerau_levenshtein(str1, str2):
2     n = len(str1)
3     m = len(str2)
4
5     if n == 0 or m == 0:
6         return abs(n - m)
7
8     cost_replacement = Cost.replacement if str1[-1] != str2[-1]
9         else 0
10
11     deletion = damerau_levenshtein(str1[:-1], str2) +
12         Cost.deletion
13     insertion = damerau_levenshtein(str1, str2[:-1]) +
14         Cost.insertion
15     replacement = damerau_levenshtein(str1[:-1], str2[:-1]) +
16         cost_replacement
17
18     if n > 1 and m > 1 and str1[-1] == str2[-2] and str1[-2] ==
19         str2[-1]:
20         transposition = damerau_levenshtein(str1[:-2], str2[:-2])
21             + Cost.transposition
22         return min(deletion, insertion, replacement,
23             transposition)
24
25     return min(deletion, insertion, replacement)
```

3.3 Функциональные тесты

В таблице 3.1 приведены функциональные тесты для функций, реализующих алгоритмы нахождения редакционного расстояния. Все тесты пройдены успешно.

Таблица 3.1 – Функциональные тесты

| № | Строка 1 | Строка 2 | Ожидаемый результат | |
|----|-----------------|-----------------|---------------------|--------------------|
| | | | Левенштейн | Дамерау-Левенштейн |
| 1 | "пустая строка" | "пустая строка" | 0 | 0 |
| 2 | "пустая строка" | привет | 6 | 6 |
| 3 | солнышко | "пустая строка" | 8 | 8 |
| 4 | улыбка | улыбк | 1 | 1 |
| 5 | рай | раааай | 3 | 3 |
| 6 | аромат | игра | 5 | 5 |
| 7 | душа | счастье | 7 | 7 |
| 8 | закат | азкта | 4 | 2 |
| 9 | мир | рим | 2 | 2 |
| 10 | оскар | окр | 2 | 2 |

3.4 Вывод

Были реализованы функции алгоритмов поиска редакционного расстояния. Было проведено функциональное тестирование указанных функций.

4 Исследовательская часть

В данном разделе будут приведены примеры работы программы, и будет проведен сравнительный анализ реализованных алгоритмов поиска редакционного расстояния по затраченному процессорному времени.

4.1 Технические характеристики

Тестирование проводилось на устройстве со следующими техническими характеристиками:

- операционная система: Ubuntu 20.04.1 Linux x86_64 [4];
- память : 8 GiB;
- процессор: AMD® Ryzen™ 3 3200u @ 2.6 GHz [5].

Тестирование проводилось на ноутбуке, включенном в сеть электропитания. Во время тестирования ноутбук был нагружен только встроенными приложениями окружения, а также непосредственно системой тестирования.

4.2 Демонстрация работы программы

На рисунке 4.1 приведен пример работы программы.

```
МЕНЮ:
Найти расстояние Левенштейна при помощи:
    1 - рекурсивного алгоритма;
    2 - матричного алгоритма;
    3 - рекурсивного алгоритма с использованием матрицы.
Найти расстояние Дамерау-Левенштейна при помощи:
    4 - рекурсивного алгоритма.
    5 - Построить графики.
    6 - Замерить время.
    7 - Найти расстояние при помощи всех алгоритмов.
    0 - Выход.
Выбор: 3
Введите строку № 1: душа
Введите строку № 2: игра

0  0  и  г  р  а
0  0  1  2  3  4
д  1  1  2  3  4
у  2  2  2  3  4
ш  3  3  3  3  4
а  4  4  4  4  3

Полученное редакционное расстояние: 3
```

Рисунок 4.1 – Пример работы программы

4.3 Время выполнения алгоритмов

Функция `process_time` из библиотеки `time` ЯП Python возвращает процессорное время в секундах - значение типа `float`.

Для замера времени:

- получить значение времени до начала выполнения алгоритма, затем после её окончания. Чтобы получить результат, необходимо вычесть из второго значения первое;

- первый шаг необходимо повторить `iters` раз (в программе `iters` равно 100), суммируя полученные значения, а затем усреднить результат.

Замеры проводились для длины слов от 0 до 8 на случайных входных строках. Результаты измерения времени приведены в таблице 4.1 (в мс).

Таблица 4.1 – Результаты замеров времени

| Длина | Рек-ый (Л.) | Матр-ый | Рек-ый с кэшем | Рек-ый (Д.-Л.) |
|-------|-------------|---------|----------------|----------------|
| 0 | 0.0022 | 0.0053 | 0.0019 | 0.0010 |
| 1 | 0.0062 | 0.0073 | 0.0047 | 0.0029 |
| 2 | 0.0214 | 0.0130 | 0.0076 | 0.0112 |
| 3 | 0.0645 | 0.0198 | 0.0148 | 0.0469 |
| 4 | 0.1988 | 0.0263 | 0.0223 | 0.2364 |
| 5 | 1.0346 | 0.0458 | 0.0320 | 1.2427 |
| 6 | 5.4673 | 0.0290 | 0.0423 | 7.0460 |
| 7 | 35.5647 | 0.0331 | 0.0639 | 39.0492 |
| 8 | 189.6622 | 0.0419 | 0.0729 | 199.8746 |

На рисунках 4.2-4.4 приведены графические результаты сравнения временных характеристик.

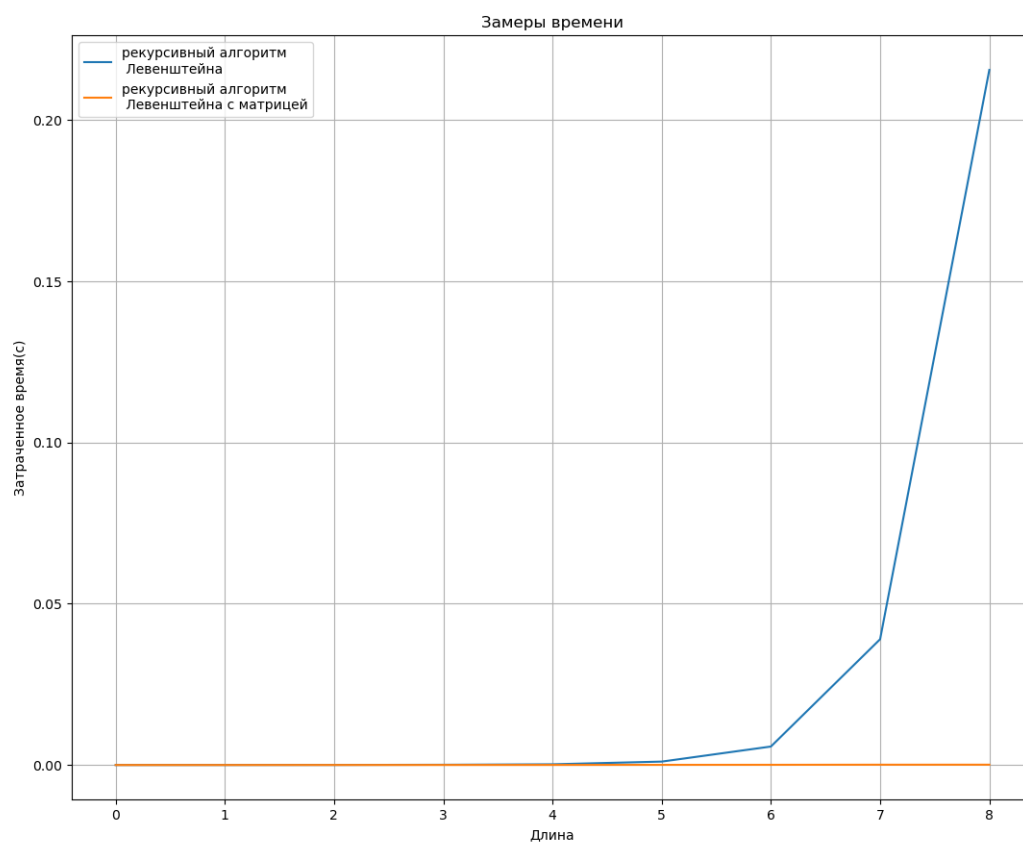


Рисунок 4.2 – Сравнение по времени рекурсивного алгоритма Левенштейна и рекурсивного алгоритма Левенштейна с использованием кэша

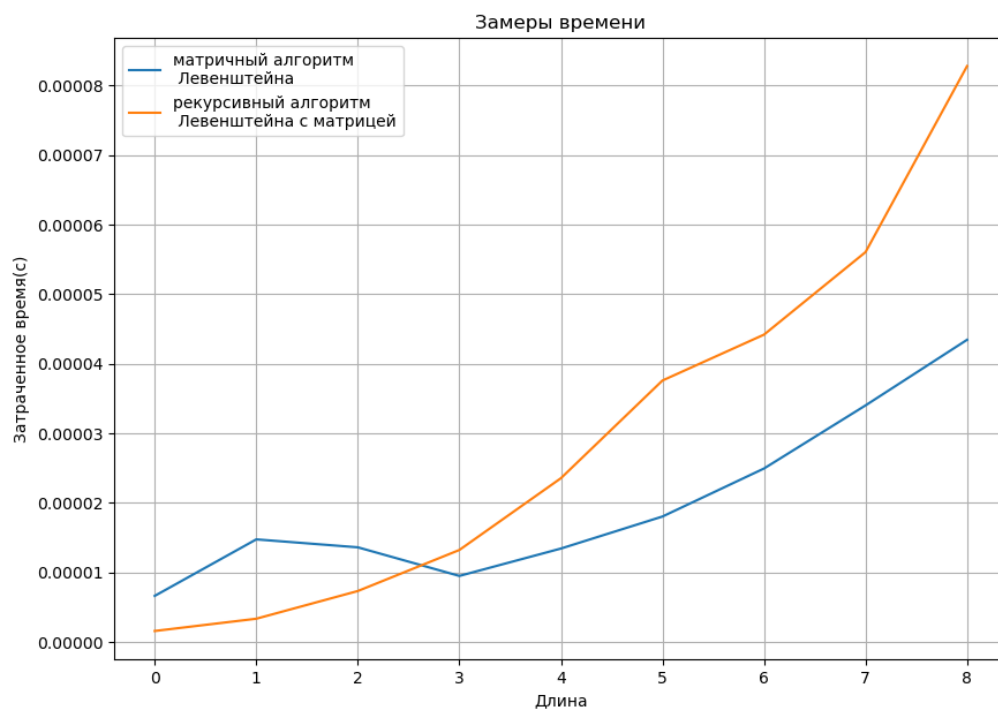


Рисунок 4.3 – Сравнение по времени матричного алгоритма Левенштейна и алгоритма Левенштейна с использованием кэша

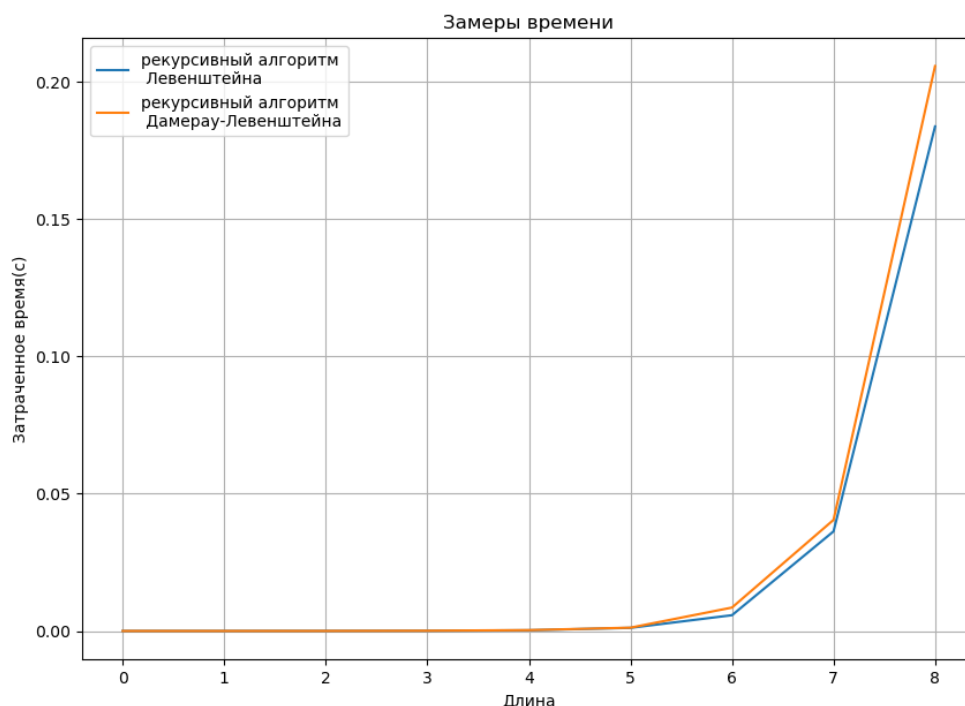


Рисунок 4.4 – Сравнение по времени рекурсивного алгоритма Левенштейна и рекурсивного алгоритма Дамерау-Левенштейна

4.4 Вывод

Исходя из замеров по памяти, итеративные алгоритмы проигрывают рекурсивным, потому что максимальный размер памяти в них растет, как произведение длин строк, а в рекурсивных - как сумма длин строк.

В результате эксперимента было получено, что при длине строк в 4 символа, рекурсивная реализация алгоритма Левенштейна становится медленнее матричной в 9 раз и с увеличением длины строк время работы растет в геометрической прогрессии. Тогда, для строк длиной более 4 символов необходимо использовать матричную версию алгоритма поиска редакционного расстояния.

Также в результате эксперимента было установлено, что при длине строк в более 5 символов, алгоритм Левенштейна работает быстрее Дамерау-Левенштейна в 1.3 раза. Можно сделать вывод, что при таких данных предпочтительно использовать алгоритм Левенштейна.

Заключение

В результате оценки алгоритмов по используемой памяти можно сделать вывод, что матричные реализации алгоритмов поиска редакционного расстояния занимают больше памяти, чем рекурсивные, так как память, используемая нерекурсивной реализацией растет как произведение длин строк, а память, используемая рекурсивной реализацией растет как сумма длин строк.

В связи с дополнительной операцией алгоритм Дамерау-Левенштейна работает в 1.3 раза медленнее, алгоритма Левенштейна. При этом среди реализаций алгоритмов подсчет расстояний Левенштейна быстрые результаты дает рекурсивная реализация с использованием матрицы за счет сохранения в ней промежуточных значений.

Цель, поставленная перед началом работы, была достигнута. В ходе лабораторной работы были решены следующие задачи:

- были изучены алгоритмы Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками;
- были разработаны изученные алгоритмы;
- был проведен сравнительный анализ реализованных алгоритмов;
- был подготовлен отчет о выполненной лабораторной работе.

Список литературы

- [1] Вычисление расстояния Левенштейна [Электронный ресурс]. Режим доступа: <https://foxford.ru/wiki/informatika/vychislenie-rasstoyaniya-levenshteyna> (дата обращения: 18.10.2021).
- [2] Welcome to Python [Электронный ресурс]. Режим доступа: <https://www.python.org> (дата обращения: 18.10.2021).
- [3] time — Time access and conversions [Электронный ресурс]. Режим доступа: <https://docs.python.org/3/library/time.html#functions> (дата обращения: 18.10.2021).
- [4] Ubuntu 20.04 LTS (Focal Fossa) Beta [Электронный ресурс]. Режим доступа: <http://old-releases.ubuntu.com/releases/20.04.1/> (дата обращения: 18.10.2021).
- [5] Мобильный процессор AMD Ryzen™ 3 3200U с графикой Radeon™ Vega 3 [Электронный ресурс]. Режим доступа: <https://www.amd.com/ru/products/apu/amd-ryzen-3-3200u> (дата обращения: 18.10.2021).