

НАПРАВЛЕНИЕ ПОДГОТОВКИ **09.03.04 Программная инженерия**

по лабораторной работе № 1

Название: Расстояния Левенштейна и Дамерау-Левенштейна

Дисциплина: Анализ алгоритмов

 (Подпись, дата)

Л.Л. Волкова

 (И.О. Фамилия)

Содержание

Введение	3
1 Аналитическая часть	4
1.1 Описание алгоритмов	5
1.1.1 Рекурсивный алгоритм Левенштейна	5
1.1.2 Матричный алгоритм Левенштейна	5
1.1.3 Рекурсивно-матричный алгоритм Левенштейна	6
2 Конструкторская часть	7
2.1 Схемы алгоритмов	8
2.1.1 Схема рекурсивной реализации алгоритма Левенштейна	8
2.1.2 Схема матричной реализации алгоритма Левенштейна	9
2.1.3 Схема рекурсивного матричного алгоритма Левенштейна	10
2.1.4 Схема алгоритма Дамерау-Левенштейна	11
3 Технологическая часть	12
3.1 Требования к программному обеспечению	12
3.2 Средства реализации	12
3.3 Листинг кода	12
3.3.1 Рекурсивный алгоритм Левенштейна	12
3.3.2 Матричный алгоритм Левенштейна	13
3.3.3 Рекурсивный матричный алгоритм Левенштейна	13
3.3.4 Алгоритм Дамерау-Левенштейна	14
3.4 Сравнительный анализ матричной и рекурсивной реализаций	15
3.4.1 Теоретический анализ затрачиваемой памяти	15
3.5 Описание тестирования	16
3.5.1 Интерфейс программы	16
3.5.2 Тесты	16
4 Исследовательская часть	17
4.1 Примеры работы	17
4.2 Замер времени	18
4.3 Сравнительный анализ на материале экспериментальных данных	18
Заключение	19

Введение

Редакционное расстояние или же расстояние Левенштейна - это минимальное кол-во редакторских операций, которое необходимо для превращения одной строки в другую.

Применяются данные алгоритмы в поисковых строках браузеров, а также в биоинформатике.

Задания данной лабораторной работы:

- 1) изучить алгоритмы Левенштейна и Дамерау-Левенштейна;
- 2) реализовать четыре алгоритма: матричный, рекурсивный, рекурсивный с заполнением матрицы и Дамерау-Левенштейна;
- 3) провести сравнительный анализ данных алгоритмов определения расстояния между строками;
- 4) предоставить экспериментальное подтверждение различий во временной эффективности алгоритмов при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализации на варьирующихся длинах строк;

1 Аналитическая часть

Задача по нахождению расстояния Левенштайна заключается в поиске минимального количества операций необходимых для преобразования одной строки в другую.

Для выполнения этой задачи используются следующие редакционные операции:

Вставка (I - Insert) - штраф 1

Удаление (D - Delete) - штраф 1

Замена (R - Replace) - штраф 1

Совпадение (M - Match) - штраф 0

В алгоритме Далмерау-Левенштайна допускается ещё одна редакционная операция:

Транспозиция(T) - штраф 1

Все операции обладают штрафом.

Штраф - это условная стоимость совершения данной операции. Цель алгоритмов подобрать набор операций, суммарный штраф которых будет минимален.

1.1 Описание алгоритмов

1.1.1 Рекурсивный алгоритм Левенштейна

Введём понятие $D(s1, s2)$ = минимальному количеству редакторских операций, с помощью которых строка $s1$ преобразуется в строку $s2$. Тогда рекурсивный алгоритм Левенштейна можно записать следующим образом:

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

1.1.2 Матричный алгоритм Левенштейна

Вводится матрица, размерностью $[Len(S_1) + 1 \times Len(S_2) + 1]$

Первая строки и столбец матрицы заполняются от 0 до $Len(S)$ (первые 3 пункта системы из предыдущего пункта).

$$A = \begin{pmatrix} & \emptyset & С & Т & О & Л & Б \\ \emptyset & 0 & 1 & 2 & 3 & 4 & 5 \\ Т & 1 & & & & & \\ Е & 2 & & & & & \\ Л & 3 & & & & & \\ О & 4 & & & & & \end{pmatrix}$$

Далее для нахождения ответа применяется последняя формула из системы, описанной в предыдущем пункте.

$$A = \begin{pmatrix} & \emptyset & С & Т & О & Л & Б \\ \emptyset & 0 & 1 & 2 & 3 & 4 & 5 \\ Т & 1 & 1 & 1 & 2 & 3 & 4 \\ Е & 2 & 2 & 2 & 2 & 3 & 4 \\ Л & 3 & 3 & 3 & 3 & 2 & 3 \\ О & 4 & 4 & 4 & 3 & 3 & \mathbf{3} \end{pmatrix}$$

Ответ в правом нижнем углу.

Чтобы определить, какая именно цепочка преобразований привела к ответу представим матрицу как карту высот: нужно спуститься на санках из клетки с ответом в левый верхний угол. В нашем случае:

I: ТЕЛО \rightarrow СТЕЛО

M: Т = Т

R: Е \rightarrow О

M: Л = Л

R: О \rightarrow Б

1.1.3 Рекурсивно-матричный алгоритм Левенштейна

Аналогичен алгоритму из предыдущего пункта с той лишь разницей, что матрица начинает заполнение "с конца". Вычисляем значение ячейки матрицы только в том случае, если значения там ещё нет (аналогично ∞ в алгоритме Дейкстры). Ответ всё так же в правом нижнем углу.

2 Конструкторская часть

Требования к вводу:

- 1) на вход подаются две строки;

Требования к программе::

- 1) Две пустые строки являются корректным вводом, который программа должна обработать.
- 2) одна и та же буква в разном регистре считается как разный символ.

2.1 Схемы алгоритмов

В данном разделе представлены схемы реализуемых алгоритмов.

2.1.1 Схема рекурсивной реализации алгоритма Левенштейна

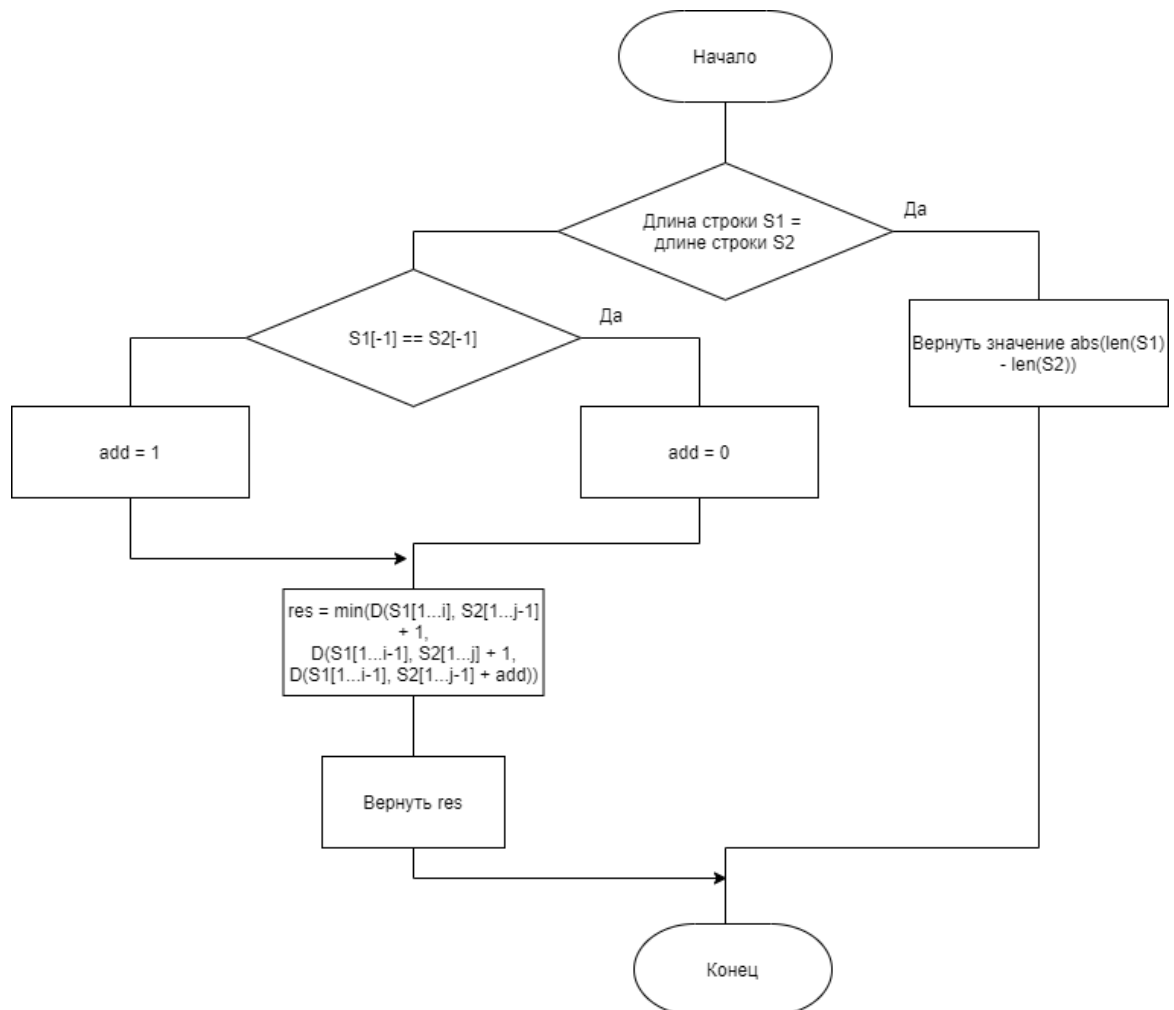


Рис. 1: Схема рекурсивного алгоритма Левенштейна

2.1.2 Схема матричной реализации алгоритма Левенштейна

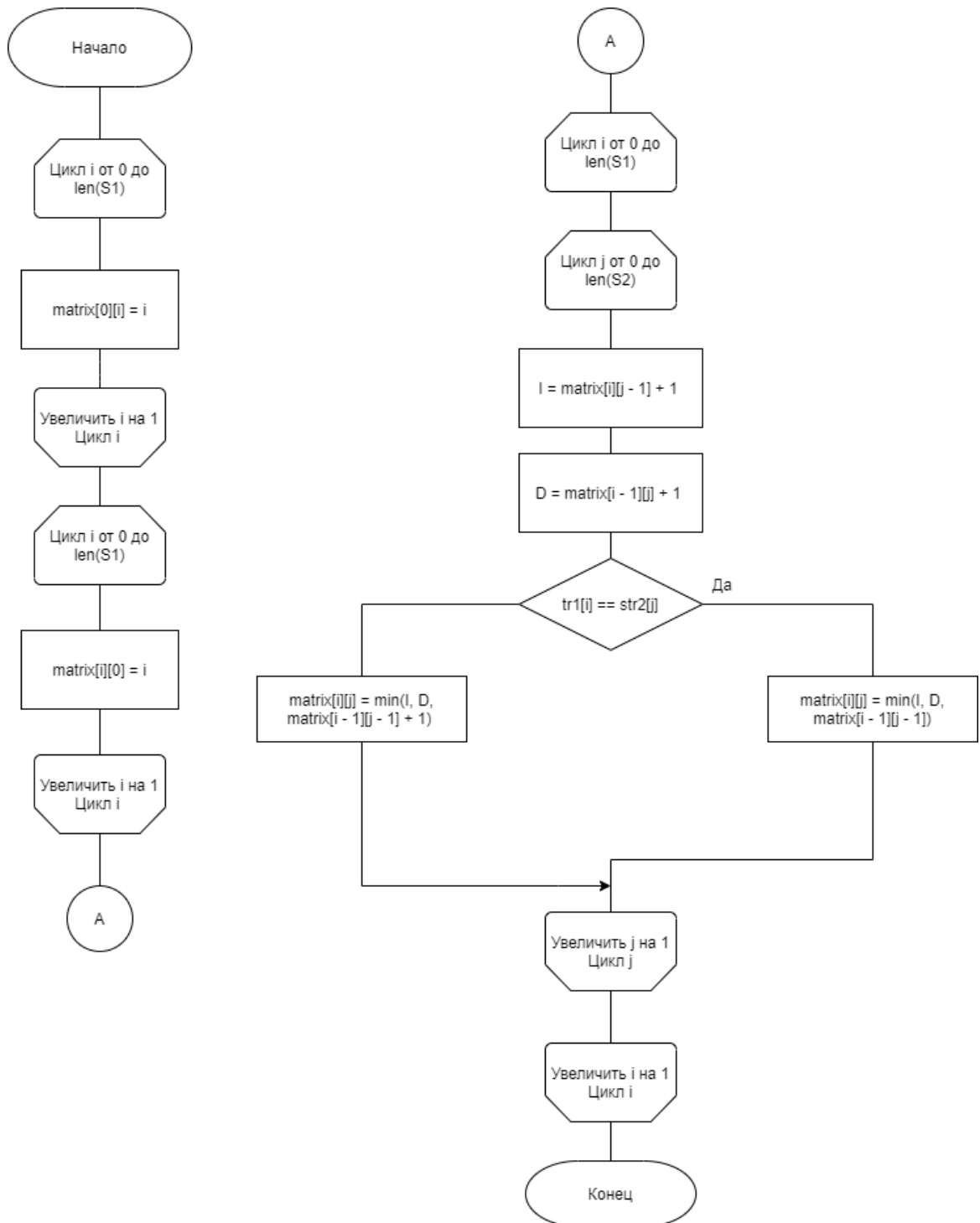


Рис. 2: Схема матричной реализации алгоритма Левенштейна

2.1.3 Схема рекурсивного матричного алгоритма Левенштейна

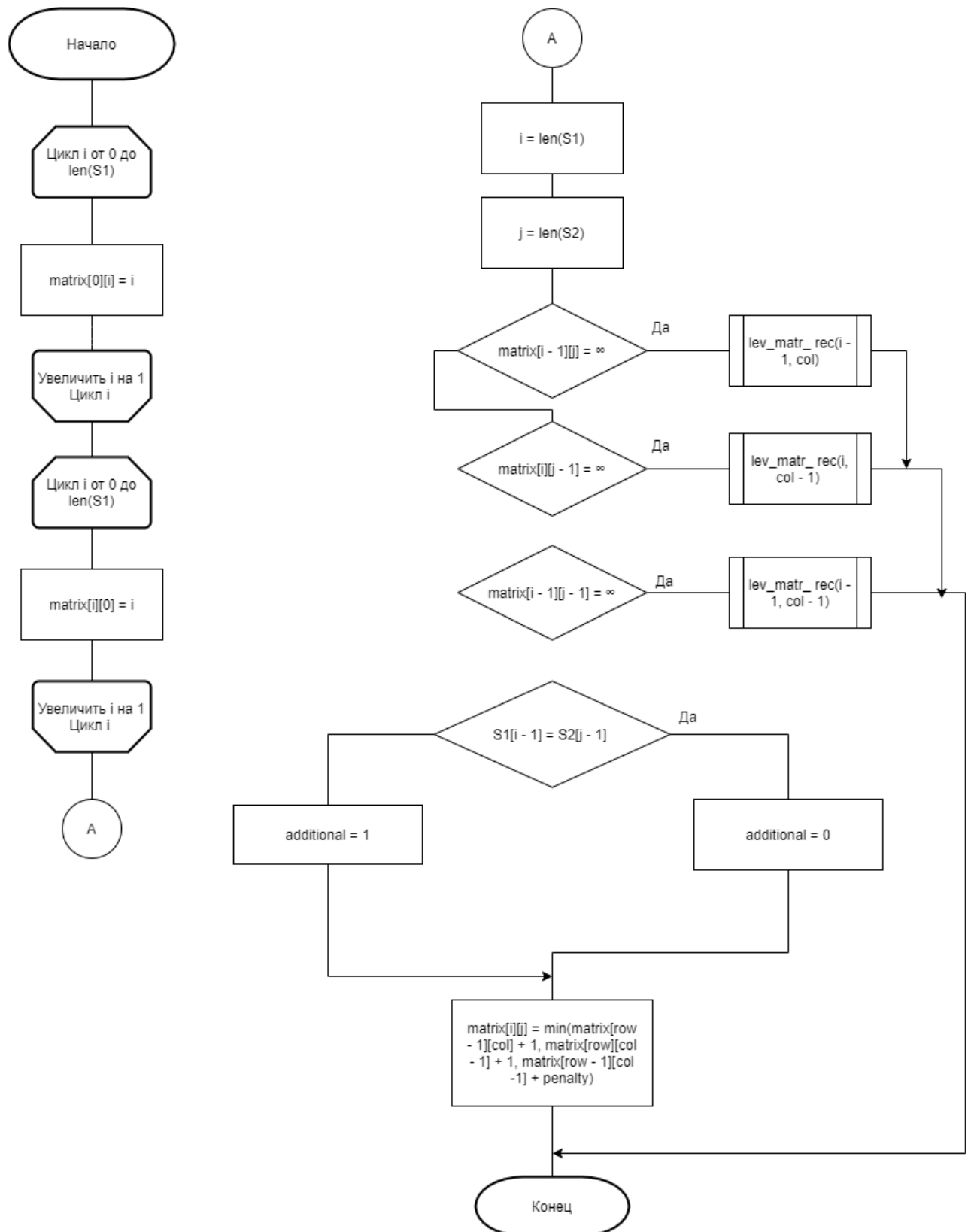


Рис. 3: Схема рекурсивной матричной реализации алгоритма Левенштейна

2.1.4 Схема алгоритма Дамерау-Левенштейна

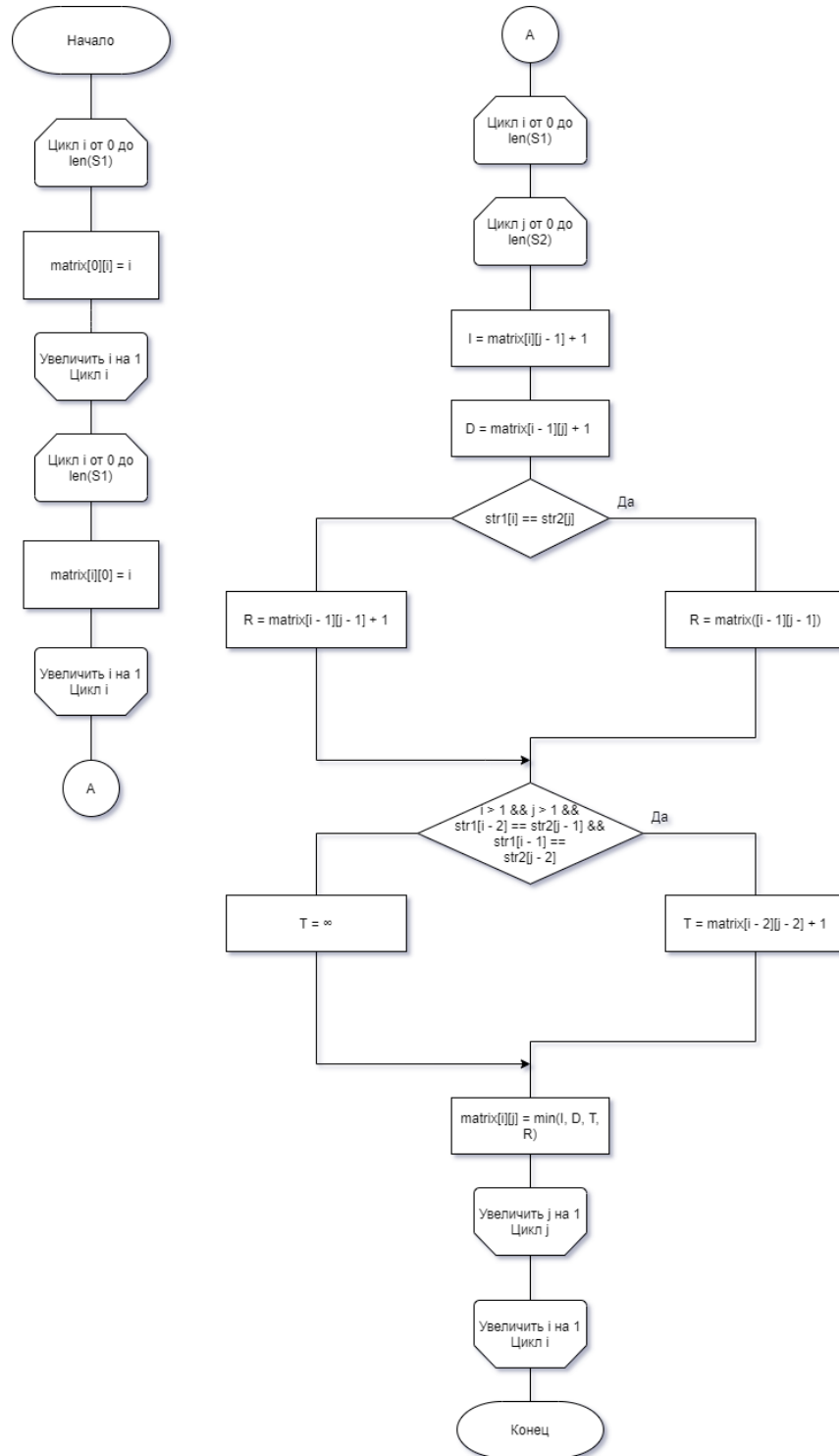


Рис. 4: Схема алгоритма Дамерау-Левенштейна

3 Технологическая часть

В данном разделе будет описана технологическая часть лабораторной работы: требования к ПО, листинг кода, сравнительный анализ всех алгоритмов.

3.1 Требования к программному обеспечению

Входные данные: две строки

Выходные данные: редакционное расстояние данных слов, а также матрица решения для матричных реализаций

Среда выполнения: Windows 10 x64

3.2 Средства реализации

Для выполнения данной лабораторной работы использовался ЯП Python 3.8.0

3.3 Листинг кода

В данном разделе представлен листинг кода разработанных алгоритмов.

3.3.1 Рекурсивный алгоритм Левенштейна

```
1 def levenshtein_recursive(s1, s2, i, j):
2     if min(i, j) == 0:
3         return max(i, j)
4     else:
5         m = 0 if s1[i-1] == s2[j-1] else 1
6         return min(levenshtein_recursive(s1, s2, i-1, j) + 1,
7                     levenshtein_recursive(s1, s2, i, j-1) + 1,
8                     levenshtein_recursive(s1, s2, i-1, j-1) + m)
```

3.3.2 Матричный алгоритм Левенштейна

```
1 def levenshtein_matrix(s1, s2):
2     n1 = len(s1)
3     n2 = len(s2)
4     matrix = [[0 for j in range(n2+1)] for i in range(n1+1)]
5
6     for i in range(n1+1):
7         matrix[i][0] = i
8
9     for j in range(n2+1):
10        matrix[0][j] = j
11
12    for i in range(1, n1+1):
13        for j in range(1, n2+1):
14            m = 0 if s1[i-1] == s2[j-1] else 1
15            matrix[i][j] = min(matrix[i-1][j] + 1,
16                               matrix[i][j-1] + 1,
17                               matrix[i-1][j-1] + m)
18
19    return matrix[-1][-1]
```

3.3.3 Рекурсивный матричный алгоритм Левенштейна

```
1 def levenshtein_recursive_matrix(s1, s2, i, j, matrix):
2     if min(i, j) == 0:
3         matrix[i][j] = max(i, j)
4     else:
5         if matrix[i][j] == -1:
6             m = 0 if s1[i-1] == s2[j-1] else 1
7             matrix[i][j] = min(levenshtein_recursive_matrix(s1, s2, i-1,
8                                                             levenshtein_recursive_matrix(s1, s2, i, j-1,
9                                                             levenshtein_recursive_matrix(s1, s2, i-1,
10
11    return matrix[i][j]
```

3.3.4 Алгоритм Дамерау-Левенштейна

```
1 def dameray_levenshtein(s1, s2):
2     n1 = len(s1)
3     n2 = len(s2)
4     matrix = [[0 for j in range(n2+1)] for i in range(n1+1)]
5
6     for i in range(n1+1):
7         matrix[i][0] = i
8
9     for j in range(n2+1):
10        matrix[0][j] = j
11
12    for i in range(1, n1+1):
13        for j in range(1, n2+1):
14            m = 0 if s1[i-1] == s2[j-1] else 1
15            if i > 1 and j > 1 and s1[i-2] == s2[j-1] and s1[i-1] == s2[j-2]:
16                matrix[i][j] = min(matrix[i-1][j] + 1,
17                                     matrix[i][j-1] + 1,
18                                     matrix[i-1][j-1] + m,
19                                     matrix[i-2][j-2] + 1)
20            else:
21                matrix[i][j] = min(matrix[i-1][j] + 1,
22                                     matrix[i][j-1] + 1,
23                                     matrix[i-1][j-1] + m)
24
25    return matrix[-1][-1]
```

3.4 Сравнительный анализ матричной и рекурсивной реализаций

Рекурсивная версия алгоритма работает значительно медленнее матричной реализации из-за многократного вызова функции. На каждый вызов необходимо производить соответствующие операции со стеком. Главным недостатком является повторное вычисление тех значений, которые были посчитаны на более ранних этапах рекурсии. В матричных же реализациях будет затрачена дополнительная память на хранение матриц и дополнительных переменных, однако скорость работы алгоритма будет значительно быстрее чем у рекурсивного.

3.4.1 Теоретический анализ затрачиваемой памяти

Рекурсивная реализация алгоритма Левенштейна. Для получения конечной оценки затрачиваемой памяти необходимо память, затрачиваемую на единичный вызов функции умножить на максимальную глубину рекурсии, то есть на $n + m$, где n и m - длины сравниваемых строк $s1$ и $s2$ соответственно.

1. ссылки на строки $s1, s2$: $(m + n) * \text{sizeof}(\text{reference})$,
2. длины строк: $2 * \text{sizeof}(\text{int})$,
3. дополнительная переменная внутри алгоритма: $\text{sizeof}(\text{int})$
4. адрес возврата

Матричная реализация алгоритма Левенштейна

1. строки: $\text{sizeof}(\text{str}) * (n + m)$
2. матрица: $\text{sizeof}(\text{int}) * (n + 1) * (m + 1)$
3. дополнительная переменная внутри алгоритма: $\text{sizeof}(\text{int})$

Рекурсивный матричный алгоритм Левенштейна. Аналогично обычному рекурсивному алгоритму для получения конечной оценки затрачиваемой памяти необходимо память, затрачиваемую на каждом рекурсивном вызове умножить на максимальную глубину рекурсии.

1. строки: $\text{sizeof}(\text{str}) * (n + m)$
2. матрица: $\text{sizeof}(\text{int}) * (n + 1) * (m + 1)$

При каждой необходимости предварительного подсчёта значения (рек. вызова)

1. передача строки и столбца: $2 * \text{sizeof}(\text{int})$
2. дополнительная переменная: $\text{sizeof}(\text{int})$
3. адрес возврата

Матричная реализация алгоритма Дамера-Левенштейна

1. строки: $\text{sizeof}(\text{str}) * (n + m)$
2. матрица: $\text{sizeof}(\text{int}) * (n + 1) * (m + 1)$
3. дополнительная переменная внутри алгоритма: $\text{sizeof}(\text{int})$

3.5 Описание тестирования

3.5.1 Интерфейс программы

При запуске программы пользователя встречает меню выбора реализаций алгоритма:

```
Меню:
1. Левенштейн с матрицей
2. Левенштейн с рекурсией
3. Левенштейн рекурсивный с матрицей
4. Дамерау-Левенштейн
5. Анализ времени
0. Выход
```

Рис.5: Меню программы

После выбора необходимой реализации пользователю предлагают ввести строки s1 и s2. После ввода программа выдаёт результат:

```
Меню:
1. Левенштейн с матрицей
2. Левенштейн с рекурсией
3. Левенштейн рекурсивный с матрицей
4. Дамерау-Левенштейн
5. Анализ времени
0. Выход

1
Введите первую строку: stolb
Введите вторую строку: telo
Расстояние между строками: 3
```

Рис.6: Ввод строк и результат работы программы

3.5.2 Тесты

Тестирование было организовано с помощью библиотеки unittest.

4 Исследовательская часть

В данной части отчета приведены примеры работы программ, а также анализ алгоритмов на основе экспериментальных данных.

4.1 Примеры работы

Пустые строки:

```
1      string_1 = ""
2      string_2 = ""
3      Recursive Levenshtein: 0
4      Table Levenshtein: 0
5      Recursive-Table Levenshtein: 0
6      Damerau-Levenshtein: 0
```

Равенство строк:

```
1      string_1 = "abc"
2      string_2 = "abc"
3      Recursive Levenshtein: 0
4      Table Levenshtein: 0
5      Recursive-Table Levenshtein: 0
6      Damerau-Levenshtein: 0
```

Удаление:

```
1      string_1 = "aac"
2      target = "aa"
3      Recursive Levenshtein: 1
4      Table Levenshtein: 1
5      Recursive-Table Levenshtein: 1
6      Damerau-Levenshtein: 1
```

Замена:

```
1      string_1 = "abc"
2      string_2 = "abd"
3      Recursive Levenshtein: 1
4      Table Levenshtein: 1
5      Recursive-Table Levenshtein: 1
6      Damerau-Levenshtein: 1
```

Вставка:

```
1      string_1 = "ab"
2      string_2 = "abc"
3      Recursive Levenshtein: 1
4      Table Levenshtein: 1
5      Recursive-Table Levenshtein: 1
6      Damerau-Levenshtein: 1
```

Перестановка:

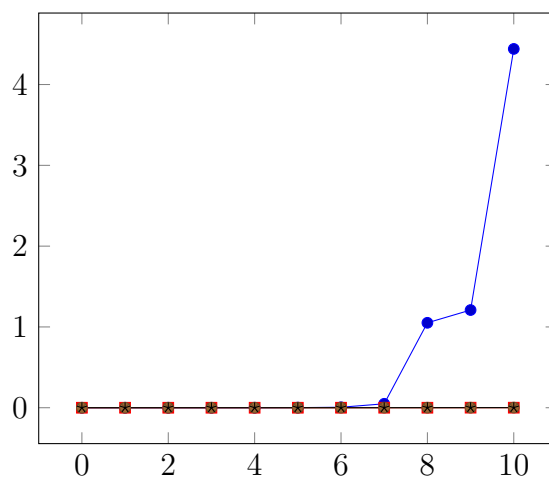
```
1      string_1 = "abc"
2      string_2 = "acb"
3      Recursive Levenshtein: 2
4      Table Levenshtein: 2
5      Recursive-Table Levenshtein: 2
6      Damerau-Levenshtein: 1
```

4.2 Замер времени

Были взяты строки 4 размерностей. Для каждой размерности было проведено 100 вызовов функции. После чего получившееся время было поделено на 100. Таким образом было получено аппроксимированное значение времени выполнения функции

Результаты замеров процессорного времени:

Длина строки	Lev(M)	Lev(R)	Lev(RM)	DamLev
3	2.0103	3.9951	2.0020	2.0096
5	1.9924	106.0962	4.0040	3.0026
7	4.0032	3012.1695	8.0068	4.9967
10	7.9977	500312.1714	16.0050	8.0165



Легенда:

Синий цвет - Рекурсивная реализация

Коричневый цвет - Рекурсивная матричная реализация

Чёрный цвет - Алгоритм Дамерау-Левенштейна

Красный цвет - Обычная матричная реализация

4.3 Сравнительный анализ на материале экспериментальных данных

Теоретические расчёты подтвердились результатами, полученными на практике: рекурсивный алгоритм из-за многократного вызова функции и пересчёта уже известных значений выполняется очень долго, рекурсивная матричная реализация выполняется быстрее, но всё равно из-за затрат времени на вызов самой себя уступает по времени обычной матричной реализации. Алгоритм Дамерау-Левенштейна незначительно уступает обычной матричной реализации ввиду дополнительной проверки.

Заключение

В ходе работы были изучены алгоритмы нахождения расстояний Левенштейна и Дамерау–Левенштейна, реализованы указанные алгоритмы, а также произведён сравнительный анализ алгоритмов. Было установлено, что обычный рекурсивный алгоритм занимает меньше памяти по сравнению с матричными реализациями, однако сильно уступает им по времени.