



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 Дamerau-Левенштейна

Студент Прохорова Л. А.

Группа ИУ7-53Б

Оценка (баллы) _____

Преподаватели Волкова Л.Л., Строганов Ю.В.

Москва — 2020 г.

Введение

Редакционные расстояния позволяют понять, насколько близки две строки. Поисковая система все равно находит то, что надо, даже если вы совершили пару опечаток в запросе, потому что для неправильно введенного слова в словаре ищется наиболее похожее на него слово. Для определения сходства между словами или, в общем случае, строками, используется расстояние по Левенштейну, или редакционное расстояние: для пары строк рассчитывается, сколько действий по редактированию одной строки нужно совершить, чтобы получить другую строку.

Алгоритмы нечеткого поиска (поиска с указанием некоторой доли сходства) являются основой систем проверки орфографии и полноценных поисковых систем вроде Google или Yandex. Например, такие алгоритмы используются для функций наподобие «Возможно вы имели в виду ...» в тех же поисковых системах, сравнения текстовых файлов утилитой diff и ей подобными, в биоинформатике для сравнения генов, хромосом и белков.

1 Аналитическая часть

1.1 Детализация задачи

В данной работе я рассмотрю алгоритмы:

- поиска расстояния Левенштейна;
- поиска расстояния Дamerau-Левенштейна.

Целью данной лабораторной работы является изучить и реализовать алгоритмы поиска расстояния Левенштейна и Дamerau-Левенштейна и сравнить разные варианты их реализации.

Задачами данной лабораторной являются:

- изучение алгоритмов поиска расстояния Левенштейна и Дamerau-Левенштейна;
- Реализация алгоритма поиска расстояния Левенштейна с заполнением таблицы по формуле, рекурсивным заполнением таблицы, с использованием рекурсии;
- реализация алгоритма поиска расстояния Дamerau-Левенштейна без использования рекурсии.
- сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
- 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;

- описание и обоснование полученных результатов в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

1.2 Расстояние Левенштейна

Расстояние Левенштейна [1] - это минимальное необходимое количество редакторских операций (вставки, удаления, замены) для преобразования одной строки в другую.

Возможные действия и их обозначения:

- D (англ. delete) — удалить,
- I (англ. insert) — вставить,
- R (replace) — заменить,
- M(match) - совпадение.

Штраф за совпадение 0, за остальные действия 1.

Расстояние Левенштейна между двумя строками a и b может быть вычислено по формуле 1, где $|a|$ означает длину строки a ; $a[i]$ — i -ый символ строки a , функция $D(i, j)$, равная расстоянию между подстроками длины

i и j , определена как:

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

В формуле 1 функция m выражается как:

$$m(a, b) = \begin{cases} 0 & \text{если } a = b, \\ 1 & \text{иначе} \end{cases}. \quad (2)$$

Рекурсивный алгоритм поиска расстояния Левенштейна

Рекурсивный алгоритм реализует формулу 1.

Для перевода из строки a в строку b требуется выполнить последовательно некоторое количество операций (удаление, вставка, замена) в некоторой последовательности. Последовательность проведения любых двух операций можно поменять, порядок проведения операций не имеет никакого значения. Полагая, что a', b' — строки a и b без последнего символа соответственно, цена преобразования из строки a в строку b может быть выражена как:

1. Сумма цены преобразования строки a в b и цены проведения операции удаления, которая необходима для преобразования a' в a ;
2. Сумма цены преобразования строки a в b и цены проведения операции вставки, которая необходима для преобразования b' в b ;

3. Сумма цены преобразования из a' в b' и операции замены, предполагая, что a и b оканчиваются разные символы;
4. Цена преобразования из a' в b' , предполагая, что a и b оканчиваются на один и тот же символ.

Минимальной ценой преобразования будет минимальное значение приведенных вариантов.

Матричный алгоритм поиска расстояния Левенштейна

Прямая реализация формулы 1 может быть малоэффективна по времени исполнения при больших i, j , т. к. множество промежуточных значений $D(i, j)$ вычисляются заново множество раз подряд. Для оптимизации нахождения расстояния Левенштейна можно использовать матрицу в целях хранения соответствующих промежуточных значений. В таком случае алгоритм представляет собой построчное заполнение матрицы $A_{|a|,|b|}$ значениями $D(i, j)$.

На рисунке 1 изображена матрица, которой можно представить весь процесс:

		А	Р	Е	С	Т	А	Н	Т
	0	1	2	3	4	5	6	7	8
Д	1	1	2	3	4	5	6	7	8
А	2	1	2	3	4	5	5	6	7
Г	3	2	2	3	4	5	6	6	7
Е	4	3	3	2	3	4	5	6	7
С	5	4	4	3	2	3	4	5	6
Т	6	5	5	4	3	2	3	4	5
А	7	6	6	5	4	3	2	3	4
Н	8	7	7	6	5	4	3	2	3

Рисунок 1 – Матрица поиска расстояния Левенштейна.

На рисунке 2 представлена работа с ячейками матрицы. Если посмотреть на процесс работы алгоритма, несложно заметить, что на каждом шаге используются только две последние строки матрицы, следовательно, потребление памяти можно уменьшить до $O(\min(m, n))$ [2].

X_z	X_y
X_v	$D_{i,j}$

Рисунок 2 – Используемые за одну итерацию ячейки матрицы при поиске расстояния Левенштейна.

Рекурсивный алгоритм поиска расстояния Левенштейна с заполнением матрицы

Рекурсивный алгоритм заполнения можно оптимизировать по времени выполнения с использованием матричного алгоритма. Суть данного метода заключается в параллельном заполнении матрицы при выполнении рекурсии. В случае, если рекурсивный алгоритм выполняет прогон для данных, которые еще не были обработаны, результат нахождения расстояния заносится в матрицу. В случае, если обработанные ранее данные встречаются снова, для них расстояние не находится и алгоритм переходит к следующему шагу.

1.3 Расстояние Дамерау–Левенштейна

Эта вариация вносит в определение расстояния Левенштейна еще одно правило — транспозиция (перестановка) двух соседних букв также учитывается как одна операция, наряду со вставками, удалениями и заменами [2]. Используемые на каждой итерации ячейки изображены на рисунке 3.

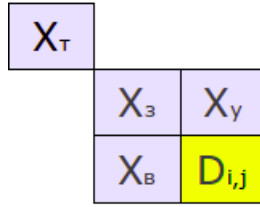


Рисунок 3 – Используемые за одну итерацию ячейки матрицы при поиске расстояния Дамерау–Левенштейна.

Расстояние Дамерау — Левенштейна может быть найдено по формуле 3, которая задана как

$$d_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ \\ \quad d_{a,b}(i, j - 1) + 1, \\ \quad d_{a,b}(i - 1, j) + 1, & \text{иначе} \\ \quad d_{a,b}(i - 1, j - 1) + m(a[i], b[j]), \\ \quad \left[\begin{array}{ll} d_{a,b}(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a[i] = b[j - 1]; \\ & b[j] = a[i - 1] \\ & \infty, & \text{иначе} \end{array} \right. \\ \} \end{cases}, \quad (3)$$

Вывод

Формулы Левенштейна и Дамерау — Левенштейна для расчета расстояния между строками задаются рекурсивно, а следовательно, алгоритмы могут быть реализованы рекурсивно или итерационно(с помощью матрицы).

2 Конструкторская часть

2.1 Требования к вводу

В разрабатываемом ПО предъявляются следующие требования ко вводу.

1. На вход подаются две строки.
2. Строчные и прописные буквы считаются разными.

2.2 Требования к выводу

Программа выводит расстояние и матрицу(если она использовалась).

2.3 Требования к программе

Две пустые строки - корректный ввод, программа не должна аварийно завершаться.

2.4 Схемы алгоритмов

На рисунках 4, 5, 6 показаны схемы алгоритмов поиска расстояния Левенштейна. На рисунке 7 показан алгоритм поиска расстояния Дameraу–Левенштейна без рекурсии.

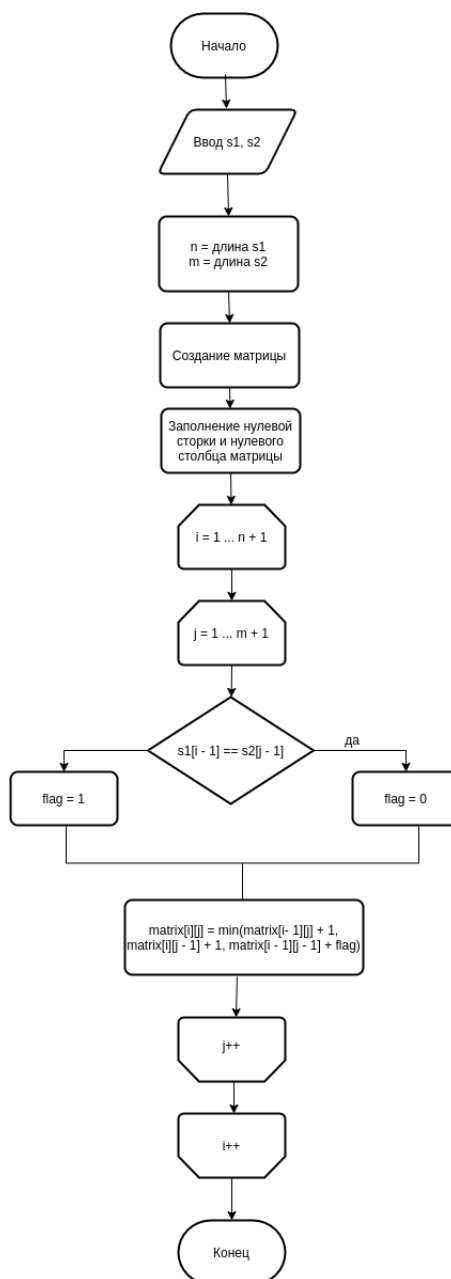


Рисунок 4 – Схема алгоритма поиска расстояния Левенштейна с помощью матрицы(без рекурсии).

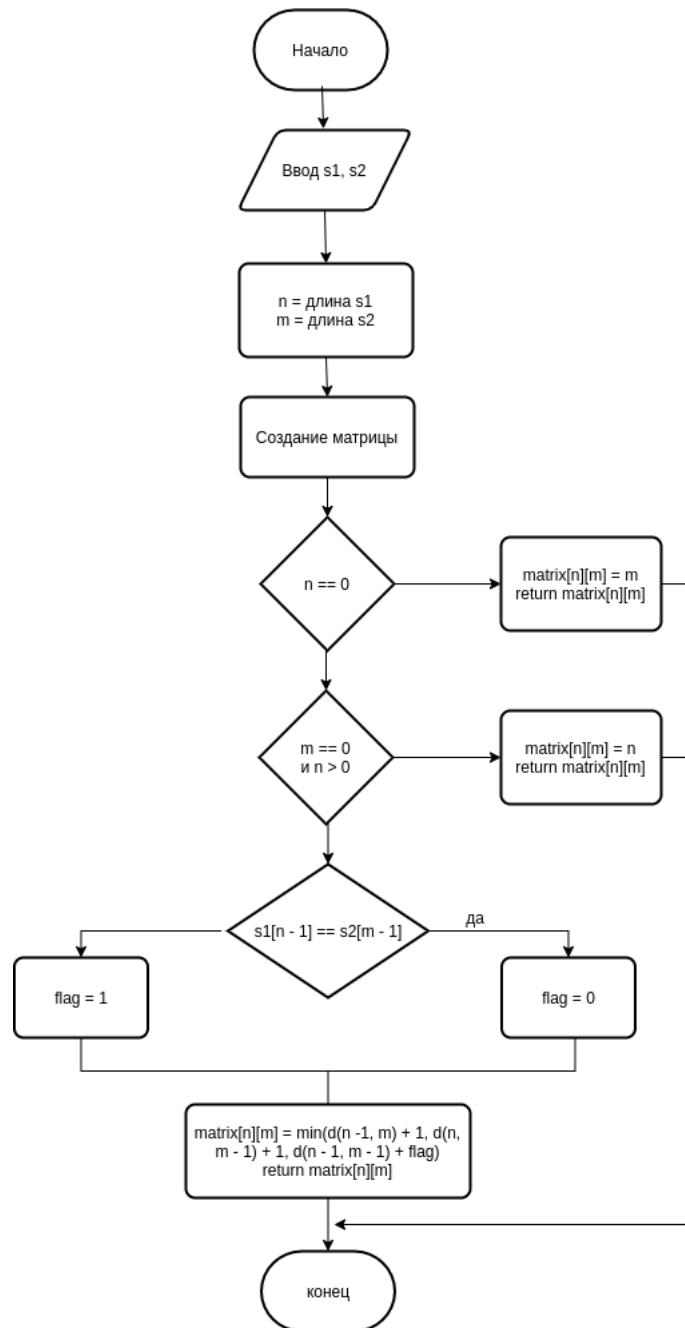


Рисунок 5 – Схема алгоритма поиска расстояния Левенштейна с помощью матрицы(с использованием рекурсии).

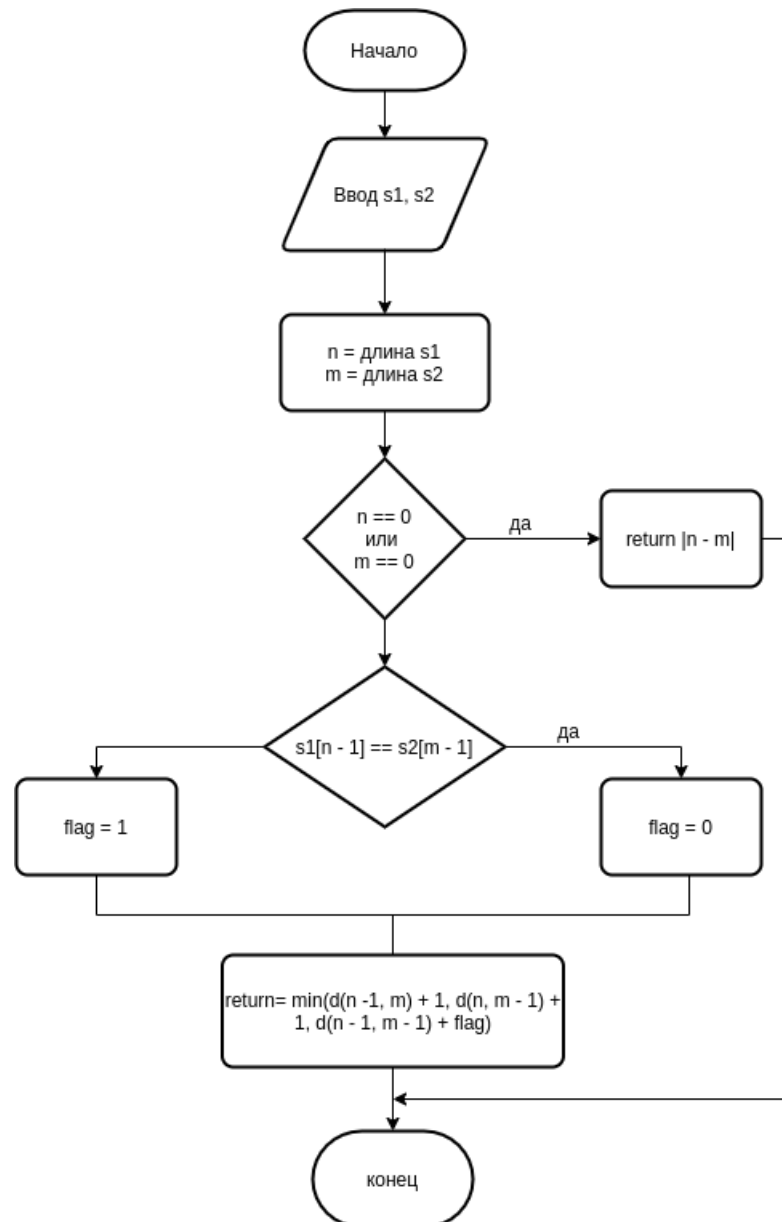


Рисунок 6 – Схема алгоритма поиска расстояния Левенштейна с использованием рекурсии.

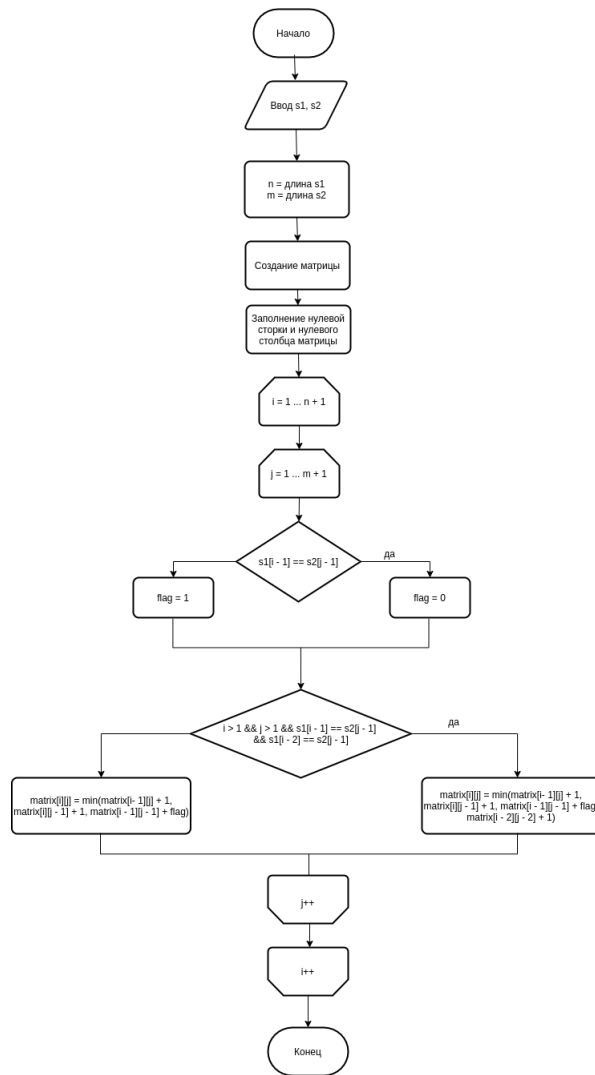


Рисунок 7 – Схема алгоритма поиска расстояния Дameraу–Левенштейна с помощью матрицы(без рекурсии).

3 Технологическая часть

3.1 Выбор языка программирования

Я выбрала Python языком программирования, потому как он достаточно удобен и гибок.

Время работы алгоритмов было замерено с помощью функции `time()` из библиотеки `time`.

3.2 Сведения о модулях программы

Программа состоит из следующих модулей:

- `main.py` - главный файл программы, в котором располагаются алгоритмы и меню
- `test.py` - файл с замерами времени

В листингах 1, 2, 3, 4 представлен код используемых алгоритмов.

Листинг 1 – Функция нахождения расстояния Левенштейна с использованием рекурсии.

```
0 def levinstein_recursive(s1, s2):
1     n = len(s1)
2     m = len(s2)
3
4     if n == 0 or m == 0:
5         return abs(n - m)
6
7     flag = 0
8     if s1[-1] != s2[-1]:  если# последние символы не равны
9         flag = 1
10
11     return min(levinstein_recursive(s1[:-1], s2) + 1,
12                levinstein_recursive(s1, s2[:-1]) + 1,
13                levinstein_recursive(s1[:-1], s2[:-1]) + flag)
```

Листинг 2 – Функция нахождения расстояния Левенштейна с использованием матрицы(не рекурсивно)

```
0 def levinstein_matrix(s1, s2):
1     n = len(s1)
2     m = len(s2)
3
4     matrix = create_matrix(n + 1, m + 1)
5
6     for i in range(1, n + 1):
7         for j in range(1, m + 1):
8             add, delete, change = matrix[i - 1][j] + 1, matrix[i][j - 1] +
9             1, matrix[i - 1][j - 1]
10            if s2[j - 1] != s1[i - 1]:
11                change += 1
12            matrix[i][j] = min(add, delete, change)
13
14    print("\Матрица расстояния Левенштейна")
15    print_matrix(s1, s2, matrix)
16    print("Расстояние Левенштейна = ", matrix[n][m])
17
18    return matrix[n][m]
```

Листинг 3 – Функция нахождения расстояния Левенштейна с использованием матрицы(рекурсивно)

```
0 def levinstein_recursive_matrix(s1, s2):
1     n, m = len(s1), len(s2)
2
3     def recursive(s1, s2, n, m, matrix):
4         if (matrix[n][m] != -1):
5             return matrix[n][m]
6
7         if (n == 0):
8             matrix[n][m] = m
9             return matrix[n][m]
10
11        if (m == 0 and n > 0):
12            matrix[n][m] = n
13            return matrix[n][m]
```

```

15     delete = recursive(s1, s2, n - 1, m, matrix) + 1 # удаление
16     add = recursive(s1, s2, n, m - 1, matrix) + 1 # вставка
17     flag = 0
18     if (s1[n - 1] != s2[m - 1]):
19         flag = 1
20     change = recursive(s1, s2, n - 1, m - 1, matrix) + flag # замена
21     matrix[n][m] = min(delete, change, add)
22
23     return matrix[n][m]
24
25 matrix = create_matrix(n + 1, m + 1)
26
27 for i in range(n + 1):
28     for j in range(m + 1):
29         matrix[i][j] = -1
30
31 recursive(s1, s2, n, m, matrix)
32 print("\Матрица_заполненная_рекурсивно:")
33 print_matrix(s1, s2, matrix)
34 print("Расстояние_Левенштейна=", matrix[n][m])
35
36 return matrix[n][m]

```

Листинг 4 – Функция нахождения расстояния Дameraу–Левенштейна с использованием матрицы(не рекурсивно)

```

0 def damerau_levinstein_matrix(s1, s2):
1     n = len(s1)
2     m = len(s2)
3
4     matrix = create_matrix(n + 1, m + 1)
5
6     for i in range(1, n + 1):
7         for j in range(1, m + 1):
8             add, delete, change = matrix[i - 1][j] + 1, matrix[i][j - 1] +
9             1, matrix[i - 1][j - 1]
10            if s2[j - 1] != s1[i - 1]:
11                change += 1
12            matrix[i][j] = min(add, delete, change)
13            if ((i > 1 and j > 1) and s1[i - 1] == s2[j - 2] and s1[i - 2]
14            == s2[j - 1]):

```



```

13         matrix[i][j] = min(matrix[i][j], matrix[i - 2][j - 2] + 1)
14     print("\Матрица расстояния ДameraуЛевенштейна - ")
15     print_matrix(s1, s2, matrix)
16     print("Расстояние ДameraуЛевенштейна -- = ", matrix[n][m])
17
18     return matrix[n][m]

```

3.3 Тесты для проверки корректности программы

В таблице 1 представлены функциональные тесты для проверки работы программы. Все тесты пройдены успешно.

Входные данные	Расстояние Левенштейна	Расстояние Дameraу– Левенштейна
Пустая строка, пустая строка	0	0
Пустая строка, "недвижимость"	12	12
"недвижимость" "недвижимость"	0	0
"недвижимость" "движимость"	2	2
"недвижимость" "двигатель"	8	8
"двигатель" "дивгатель"	2	1
"сад" "огурец"	6	6

Таблица 1 – Функциональные тесты

3.4 Сравнительный анализ памяти, затрачиваемой алгоритмами

Алгоритмы Левенштейна и Дameraу — Левенштейна не отличаются друг от друга с точки зрения использования памяти. Рассмотрим разницу между рекурсивной и матричной реализациями.

Максимальная глубина стека вызовов при рекурсивной реализации равна сумме длин входящих строк, соответственно, максимальный расход памяти (4)

$$(\mathcal{C}(S_1) + \mathcal{C}(S_2)) \cdot (2 \cdot \mathcal{C}(\text{string}) + 3 \cdot \mathcal{C}(\text{int})), \quad (4)$$

где \mathcal{C} — оператор вычисления размера, S_1, S_2 — строки, int — целочисленный тип, string — строковый тип.

Использование памяти при итеративной реализации теоретически равно

$$(\mathcal{C}(S_1) + 1) \cdot (\mathcal{C}(S_2) + 1) \cdot \mathcal{C}(\text{int}) + 7 \cdot \mathcal{C}(\text{int}) + 2 \cdot \mathcal{C}(\text{string}). \quad (5)$$

Вывод

Были разработаны и протестированы спроектированные алгоритмы: вычисления расстояния Левенштейна рекурсивно, с заполнением матрицы и рекурсивно с заполнением матрицы, а также вычисления расстояния Дамерау – Левенштейна с заполнением матрицы.

4 Исследовательская часть

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялось тестирование:

- Операционная система Ubuntu 18.04 64-bit.
- Память 8 GiB Процессор Intel(R) Core(TM) i7-7700HQ CPU @ 2.80GHz

4.2 Демонстрация работы программы

На рисунке 8 представлен результат работы программы.

```
Введите первую строку: стол
Введите вторую строку: стул

Матрица расстояния Левенштейна
0 0 с т у л
0 0 1 2 3 4
с 1 0 1 2 3
т 2 1 0 1 2
о 3 2 1 1 2
л 4 3 2 2 1
Расстояние Левенштейна = 1

Матрица заполненная рекурсивно:
0 0 с т у л
0 0 1 2 3 4
с 1 0 1 2 3
т 2 1 0 1 2
о 3 2 1 1 2
л 4 3 2 2 1
Расстояние Левенштейна = 1

Расстояние Левенштейна, полученное с использованием рекурсии: 1

Матрица расстояния Дамерау-Левенштейна
0 0 с т у л
0 0 1 2 3 4
с 1 0 1 2 3
т 2 1 0 1 2
о 3 2 1 1 2
л 4 3 2 2 1
Расстояние Дамерау-Левенштейна = 1
```

Рисунок 8 – Результат работы программы.

4.3 Время работы алгоритмов.

Алгоритмы тестировались при помощи функции `process_time()` из библиотеки `time` языка Python. `time.process_time()` всегда возвращает значение времени типа `float` в секундах. Возвращает значение (в долях секунды) суммы системного и пользовательского процессорного времени текущего процесса. Не включает время, прошедшее во время сна. Контрольная точка возвращаемого значения не определена, поэтому допустима только разница между результатами последовательных вызовов. Замеры времени для каждой длины слов проводились 100 раз. В качестве результата взято среднее время работы алгоритма на данной длине слова.

Результат замера времени представлен в таблице ??(время в мкс).

4.4 Функциональные тесты

Таблица 1 – Замеры времени

Длина	Л.(матр.)	Л.(рек с матр.)	Л.(рек)	Д.-Л.(матр.)
0	6.55026	8.18477	2.18068	6.66018
1	6.44065	9.33003	4.08478	6.34225
2	21.11541	31.87889	27.97856	23.62137
3	19.70278	31.66337	78.40595	22.21815
4	25.83701	46.96368	375.43730	33.89747
5	36.70361	65.01132	2004.88116	49.30683
6	46.13962	86.76184	10236.31448	66.81548
7	60.47135	113.15406	58054.26190	88.57176
8	80.86100	145.18792	306364.39168	108.38893
9	89.60874	169.33969	1669072.17389	127.26931

Вывод

Рекурсивный алгоритм Левенштейна работает дольше итеративных реализаций, время его работы увеличивается в геометрической прогрессии. При увеличении длины строк становится очевидна выигрышность по времени матричного варианта. Уже при длине в 7 символов матричная реализация в 950 раз быстрее. Рекурсивный алгоритм с заполнением матрицы превосходит простой рекурсивный на аналогичных данных в 9876 раз. Алгоритм Дамерау — Левенштейна по времени выполнения сопоставим с алгоритмом Левенштейна. В нём добавлены дополнительные проверки, и по сути он является алгоритмом другого смыслового уровня.

По расходу памяти итеративные алгоритмы проигрывают рекурсивному: максимальный размер используемой памяти в них растёт как произведение длин строк, в то время как у рекурсивного алгоритма — как сумма длин строк.

Заключение

В ходе выполнения работы решены следующие задачи.

- Были изучены алгоритмы поиска расстояни Левенштейна и Дameraу–Левенштейна.
- Реализованны алгоритмы поиска расстояния Левенштейна с заполнением таблицы по формуле, рекурсивным заполнением таблицы, с использованием рекурсии.
- Реализован алгоритм поиска расстояния Дameraу–Левенштейна без использования рекурсии.
- Проведён сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти).
- Проведено 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.
- Были описаны и обоснованы полученные результаты в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

Экспериментально было подтверждено различие 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.

В результате исследований я пришла к выводу, что матричная реализация данных алгоритмов заметно выигрывает по времени при росте длины строк, но проигрывает по количеству затрачиваемой памяти.

Список литературы

- [1] В. И. Левенштейн. Двоичные коды с исправлением выпадений, вставок и замещений символов. Доклады Академий Наук СССР, 1965. 163.4:845-848.
- [2] Нечёткий поиск в тексте и словаре [Электронный ресурс]. Режим доступа:
<https://habr.com/ru/post/114997/> (Дата обращения: 12.09.20)