



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"

Студент _____ Сукочева Алис
Группа _____ ИУ7-53Б
Название предприятия _____ МГТУ им. Н. Э. Баумана, каф. ИУ7
Тема _____ Расстояние Левенштейна

Студент:	_____	Сукочева А.
	подпись, дата	Фамилия, И.О.
Преподаватель:	_____	Волкова Л.Л.
	подпись, дата	Фамилия, И. О.
Преподаватель:	_____	Строганов Ю.В.
	подпись, дата	Фамилия, И. О.

Содержание

Введение	3
1 Аналитический раздел	4
1.1 Некоторые теоретические сведения	4
1.2 Расстояние Левенштейна	4
1.3 Расстояние Дамерау-Левенштейна	4
1.4 Вывод	5
2 Конструкторский раздел	6
2.1 Разработка алгоритмов	6
2.2 Вывод	6
3 Технологический раздел	8
3.1 Выбор ЯП	8
3.2 Требования к программному обеспечению	8
3.3 Сведения о модулях программы	8
3.4 Тестирование	11
3.5 Вывод	11
4 Экспериментальная часть	13
4.1 Временные характеристики	13
4.2 Характеристики по памяти	14
4.3 Сравнительный анализ алгоритмов	14
4.4 Вывод	15
Заключение	16
Список использованных источников	17

Введение

В данной лабораторной работе будет рассмотрено расстояние Левенштейна. Данное расстояние показывает минимальное количество редакторских операций (вставки, замены и удаления), которые необходимы для перевода одной строки в другую. Это расстояние помогает определить схожесть двух строк.

Упомянутое расстояние используется в задачах автозамены. В частности, учитываются ошибки, которые человек может допускать при наборе текста. Кроме упомянутых трех ошибок (вставка лишнего символа, пропуск символа, замена одного символа другим), его пальцы могут нажимать на нужные клавиши не в том порядке. С этой проблемой поможет справиться расстояние Дameraу-Левенштейна. Данное расстояние задействует еще одну редакторскую операцию – транспозицию, или перестановку.

Практическое применение расстояние Левенштейна:

- сравнение введенной строки со словарными словами в поисковой системе, такой как 'yandex' или 'google';
- помогает найти разницу двух ДНК, дать оценку мутации.

Целью данной работы является разбор и реализация алгоритма Дameraу-Левенштейна и Левенштейна.

В рамках выполнения работы необходимо решить следующие задачи.

- а) Изучить расстояния Дameraу-Левенштейна и Левенштейна.
- б) Реализовать алгоритмы поиска расстояний, в частности, нерекурсивный (матричный) и рекурсивные алгоритмы поиска расстояния Левенштейна и нерекурсивный алгоритм поиска расстояния Дameraу-Левенштейна.
- в) Подсчет времени поиска расстояния.
- г) Сравнить временные характеристики, а также затраченную память.

1 Аналитический раздел

1.1 Некоторые теоретические сведения

При преобразовании одного слова в другое можно использовать следующие операции:

- а) D (*от англ. delete*) - удаление.
- б) I (*от англ. insert*) - вставка.
- в) R (*от англ. replace*) - замена.

Будем считать стоимость каждой вышеизложенной операции - 1.

Введем понятие совпадения - M (*от англ. match*). Его стоимость будет равна 0.

1.2 Расстояние Левенштейна

Имеем две строки S_1 и S_2 , длиной M и N соответственно. Расстояние Левенштейна рассчитывается по рекуррентной формуле (1.1).

$$D(S_1[1...i], S_2[1...j]) = \begin{cases} j, \text{ если } i == 0 \\ i, \text{ если } j == 0 \\ \min(\\ D(S_1[1...i], S_2[1...j-1]) + 1, \\ D(S_1[1...i-1], S_2[1...j]) + 1, \quad j>0, i>0 \\ D(S_1[1...i-1], S_2[1...j-1]) + \\ \left[\begin{array}{l} 0, \text{ если } S_1[i] == S_2[j] \\ 1, \text{ иначе} \end{array} \right. \\ \left. \right) \end{cases} \quad (1.1)$$

1.3 Расстояние Дамерау-Левенштейна

Как было написано выше, в расстоянии Дамерау-Левенштейна задействует еще одну редакторскую операцию - транспозицию T (*от англ. transposition*). Расстояние Дамерау-Левенштейна рассчитывается по рекуррентной формуле. (1.2).

$$D(S_1[1...i], S_2[1...j]) = \begin{cases} j, \text{ если } i == 0 \\ i, \text{ если } j == 0 \\ \min(\\ D(S_1[1...i], S_2[1...j-1]) + 1, \\ D(S_1[1...i-1], S_2[1...j]) + 1, \\ D(S_1[1...i-1], S_2[1...j-1]) + \\ \left[\begin{array}{l} 0, \text{ если } S_1[i] == S_2[j] \\ 1, \text{ иначе,} \end{array} \right. \\ \left. \begin{array}{l} D(S_1[1...i-2], S_2[1...j-2]) + 1, \text{ если } i, j > 1, a_i = b_{j-1}, b_j = a_{i-1} \\ \infty, \text{ иначе } \end{array} \right) \end{cases} \quad j > 0$$

(1.2)

1.4 Вывод

Были рассмотрены основополагающие материалами, которые в дальнейшем потребуются при реализации алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна.

2 Конструкторский раздел

В данном разделе мы рассмотрим схемы вышеизложенных алгоритмов.

2.1 Разработка алгоритмов

На рис. (2.1) представлен алгоритм поиска расстояние Левенштейна.

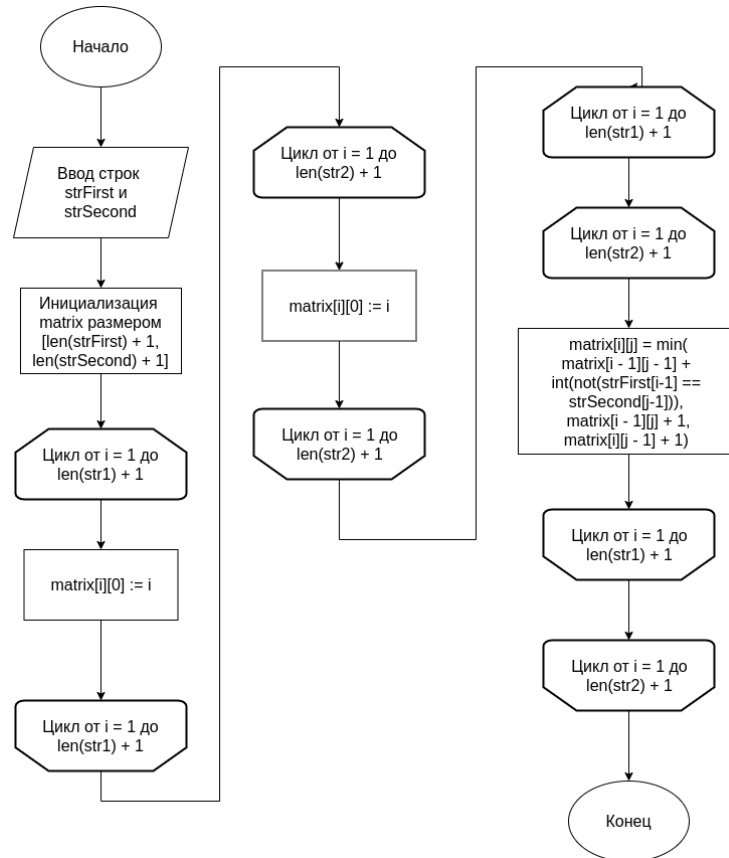


Рисунок 2.1 — Схема алгоритма поиска расстояния Левенштейна

На рис. (2.2) представлен алгоритм поиска расстояния Дамерау-Левенштейна.

На рис. (2.3) представлен рекурсивный алгоритм поиска расстояния Левенштейна.

2.2 Вывод

В данном разделе мы рассмотрели схемы алгоритмов Левенштейна (2.1). и Дамерау-Левенштейна (2.2).

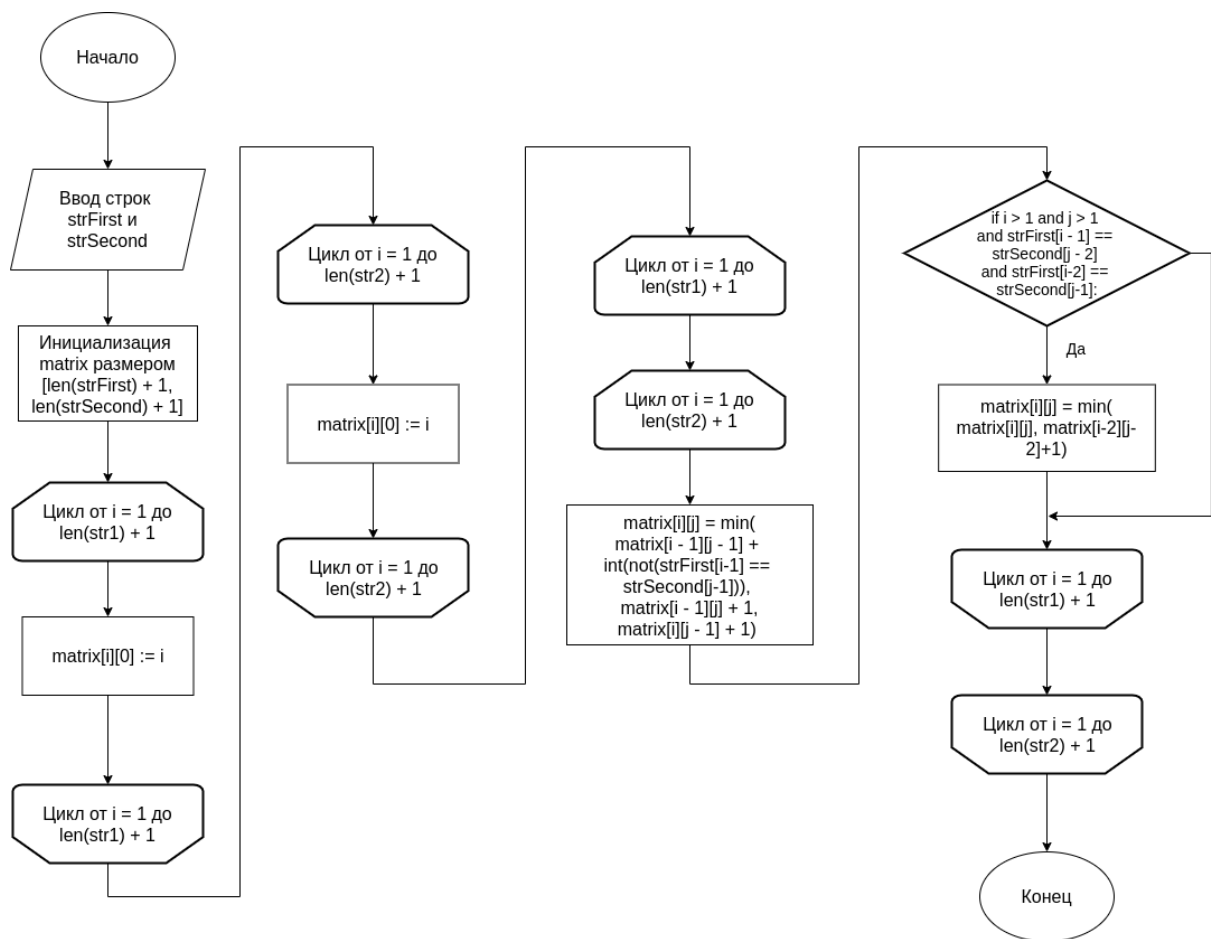


Рисунок 2.2 — Схема алгоритма поиска расстояния Дameraу-Левенштейна

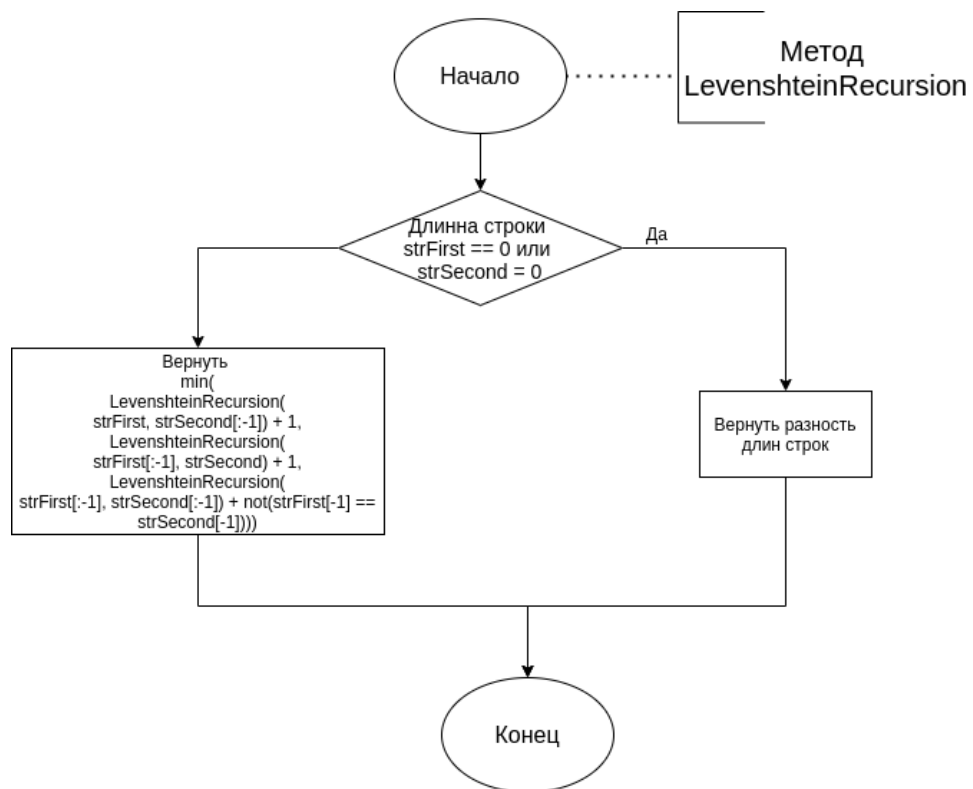


Рисунок 2.3 — Схема рекурсивного алгоритма поиска расстояния Левенштейна

3 Технологический раздел

3.1 Выбор ЯП

В данной лабораторной работе использовался язык программирования - python [1]. Данный язык простой и понятный, также я знакома с ним. Поэтому данный язык был выбран. В качестве среды разработки я использовала Visual Studio Code [2], т.к. считаю его достаточно удобным и легким. Visual Studio Code подходит не только для Windows [3], но и для Linux [4], это еще одна причина, по которой я выбрала VS code, т.к. у меня установлена ОС Ubuntu 18.04.4 [5]. При замерах времени использовалась функция `process_time` [6]

3.2 Требования к программному обеспечению

Входными данными являются две строки. Строки регистрозависимые. На выходе имеется матрица и дистанция, полученная алгоритмом Левенштейна. Результат алгоритма Дамерау-Левенштейна выводится только в том случае, когда его результат не совпадает с результатом алгоритма Левенштейна. Также требуется замерить время работы каждой реализации.

3.3 Сведения о модулях программы

Данная программа разбита на модули:

- `main.py` - Файл, содержащий точку входа в программу. В нем происходит общение с пользователем и вызов алгоритмов;
- `algorithms.py` - Файл содержит непосредственно сами алгоритмы;
- `test_time.py` - Файл с замером времени работы алгоритмов.

На листингах представлены коды разобранных ранее алгоритмов.

Листинг 3.1 — Главная функция `main`

```
1 def main():
2     output("Enter the first string", GREEN) # Column
3     strFirst = input()
4     output("Enter the second string:", GREEN) # Row
5     strSecond = input()
6     distanceLev = Levenshtein(strFirst, strSecond, True)
7     distanceDamLev = DamerauLevenshtein(strFirst, strSecond)
8     output("distance Levenshtein: " + str(distanceLev), GREEN)
9     if distanceLev != distanceDamLev:
10        output("distance Damerau-Levenshtein: " + str(distanceDamLev),
                GREEN)
```

Листинг 3.2 — Функция нахождения расстояния Левенштейна матрично


```

1 def Levenshtein(strFirst, strSecond, flag=False):
2     n, m = len(strFirst), len(strSecond)
3     matrix = np.full((n + 1, m + 1), 0) # math.inf)
4
5     matrix[0][0] = 0
6     for i in range(1, n + 1):
7         matrix[i][0] = i
8     for i in range(1, m + 1):
9         matrix[0][i] = i
10
11    for i in range(1, n + 1):
12        for j in range(1, m + 1):
13            matrix[i][j] = min(
14                matrix[i - 1][j - 1] +
15                int(not(strFirst[i - 1] == strSecond[j - 1])), # R
16                matrix[i - 1][j] + 1, # D
17                matrix[i][j - 1] + 1) # I
18
19    if flag:
20        OutputMatrix(matrix, strFirst, strSecond)
21
22    return matrix[-1][-1]

```

Листинг 3.3 — Рекурсивная функция нахождения расстояния Левенштейна

```

1 def LevenshteinRecursion(strFirst, strSecond):
2     if (strFirst == "" or strSecond == ""):
3         return abs(len(strFirst) - len(strSecond))
4
5     temp = 0 if strFirst[-1] == strSecond[-1] else 1
6     return min(
7         LevenshteinRecursion(strFirst, strSecond[:-1]) + 1, # I
8         LevenshteinRecursion(strFirst[:-1], strSecond) + 1, # D
9         LevenshteinRecursion(strFirst[:-1], strSecond[:-1]) + temp # R
10    )

```

Листинг 3.4 — Функция нахождения расстояния Дамерау-Левенштейна матрично

```

1 def DamerauLevenshtein(strFirst, strSecond, flag=False):
2     n, m = len(strFirst), len(strSecond)
3     matrix = np.full((n + 1, m + 1), 0) # math.inf)
4
5     matrix[0][0] = 0
6     for i in range(1, n + 1):
7         matrix[i][0] = i
8     for i in range(1, m + 1):
9         matrix[0][i] = i

```

```

10
11     for i in range(1, n + 1):
12         for j in range(1, m + 1):
13             matrix[i][j] = min(
14                 matrix[i - 1][j - 1] +
15                 int(not(strFirst[i - 1] == strSecond[j - 1])), # R
16                 matrix[i - 1][j] + 1, # D
17                 matrix[i][j - 1] + 1) # I
18             if i > 1 and j > 1 and strFirst[i - 1] == strSecond[j - 2] \
19                 and strFirst[i - 2] == strSecond[j - 1]:
20                 matrix[i][j] = min(matrix[i][j], matrix[i - 2][j - 2] + 1) # T
21
22     if flag:
23         OutputMatrix(matrix, strFirst, strSecond)
24     return matrix[-1][-1]

```

Листинг 3.5 — Рекурсивная функция нахождения расстояния
Дамерау-Левенштейна

```

1 def DamerauLevenshteinRecursion(strFirst, strSecond):
2     if (strFirst == "" or strSecond == ""):
3         return abs(len(strFirst) - len(strSecond))
4
5     temp = 0 if strFirst[-1] == strSecond[-1] else 1
6     result = min(
7         DamerauLevenshteinRecursion(
8             strFirst, strSecond[:-1]) + 1, # I
9         DamerauLevenshteinRecursion(
10            strFirst[:-1], strSecond) + 1, # D
11         DamerauLevenshteinRecursion(
12            strFirst[:-1], strSecond[:-1]) + temp # R
13     )
14
15     if len(strFirst) > 1 and len(strSecond) > 1 and \
16         strFirst[-1] == strSecond[-2] and strFirst[-2] == strSecond[-1]:
17         result = min(result, DamerauLevenshteinRecursion(
18             strFirst[:-2], strSecond[:-2]) + 1) # T
19
20     return result

```

Листинг 3.6 — функция замера времени

```

1 def TimeTest(strFirst, strSecond, countOperations):
2     t1 = time.process_time()
3     for _ in range(countOperations):
4         Levenshtein(strFirst, strSecond)
5     t2 = time.process_time()

```

```

6      print("Levenshtein = ", t2 - t1)
7
8      t1 = time.process_time()
9      for _ in range(countOperations):
10         LevenshteinRecursion(strFirst , strSecond)
11      t2 = time.process_time()
12      print("Levenshtein (Recursion)= ", t2 - t1)
13
14      t1 = time.process_time()
15      for _ in range(countOperations):
16         DamerauLevenshtein(strFirst , strSecond)
17      t2 = time.process_time()
18      print("DamerauLevenshtein = ", t2 - t1)
19
20      t1 = time.process_time()
21      for _ in range(countOperations):
22         DamerauLevenshteinRecursion(strFirst , strSecond)
23      t2 = time.process_time()
24      print("DamerauLevenshtein (Recursion)= ", t2 - t1)

```

3.4 Тестирование

В данном разделе будет приведена таблица 3.1, в которой четко отражено тестирование программы. Первый и второй столбец отвечают за введенные пользователем слова.

Таблица 3.1 — Таблица тестов

Слово 1	Слово 2	Ожидаемый вывод	Вывод программы
сито	столб	3	3
exponential	polynomial	6	6
Alice	Alice	0	0
Alice	alice	1	1
ma	am	2 1	2 1
		0	0
abc	cab	2	2

Все тесты пройдены успешно.

3.5 Вывод

В данном разделе были рассмотрены листинги кода, обоснован выбор использованного в данной работе языка программирования и среды разработки, а также

была произведена проверка корректной работы программы, благодаря таблице 3.1. Сравнив представленные листинги можно сказать, что написание рекуррентных подпрограмм проще, чем матричных.

4 Экспериментальная часть

В данном разделе сравним работу каждого алгоритма.

4.1 Временные характеристики

Сравним матричный алгоритм поиска расстояния Левенштейна и Дамерау-Левенштейна. Для сравнения возьмем строки длиной [10, 20, 30, 50, 100, 200]. Так как подсчет расстояния считается короткой задачей, воспользуемся усреднением массового эксперимента. Для этого сложим результат работы алгоритма n раз ($n \geq 10$), после чего поделим на n . Тем самым получим достаточно точные характеристики времени. Сравнение произведем при $n = 500$. Результат можно увидеть на рис. 4.1. При короткой длине разница по времени минимальна, при увеличении длины строки алгоритм поиска расстояния Левенштейна с небольшим опережением вырывается вперед. Это обосновывается тем, что у алгоритма поиска расстояния Дамерау-Левенштейна задействуется еще одна операция, которая замедляет алгоритм.

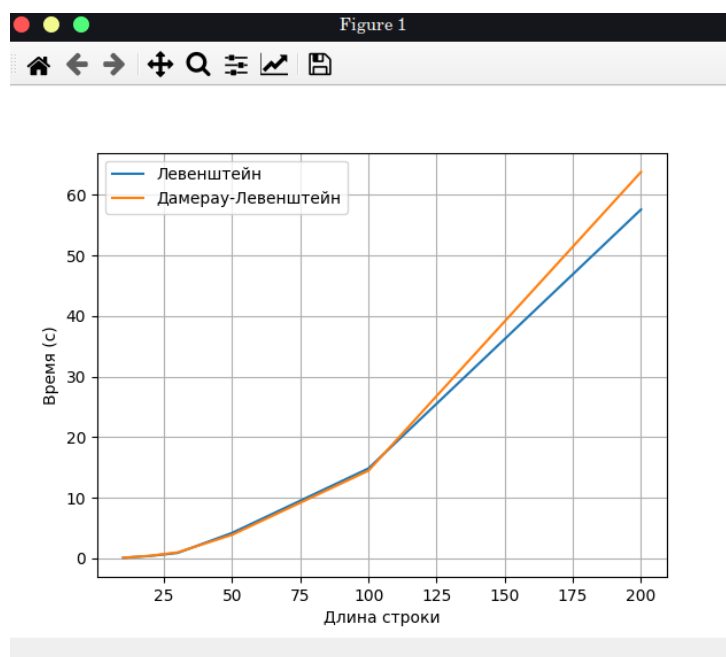


Рисунок 4.1 — Сравнение времени работы алгоритма поиска расстояния Левенштейна и Дамерау-Левенштейна

Далее проведем сравнительный анализ временных характеристик рекурсивной и матричной реализаций алгоритма Левенштейна. Возьмем строки длиной [2, 3, 5, 7, 8], n положим равным 50. Результат можно увидеть на рис. 4.2. Такая большая разница во времени объясняется тем, что в рекурсивном алгоритме Левенштейна много рекурсивных вызовов с однотипными параметрами.

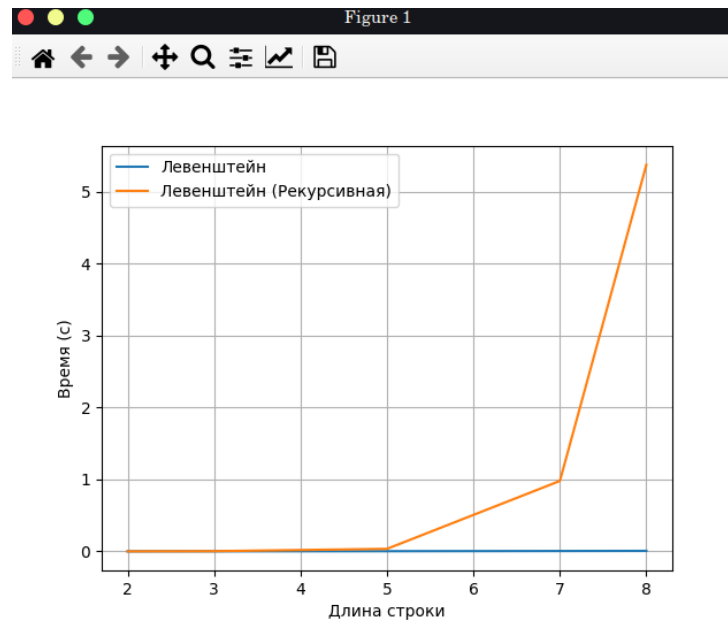


Рисунок 4.2 — Сравнение времени работы рекурсивной и матричной реализаций алгоритма Левенштейна.

4.2 Характеристики по памяти

На рисунке 4.3 представлено дерево вызовов рекурсивного алгоритма Левенштейна. Видно, что на третьем уровне встречаются повторные вызовы. Чем больше будет уровень, тем чаще будут вызываться функции с однотипными аргументами, что может привести к превышению максимальной глубины рекурсии. При строках длиной 2 подпрограмма вызовется 18 раз. Каждый вызов задействует 32 мегабайт (замеры проведены с помощью библиотеки `memogu_profiler` [7]). В итоге нам требуется 576 мегабайт для рекурсивных вызовов, в то время, когда в матричном алгоритме используется 42 мегабайта.

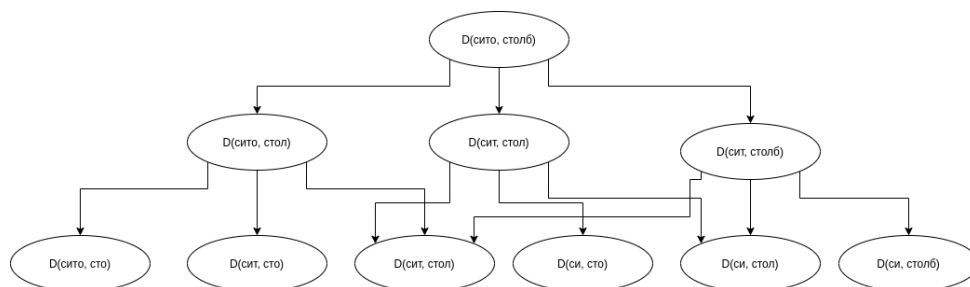


Рисунок 4.3 — Сравнение времени работы рекурсивной и матричной реализаций алгоритма Левенштейна.

4.3 Сравнительный анализ алгоритмов

Приведенные характеристики показывают нам, что рекурсивная реализация алгоритма очень сильно проигрывает по времени и по памяти.

Во время печати очень часто возникают ошибки связанные с транспозицией букв, поэтому алгоритм поиска расстояния Дамерау-Левенштейна предпочтительнее, не смотря на то, что он проигрывает по времени алгоритму Левенштейна.

По аналогии с первым абзацем можно сделать вывод о том, что рекуррентный алгоритм поиска расстояния Дамерау-Левенштейна будет более затратный, как по памяти, так и по времени по сравнению с матричной реализацией алгоритма поиска расстояния Дамерау-Левенштейна.

4.4 Вывод

В данном разделе было произведено сравнение количества затраченного времени и памяти вышеизложенных алгоритмов. Самым быстрым оказался матричный алгоритм нахождения расстояния Левенштейна.

Заключение

Алгоритмы Левенштейна и Дамерау-Левенштейна являются самыми популярными алгоритмами, которые помогают найти редакционное расстояние.

В этой лабораторной работе мы познакомились с алгоритмами Левенштейна (Формула 1.1) и Дамерау-Левенштейна (Формула 1.2). Построили схемы (Рисунок 2.1, Рисунок 2.2), соответствующие данным алгоритмам, также разобрали рекуррентные реализации (Рисунок 2.3). Написали полностью готовый и протестированный (Таблица 3.1) программный продукт, который считает дистанцию 4 способами.

В рамках выполнения работы решены следующие задачи.

- а) Изучены алгоритмы Дамерау-Левенштейна и Левенштейна.
- б) Реализовали изученные алгоритмы, а также матричную и рекурсивную реализации алгоритма.
- в) Проиллюстрировали алгоритмы на схемах.
- г) Сравнили временные характеристики, а также затраченную память.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. *Майкл, Доусон*. Python Programming for the Absolute Beginner, 3rd Edition / Доусон Майкл. — Прогресс книга, 2019. — Р. 416.
2. Visual Studio Code. — Microsoft, 2005. <https://code.visualstudio.com/>.
3. Windows. — Microsoft, 1985. <https://www.microsoft.com/ru-ru/windows>.
4. Linux. — 1991. <https://www.linux.org.ru/>.
5. Ubuntu 18.04. — 2018. <https://releases.ubuntu.com/18.04/>.
6. process time function in Python. https://www.geeksforgeeks.org/time-process_time-function-in-python/.
7. Замер памяти в Python3. — 2019. <https://pypi.org/project/memory-profiler/>.