

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
ИМ. Н. Э. БАУМАНА

УДК _____

№ госрегистрации _____

Инв. № _____

УТВЕРЖДАЮ

головной исполнитель НИР

«_____» _____ 2019 г.

ОТЧЁТ ПО ДИСЦИПЛИНЕ "АНАЛИЗ АЛГОРИТМОВ"
О ЛАБОРАТОРНОЙ РАБОТЕ №1

по теме:

"Расстояние Левенштейна и Дамерау-Левенштейна"

(промежуточный)

Студент ИУ7-53Б

___ Пудов Дмитрий Юрьевич

Москва 2019

СОДЕРЖАНИЕ

Введение	3
1 Аналитическая часть.....	4
1.1 Описание алгоритмов.....	4
2 Конструкторская часть	6
2.1 Разработка алгоритмов	6
2.2 Сравнительный анализ рекурсивной и нерекурсивной реализаций	14
3 Технологическая часть.....	15
3.1 Требования к программному обеспечению	15
3.2 Средства реализации	15
3.3 Листинг кода.....	15
3.4 Описание тестирования	18
4 Экспериментальная часть	19
4.1 Примеры работы	19
4.2 Результаты тестирования	20
4.3 Постановка эксперимента по замеру времени и памяти	23
4.4 Сравнительный анализ на материале экспериментальных данных	23
4.5 Выводы	24
Заключение	25

ВВЕДЕНИЕ

Цель работы: изучение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна.

Постановка задачи:

- изучить метод динамического программирования на материалах алгоритмов Левенштейна и Дамерау-Левенштейна;
- применить его;
- получить практические навыки реализации указанных алгоритмов.

1 Аналитическая часть

В данной части будут описаны суть задач нахождения расстояния Левенштейна и Дameraу-Левенштейна, а также математические способы их решения.

Поиск расстояния Левенштейна заключается в определении минимального количества редакционных операций (вставка, удаление, замена одного символа), необходимых для трансформации одной строки в другую.

В задаче о расстоянии Дameraу-Левенштейна во множество редакционных операций добавляется транспозиция - перестановка двух соседних символов.

1.1 Описание алгоритмов

Пусть a и b - строки над некоторым алфавитом длины M и N соответственно. Тогда расстояние Левенштейна определяется формулой $d(a, b) = D(M, N)$, где

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

Расстояние Дамерау-Левенштейна для данных строк определяется как $d(a, b) = D(M, N)$, где

$$D(i, j) = \begin{cases} \max(i, j) & \min(i, j) = 0 \\ \min \begin{cases} D(i, j) + 1 \\ D(i, j - 1) + 1 \\ D(i - 1, j - 1) + 1_{a[i] \neq b[j]} \\ D(i - 2, j - 2) + 1 \end{cases} & i > 1 \text{ и } j > 1 \\ \text{и } a[i] = b[j - 1] \text{ и } a[i - 1] = b[j] \\ \min \begin{cases} D(i - 1, j) + 1 \\ D(i, j - 1) + 1 \\ D(i - 1, j - 1) + 1_{a[i] \neq b[j]} \end{cases} & \text{иначе.} \end{cases}$$

2 Конструкторская часть

В данной части будут приведены схемы алгоритмов Левенштейна в рекурсивной и матричной реализации и Дамерау-Левенштейна.

2.1 Разработка алгоритмов

Далее указаны разработанные схемы алгоритмов Левенштейна и Дамерау-Левенштейна. Будем считать, что известны следующие функции: определения длины строки, поиска максимума и минимума среди нескольких чисел. Для матричных реализаций требуется наличие функции, динамически выделяющей память.

На рис. 2.1 представлена схема алгоритма определения расстояния Левенштейна в рекурсивной реализации.

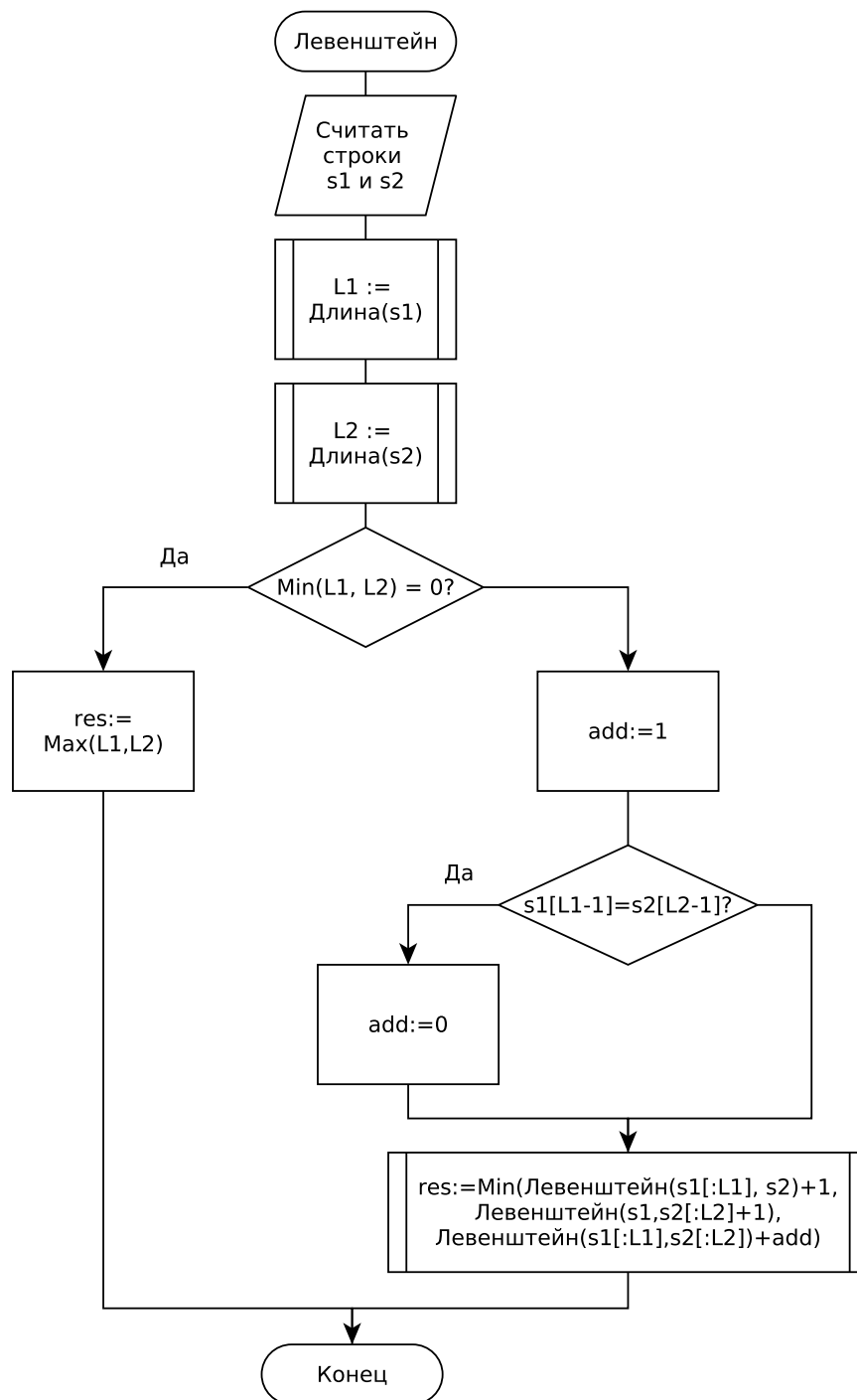


Рисунок 2.1 — Схема алгоритма определения расстояния Левенштейна в рекурсивной реализации.

На рис. 2.2 и на рис. 2.3 представлена 1 часть схемы алгоритма определения расстояния Левенштейна в матричной реализации.

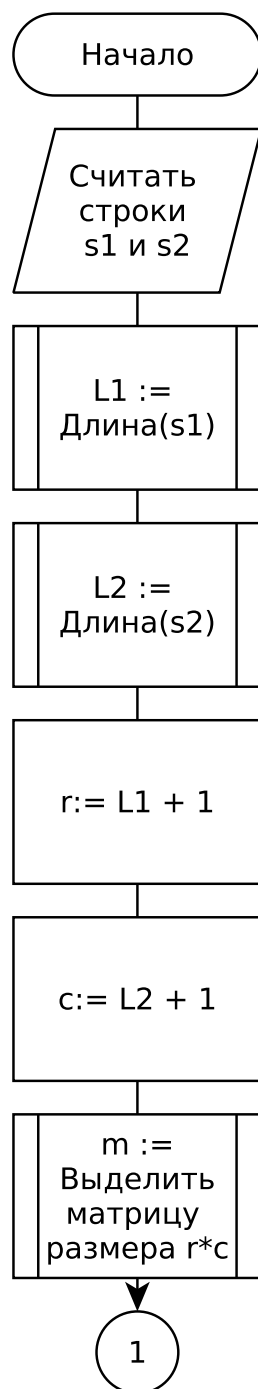


Рисунок 2.2 — Схема алгоритма определения расстояния Левенштейна в матричной реализации. Часть 1.

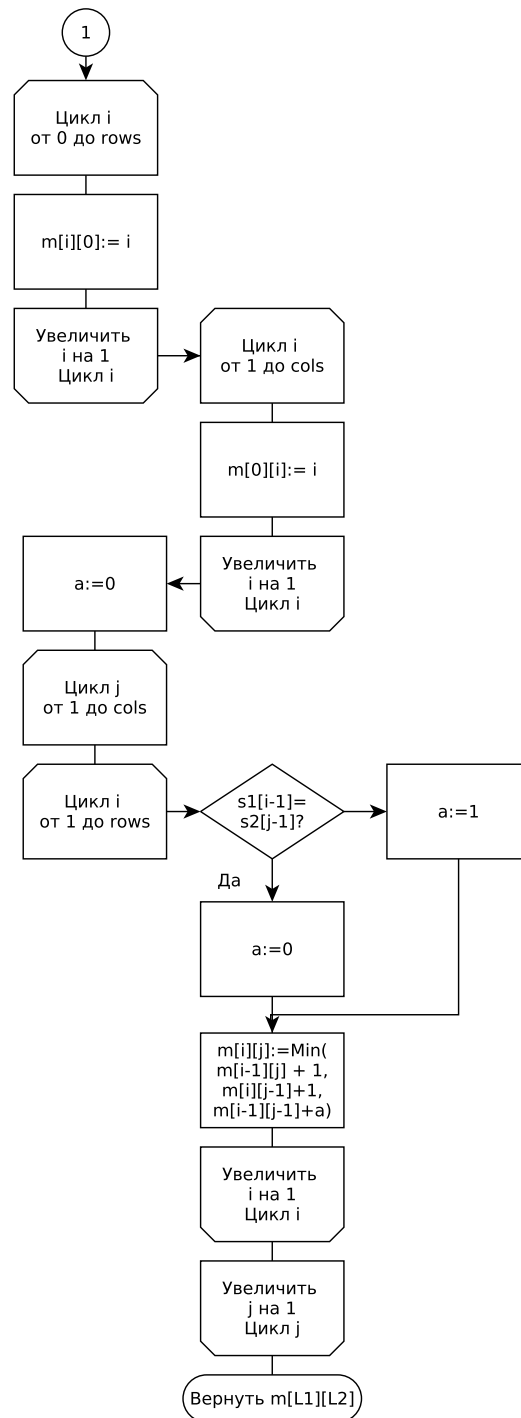


Рисунок 2.3 — Схема алгоритма определения расстояния Левенштейна в матричной реализации. Часть 2.

На рис. 2.4 и рис. 2.5 представлена схема алгоритма определения расстояния Дамерау-Левенштейна в матричной реализации.

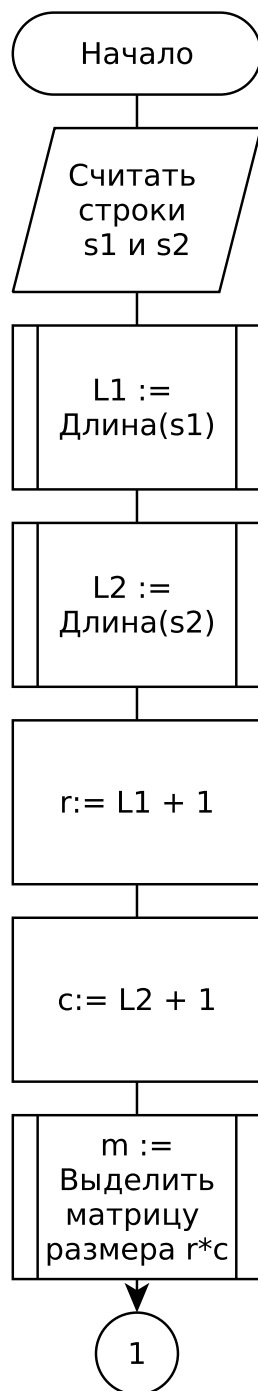


Рисунок 2.4 — Схема алгоритма определения расстояния Дамерау-Левенштейна.
Часть 1.

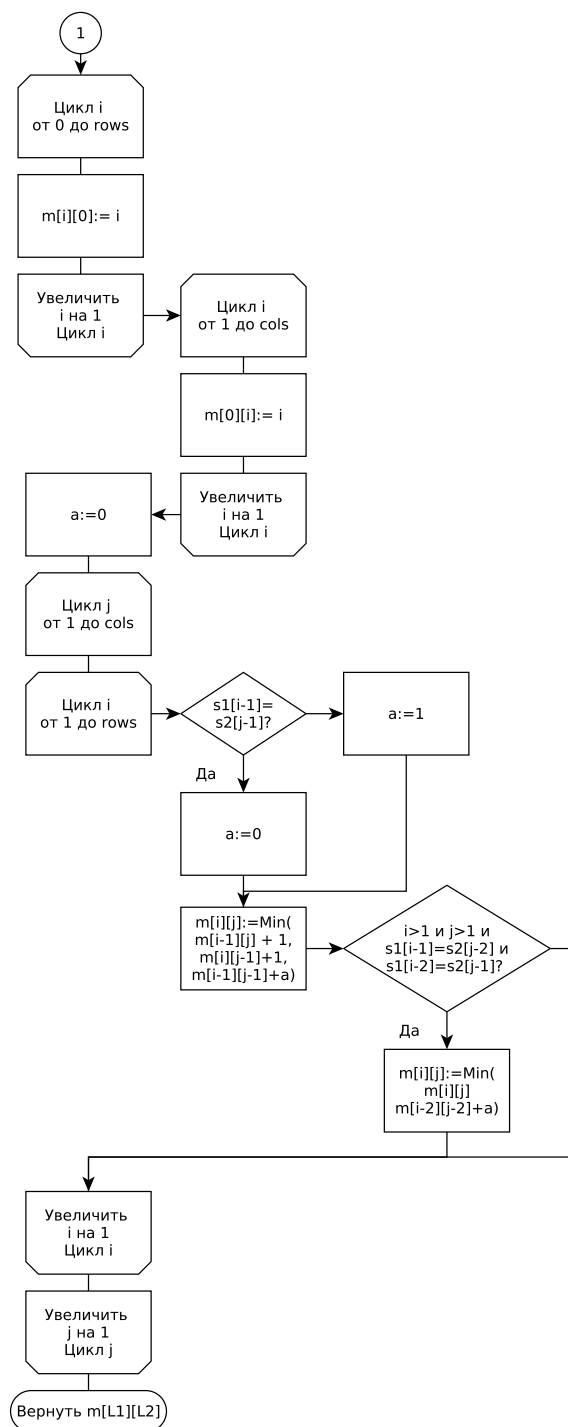


Рисунок 2.5 — Схема алгоритма определения расстояния Дамерау-Левенштейна.
Часть 2.

2.2 Сравнительный анализ рекурсивной и нерекурсивной реализаций

Рекурсивная реализация по сравнению с матричной будет иметь большую сложность. В матричной количество итераций заранее известно и будет равно $m * n$, где m - число строк, а n - число столбцов.

При этом в рекурсивной версии сложность может достигать $3^{max(m,n)}$, так как на каждый вызов функции может требоваться ещё 3. При этом в данном дереве рекурсии на некоторых ветвях будут производиться повторные вычисления, что неэффективно.

3 Технологическая часть

В данной части будет описан подход для программной реализации ранее указанных алгоритмов.

3.1 Требования к программному обеспечению

Разработанное ПО должно предоставлять возможность замеров процессорного времени выполнения каждого алгоритма. Требуется вводить 2 строки и выводить матрицу (кроме рекурсивной реализации) и значения расстояний, полученных различными реализациями.

3.2 Средства реализации

Для реализации ПО был выбран язык Go, поскольку в его основном пакете присутствуют удобные средства замера процессорного времени и памяти модуля testing.

3.3 Листинг кода

В данном разделе приведены листинги реализаций ранее указанных алгоритмов на языке Go.

Листинг 3.1 — Матричная версия алгоритма Дамерау-Левенштейна

```
1 func LevenshteinDamerau(first string, second string) (int, [][] int) {
2     lenFirst := len(first)
3     lenSecond := len(second)
4     rows := lenFirst + 1
5     cols := lenSecond + 1
6     matrix := allocateMatrix(rows, cols)
7
8     for i := 0; i < rows; i++ {
9         matrix[i][0] = i
10    }
11    for i := 1; i < cols; i++ {
12        matrix[0][i] = i
13    }
14
15    add := 0
```

```

16     for j := 1; j < cols; j++ {
17         for i := 1; i < rows; i++ {
18             if first[i-1] == second[j-1] {
19                 add = 0
20             } else {
21                 add = 1
22             }
23             matrix[i][j] = MinThree(matrix[i-1][j]+1,
24                                     matrix[i][j-1]+1,
25                                     matrix[i-1][j-1]+add)
26
27             if i > 1 && j > 1 && first[i-1] == second[j-2] && first[i-2] ==
                second[j-1] {
28                 matrix[i][j] = Min(matrix[i][j], matrix[i-2][j-2]+add)
29             }
30         }
31     }
32
33     return matrix[lenFirst][lenSecond], matrix
34 }

```

Листинг 3.2 — Матричная реализация алгоритма Левенштейна

```

1  func LevenshteinIterative(first string, second string) (int, [][]int) {
2      lenFirst := len(first)
3      lenSecond := len(second)
4      rows := lenFirst + 1
5      cols := lenSecond + 1
6      matrix := allocateMatrix(rows, cols)
7
8      for i := 0; i < rows; i++ {
9          matrix[i][0] = i
10     }
11     for i := 1; i < cols; i++ {
12         matrix[0][i] = i
13     }
14
15     add := 0
16     for j := 1; j < cols; j++ {
17         for i := 1; i < rows; i++ {
18             if first[i-1] == second[j-1] {
19                 add = 0
20             } else {
21                 add = 1
22             }
23             matrix[i][j] = MinThree(matrix[i-1][j]+1,
24                                     matrix[i][j-1]+1,

```



```

25             matrix[i-1][j-1]+add)
26         }
27     }
28
29     return matrix[lenFirst][lenSecond], matrix
30 }

```

Листинг 3.3 — Рекурсивная версия алгоритма Левенштейна

```

1  func LevenshteinRecursive(first string , second string) int {
2      lenFirst := len(first)
3      lenSecond := len(second)
4      if Min(lenFirst , lenSecond) == 0 {
5          return Max(lenFirst , lenSecond)
6      } else {
7          add := 1
8          if first[lenFirst-1] == second[lenSecond-1] {
9              add = 0
10         }
11
12         return MinThree(
13             LevenshteinRecursive(first[:lenFirst-1], second)+1,
14             LevenshteinRecursive(first , second[:lenSecond-1])+1,
15             LevenshteinRecursive(first[:lenFirst-1],
16                                     second[:lenSecond-1])+add)
17     }
18 }

```

Листинг 3.4 — Оптимизированная версия рекурсивного алгоритма Левенштейна

```

1  func levenshteinRecursiveModule(first , second string , result , minval int)
2      int {
3      lenFirst := len(first)
4      lenSecond := len(second)
5      if result >= minval {
6          return minval
7      } else if lenFirst == 0 {
8          return result + lenSecond
9      } else if lenSecond == 0 {
10         return result + lenFirst
11     } else {
12         add := 1
13         if first[lenFirst-1] == second[lenSecond-1] {
14             add = 0
15         }
16         r1 := levenshteinRecursiveModule(first[:lenFirst-1],
17             second[:lenSecond-1], result+add , minval)

```

```

16         r2 := levenshteinRecursiveModule(first, second[:lenSecond-1],
           result+1, Min(r1, minval))
17         r3 := levenshteinRecursiveModule(first[:lenFirst-1], second,
           result+1, MinThree(r1, r2, minval))
18         return MinThree(r1, r2, r3)
19     }
20 }
21
22 func LevenshteinRecursiveOptimized(first string, second string) int {
23     minval := Max(len(first), len(second))
24     result := 0
25     return levenshteinRecursiveModule(first, second, result, minval)
26 }

```

3.4 Описание тестирования

Тестирование будет проведено для каждой из реализаций со следующими входными данными:

- пустые строки;
- идентичные строки;
- строки одинаковой длины, но с разными символами;
- строки различной длины, но отличных лишь в некотором числе вставок или удалений символа;
- строки различной длины и требующих замены символа.

Для расстояния Дameraу-Левенштейна дополнительный тест будет заключаться в перестановке двух соседних символов.

4 Экспериментальная часть

В данной части будет проведена апробация и тестирование разработанной программы.

4.1 Примеры работы

Программа имеет консольный интерфейс. Далее будут приведены скриншоты.

На рис. 4.1 показаны возможные сценарии запуска программы.

```
dpudov@dpudov-Inspiron-5547:/media/dpudov/media/workspace/fifth_semester/algorithm-analysis/levenshtein-distance/src(master)$ ./levenshtein-distance
You may run program with those arguments:
once Input two strings, get result. May be combined with 'visualize'
visualize Prints result matrices
many Takes several inputs until kill signal
```

Рисунок 4.1 — Использование программы.

На рис. 4.2 показан пример использования программы с однократным вводом входных данных.

```
dpudov@dpudov-Inspiron-5547:/media/dpudov/media/workspace/fifth_semester/algorithm-analysis/levenshtein-distance/src(master)$ ./levenshtein-distance once
Please input string:
abcd
Please input string:
abcde
Result is...
Strings: abcd abcde
For Levenshtein recursive: 1
For Levenshtein iterative: 1
For Levenshtein-Damerau: 1
```

Рисунок 4.2 — Однократное использование.

На рис. 4.3 показан пример использования программы с выводом матриц.

```
dpudov@dpudov-Inspiron-5547:/media/dpudov/media/workspace/fifth_semester/algorithm-analysis/levenshtein-distance/src(master)$ ./levenshtein-distance once visualize
Please input string:
abcd
Please input string:
abcde
Result is...
Strings: abcd abcde
For Levenshtein recursive: 1
For Levenshtein iterative: 1
s1/s2 a b c d e
  0 1 2 3 4 5
a 1 0 1 2 3 4
b 2 1 0 1 2 3
c 3 2 1 0 1 2
d 4 3 2 1 0 1
For Levenshtein-Damerau: 1
s1/s2 a b c d e
  0 1 2 3 4 5
a 1 0 1 2 3 4
b 2 1 0 1 2 3
c 3 2 1 0 1 2
d 4 3 2 1 0 1
```

Рисунок 4.3 — Использование с выводом матриц.

На рис. 4.4 показан пример использования программы со множеством ВХОДНЫХ ДАННЫХ.

```
dpudov@dpudov-Inspiron-5547:/media/dpudov/media/workspace/fifth_semester/algorithm-analysis/levenshtein-distance/src(master)$ ./levenshtein-distance many
Please input string:
abcd
Please input string:
abcde
Result is...
Strings: abcd abcde
For Levenshtein recursive: 1
For Levenshtein iterative: 1
For Levenshtein-Damerau: 1
Please input string:
dcbe
Please input string:
abc
Result is...
Strings: dcbe abc
For Levenshtein recursive: 3
For Levenshtein iterative: 3
For Levenshtein-Damerau: 3
Please input string:
```

Рисунок 4.4 — Многократное использование.

4.2 Результаты тестирования

В данном разделе будут расположены результаты тестирования ранее описанных реализаций.

В табл. 4.1 отражены результаты тестирования рекурсивной реализации алгоритма поиска расстояния Левенштейна.

Таблица 4.1 — Результаты тестирования рекурсивной реализации.

S1	S2	Ожидаемый результат	Фактический результат
Пустая строка	Пустая строка	0	0
"abc"	"abc"	0	0
"abc"	"def"	3	3
"abcdef"	"abc"	3	3
"abc"	"abcdef"	3	3
"xyzabc"	"abcdef"	6	6
"abkdef"	"abcdefg"	2	2
"abc"	"acb"	2	2

В табл. 4.2 отражены результаты тестирования матричной реализации алгоритма поиска расстояния Левенштейна.

В табл. 4.3 отражены результаты тестирования матричной реализации алгоритма поиска расстояния Дамерау-Левенштейна.

Таблица 4.2 — Результаты тестирования матричной реализации.

S1	S2	Ожидаемый результат	Фактический результат
Пустая строка	Пустая строка	0	0
"abc"	"abc"	0	0
"abc"	"def"	3	3
"abcdef"	"abc"	3	3
"abc"	"abcdef"	3	3
"xyzabc"	"abcdef"	6	6
"abkdef"	"abcdefg"	2	2
"abc"	"acb"	2	2

Таблица 4.3 — Результаты тестирования матричной реализации Дамерау-Левенштейна.

S1	S2	Ожидаемый результат	Фактический результат
Пустая строка	Пустая строка	0	0
"abc"	"abc"	0	0
"abc"	"def"	3	3
"abcdef"	"abc"	3	3
"abc"	"abcdef"	3	3
"xyzabc"	"abcdef"	6	6
"abkdef"	"abcdefg"	2	2
"abc"	"acb"	1	1

4.3 Постановка эксперимента по замеру времени и памяти

Для матричных алгоритмов эксперимент ставился со следующими условиями:

- число повторов эксперимента - 100
- диапазоны сравниваемых строк: от 100 до 1000 с шагом 100
- все строки получены случайным образом.

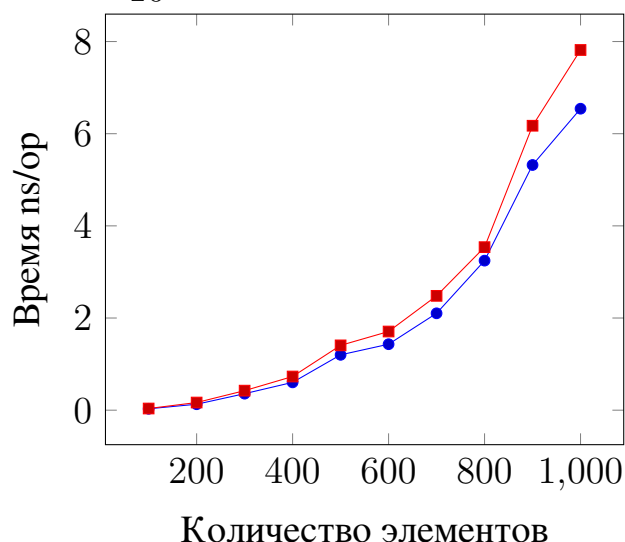
Для рекурсивного алгоритма время на одну итерацию достигает 1.75 мс, при длине строки в 12 символов. Аналогичный эксперимент для матричных реализаций занимает 4 мкс. Поэтому рекурсивная реализация была исключена из дальнейшего исследования.

4.4 Сравнительный анализ на материале экспериментальных данных

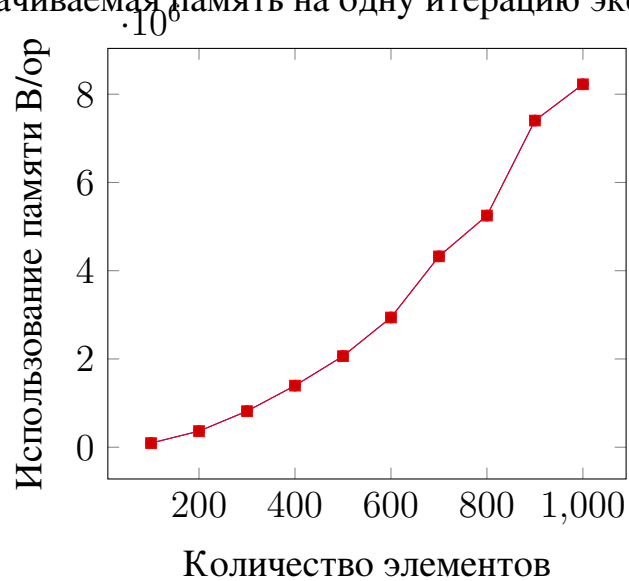
В данном разделе будут приведены результаты экспериментов.

Далее приведены сравнительные характеристики матричных реализаций.

Затрачиваемое время на одну итерацию эксперимента



Затрачиваемая память на одну итерацию эксперимента



4.5 Выводы

Использование памяти матричными реализациями алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна примерно одинаково, поскольку размерность матрицы одна и та же.

По времени матричный алгоритм поиска расстояния Дамерау-Левенштейна медленнее матричного алгоритма Левенштейна за счёт дополнительного условия на штраф за транспозицию.

ЗАКЛЮЧЕНИЕ

В результате работы был изучен метод динамического программирования на материала алгоритмов Левенштейна и Дамерау-Левенштейна.

Была разработана программа реализующая алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна и получены экспериментальные данные сравнения этих алгоритмов по времени и памяти.