

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ «МОСКОВСКИЙ
ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ ИМЕНИ Н.Э.
БАУМАНА (НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ)»
(МГТУ им. Н.Э. БАУМАНА)

ЛАБОРАТОРНАЯ РАБОТА №1

ПО КУРСУ: "АНАЛИЗ АЛГОРИТМОВ"

Расстояние Левенштейна

Работу выполнил: Саркисов Артём, ИУ7-53Б

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

Москва, 2020

Оглавление

Введение	2
1 Аналитическая часть	4
2 Конструкторская часть	6
2.1 Схемы алгоритмов	6
3 Технологическая часть	13
3.1 Выбор ЯП	13
3.2 Реализация алгоритма	13
4 Исследовательская часть	17
4.1 Сравнительный анализ на основе замеров времени работы алгоритмов	17
4.2 Тесты	18
Заключение	20
Литература	20

Введение

Расстояние Левенштейна - минимальное количество операций вставки одного символа, удаления одного символа и замены одного символа на другой, необходимых для превращения одной строки в другую.

Расстояние Левенштейна применяется в теории информации и компьютерной лингвистике в следующих случаях:

- Исправления ошибок в слове(поисковая строка браузера);
- Сравнения текстовых файлов утилитой diff;
- В биоинформатике для сравнения генов, хромосом и белков.

Целью данной лабораторной работы является изучение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна.

Задачами данной лабораторной являются:

1. Изучение алгоритмов Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками;
2. Применение метода динамического программирования для матричной реализации указанных алгоритмов;
3. Получение практических навыков реализации указанных алгоритмов: двух алгоритмов в матричной версии и двух алгоритмов в рекурсивной версии;
4. Сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
5. 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;

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

1. Аналитическая часть

Задача по нахождению расстояния Левенштейна заключается в поиске минимального количества операций вставки/удаления/замены для превращения одной строки в другую.

При нахождении расстояния Дameraу — Левенштейна добавляется операция транспозиции (перестановки соседних символов).

ействия обозначаются следующим образом:

1. D (англ. delete) — удалить;
2. I (англ. insert) — вставить;
3. R (replace) — заменить;
4. M(match) - совпадение.

Пусть S_1 и S_2 — две строки (длиной M и N соответственно) над некоторым алфавитом, тогда расстояние Левенштейна можно подсчитать по следующей рекуррентной формуле (1.1):

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min(& \\ D(i, j - 1) + 1, & \\ D(i - 1, j) + 1, & j > 0, i > 0 \\ D(i - 1, j - 1) + m(S_1[i], S_2[j]) & \\), & \end{cases} \quad (1.1)$$

где $m(a, b)$ равна нулю, если $a = b$ и единице в противном случае; $\min\{a, b, c\}$ возвращает наименьший из аргументов.

Расстояние Дameraу-Левенштейна вычисляется по следующей рекуррентной формуле (1.2):

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & i > 0, j = 0 \\ j, & i = 0, j > 0 \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[i]), \\ D(i - 2, j - 2) + m(S_1[i], S_2[i]), \end{cases} & \begin{array}{l} \text{, если } i, j > 0 \\ \text{и } S_1[i] = S_2[j - 1] \\ \text{и } S_1[i - 1] = S_2[j] \end{array} \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + m(S_1[i], S_2[i]), \end{cases} & \text{, иначе} \end{cases} \quad (1.2)$$

Вывод

В данном разделе были рассмотрены алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна, который является модификаций первого, учитывающего возможность перестановки соседних символов.

2. Конструкторская часть

Требования к программе:

1. На вход подаются две строки.
2. Буквы в верхнем регистре и в нижнем регистре считаются разными.
3. Две пустые строки - корректный ввод, программа не должна аварийно завершаться.
4. Для всех алгоритмов выводиться процессорное время исполнения.
5. Для всех алгоритмов кроме Левенштейна с рекурсивной реализацией должна выводиться матрица.

2.1 Схемы алгоритмов

В данной части будут рассмотрены схемы следующих алгоритмов Левенштейна [1]: рекурсивного алгоритма Левенштейна (рис. 2.1), рекурсивного алгоритма Левенштейна с матричной оптимизацией (рис. 2.2), матричного алгоритма Левенштейна (рис. 2.3), матричного алгоритма Левенштейна (рис. 2.4).

Вывод

В данном разделе были рассмотрены требования к работе, а также схемы алгоритмов нахождения редакционного расстояния с помощью алгоритмов Левенштейна и Дамерау-Левенштейна.

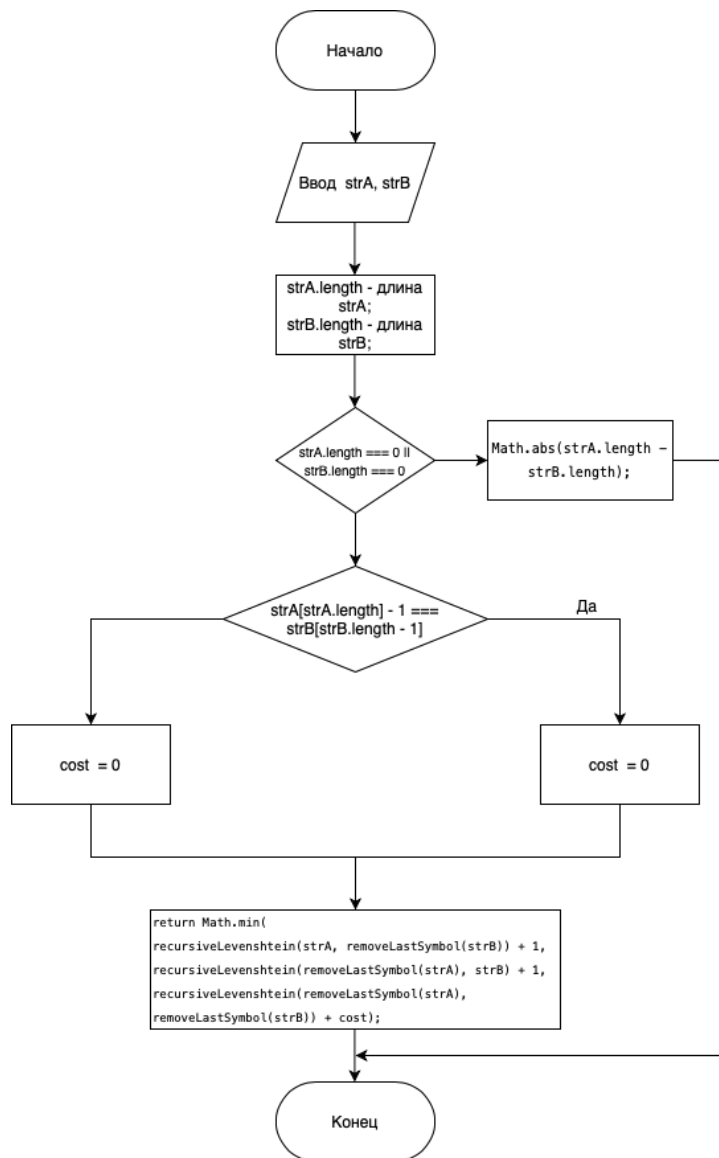


Рис. 2.1: Схема рекурсивного алгоритма нахождения расстояния Левенштейна

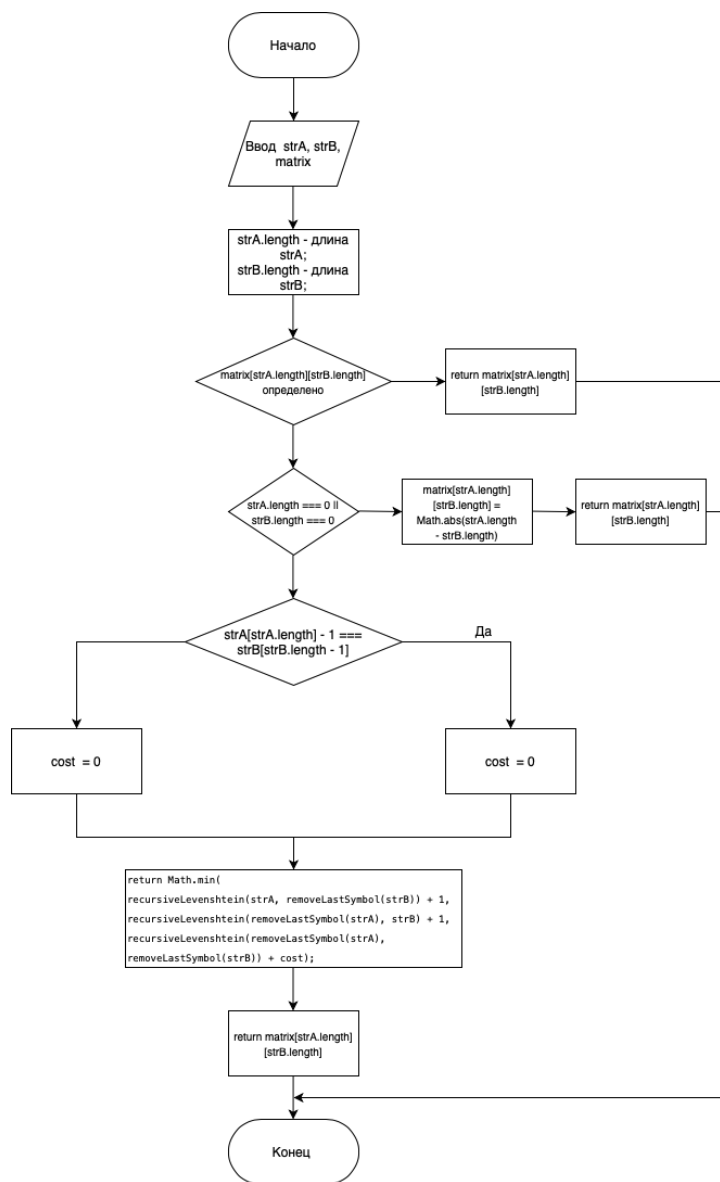
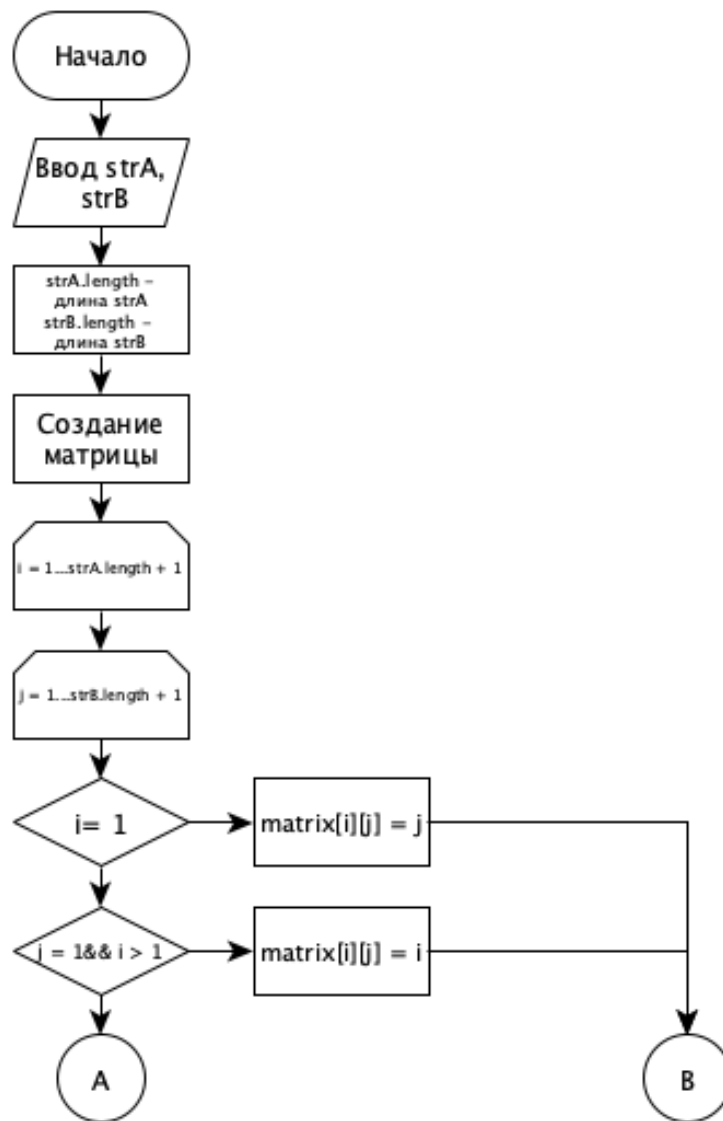


Рис. 2.2: Схема рекурсивного алгоритма нахождения расстояния Левенштейна с матричной оптимизацией



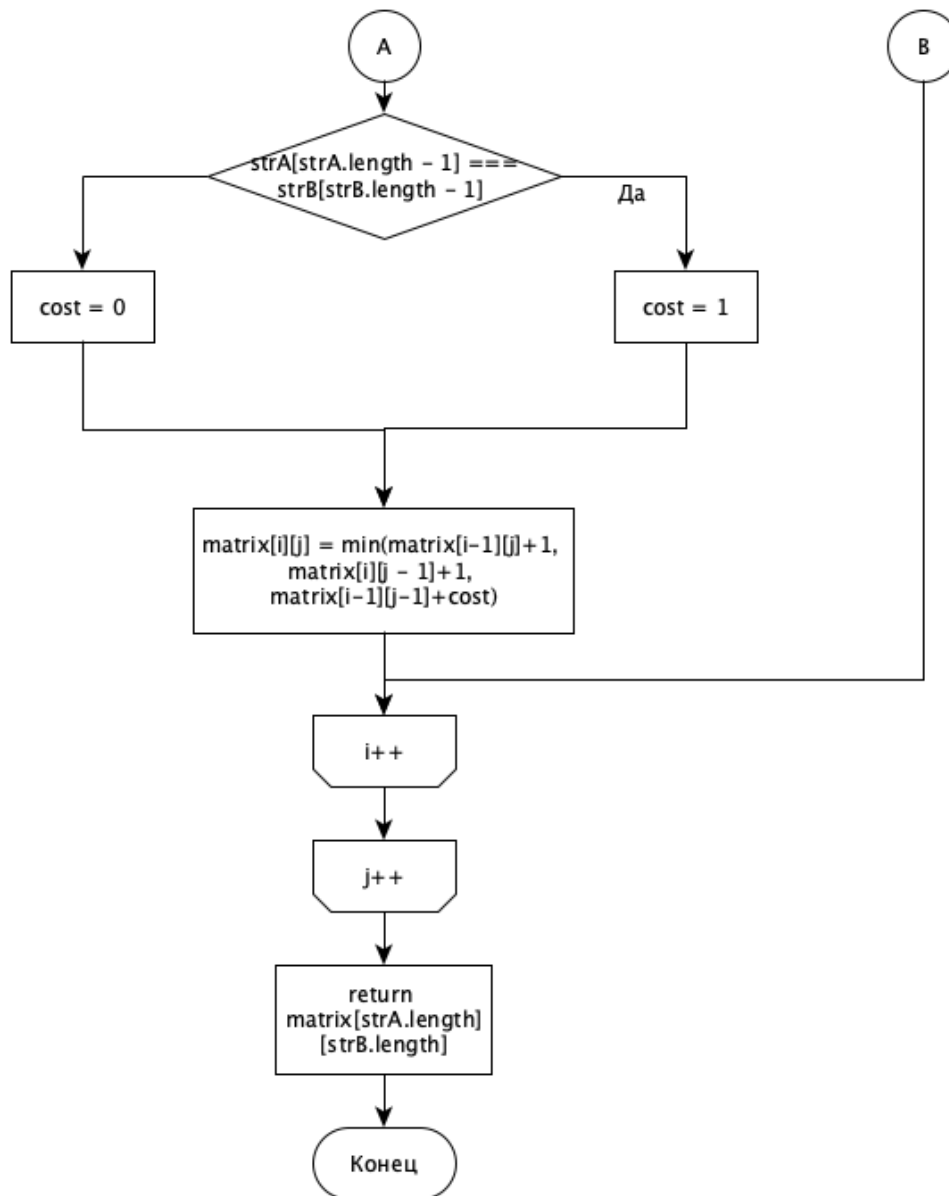
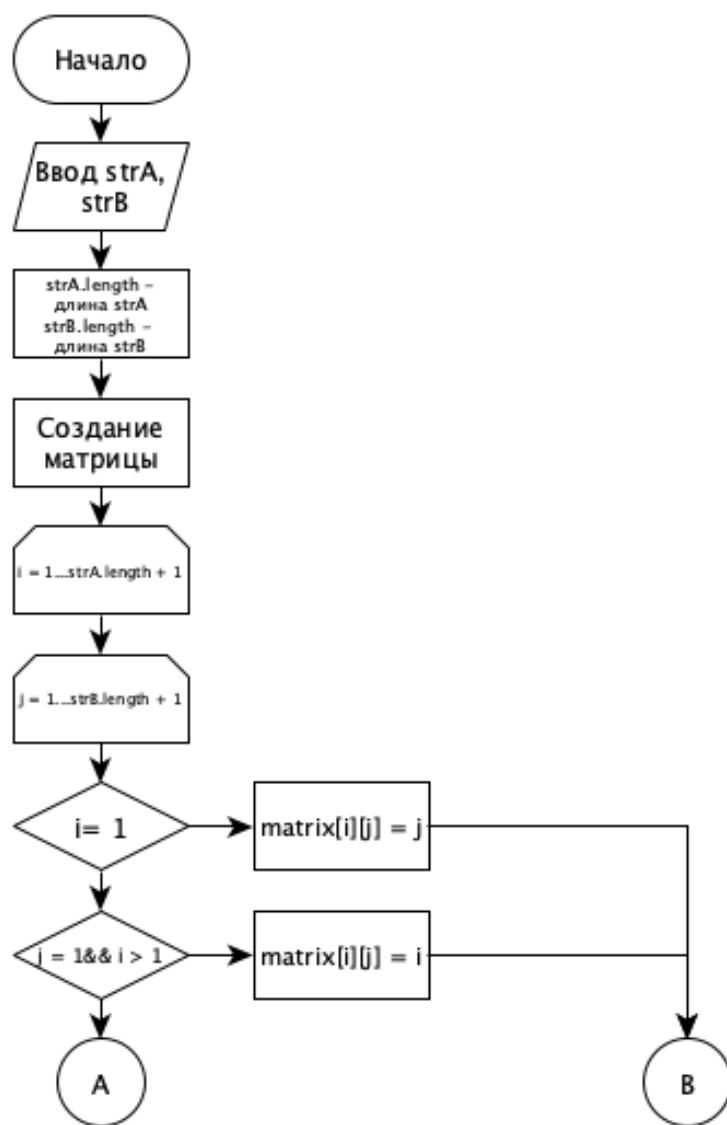


Рис. 2.3: Схема матричного алгоритма нахождения расстояния Левенштейна



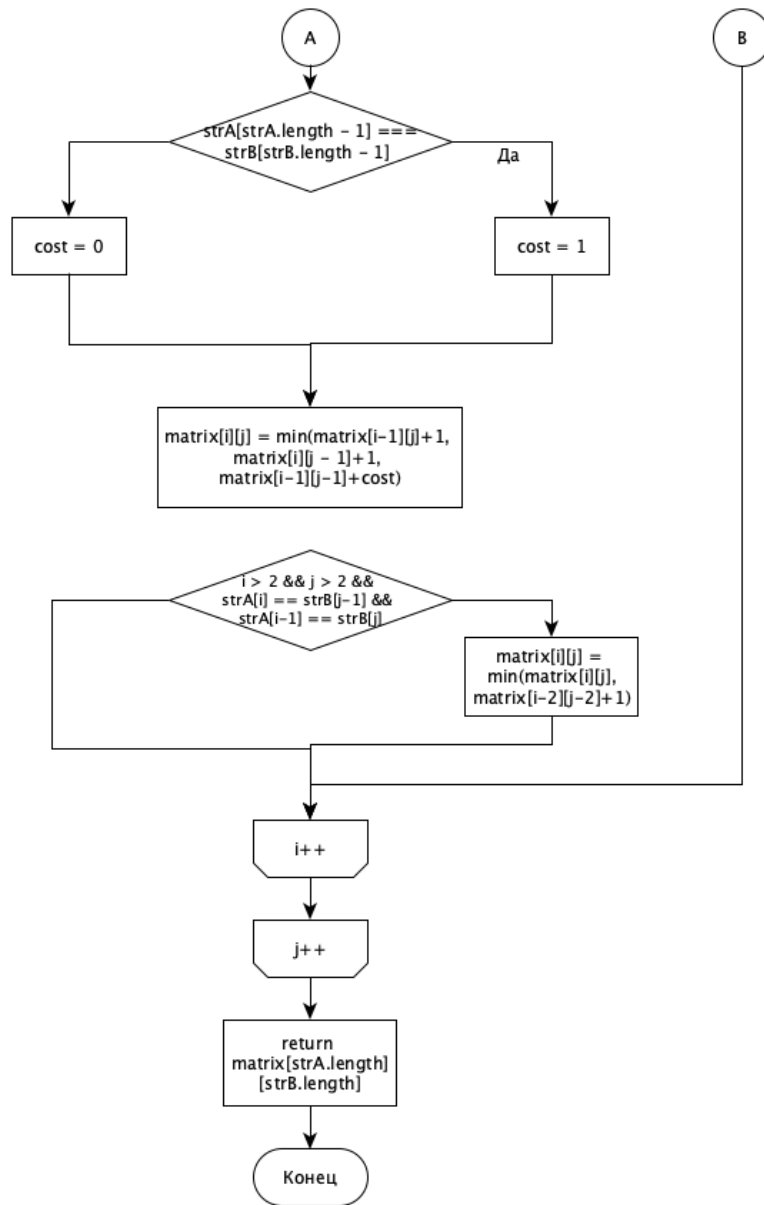


Рис. 2.4: Схема матричного алгоритма нахождения расстояния Дameraу-Левенштейна

3. Технологическая часть

3.1 Выбор ЯП

Для реализации программы был выбран язык программирования JavaScript [2] в связи с потребностью практики разработки на нем. Среда разработки - VS Code [3].

3.2 Реализация алгоритма

Листинг 3.1: Функция нахождения расстояния Левенштейна рекурсивно

```
1 function recursiveLevenshtein(strA, strB) {  
2     if (strA.length === 0 || strB.length === 0) {  
3         return Math.abs(strA.length - strB.length);  
4     }  
5     let cost = (strA.charAt(strA.length - 1) === strB.  
6         charAt(strB.length - 1)) ? 0 : 1;  
7     return Math.min(recursiveLevenshtein(strA,  
8         removeLastSymbol(strB)) + 1,  
9         recursiveLevenshtein(removeLastSymbol(  
10         strA), strB) + 1,  
11         recursiveLevenshtein(removeLastSymbol(  
12         strA), removeLastSymbol(strB)) +  
13         cost);  
14 }
```

Листинг 3.2: Функция удаления последнего символа в строке

```
1 function removeLastSymbol(str) {  
2     return str.slice(0, -1);  
3 }
```

Листинг 3.3: Функция нахождения расстояния Левенштейна рекурсивно с матрицей

```
1 function recursiveOptimizedLevenshtein(strA, strB, matrix)  
2 {
```

```

2      if (typeof(matrix[strA.length][strB.length]) !== '
        undefined') {
3          return matrix[strA.length][strB.length];
4      }
5      if (strA.length === 0 || strB.length === 0) {
6          matrix[strA.length][strB.length] = Math.abs(strA.
            length - strB.length);
7          return matrix[strA.length][strB.length];
8      }
9      let cost = (strA.charAt(strA.length - 1) === strB.
        charAt(strB.length - 1)) ? 0 : 1;
10     matrix[strA.length][strB.length] = Math.min(
11         recursiveOptimizedLevenshtein(strA,
            removeLastSymbol(strB), matrix) + 1,
12         recursiveOptimizedLevenshtein(removeLastSymbol(strA
            ), strB, matrix) + 1,
13         recursiveOptimizedLevenshtein(removeLastSymbol(strA
            ), removeLastSymbol(strB), matrix) + cost);
14     return matrix[strA.length][strB.length];
15 }

```

Листинг 3.4: Функция нахождения расстояния Левенштейна матрично

```

1 function matrixLevenshtein(strA, strB, matrix) {
2     for (let i = 0; i < strA.length + 1; i++) {
3         matrix[i] = [];
4         for (let j = 0; j < strB.length + 1; j++) {
5             if (i === 0) {
6                 matrix[i][j] = j;
7             } else if (j === 0 && i > 0) {
8                 matrix[i][j] = i;
9             } else {
10                let cost = (strA.charAt(i - 1) === strB.
                    charAt(j - 1)) ? 0 : 1;
11                matrix[i][j] = Math.min(matrix[i][j - 1] +
                    1,
12                                         matrix[i - 1][j] +
13                                         1,
14                                         matrix[i - 1][j -
                        1] + cost);
15            }
16        }
17    }
18 }

```

```

15     }
16   }
17   return matrix[strA.length][strB.length];
18 }

```

Листинг 3.5: Функция нахождения расстояния Дамерау-Левенштейна матрично

```

1 function matrixDamerauLevenshtein(strA, strB, matrix) {
2   for (let i = 0; i < strA.length + 1; i++) {
3     matrix[i] = [];
4     for (let j = 0; j < strB.length + 1; j++) {
5       if (i === 0 || j === 0) {
6         matrix[i][j] = Math.max(i, j);
7       } else {
8         let cost = (strA.charAt(i - 1) === strB.
9           charAt(j - 1)) ? 0 : 1;
10        matrix[i][j] = Math.min(matrix[i][j - 1] +
11          1,
12          matrix[i - 1][j] + 1,
13          matrix[i - 1][j - 1] + cost);
14        if (j > 1 && i > 1 &&
15          strA.charAt(i) === strB.charAt(j - 1)
16          &&
17          strA.charAt(i - 1) === strB.charAt(j))
18          {
19            matrix[i][j] = Math.min(matrix[i][j]
20              , matrix[i - 2][j - 2] + 1);
21          }
22        }
23      }
24    }
25    return matrix[strA.length][strB.length];
26  }

```

Листинг 3.6: Функция проверки транпозиции

```

1 func isTransposition(_ strA: String, _ strB: String, _ i:
2   Int, _ j: Int) -> Bool {
3   return (strA[strA.index(strA.startIndex, offsetBy: i - 1)
4     ] == strB[strB.index(strB.startIndex, offsetBy: j -

```

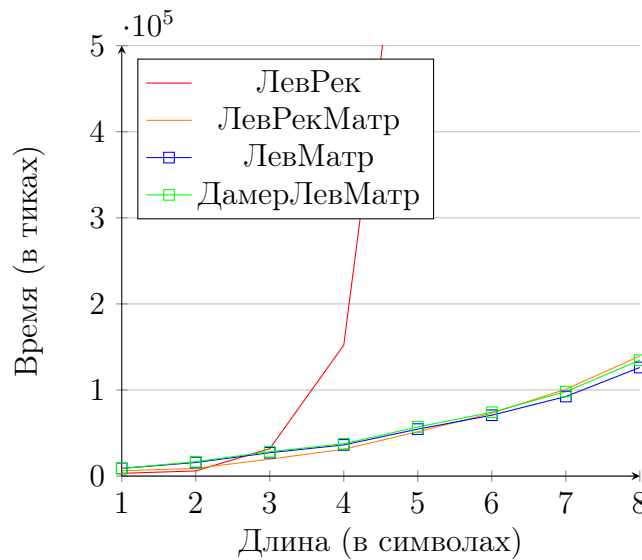


```
        2)] &&  
3    strA[strA.index(strA.startIndex, offsetBy: i - 2)] ==  
        strB[strB.index(strB.startIndex, offsetBy: j - 1)])  
4 }
```

4. Исследовательская часть

4.1 Сравнительный анализ на основе замеров времени работы алгоритмов

Был проведен замер времени работы каждого из алгоритмов на машине MacBook Pro 13 [4]. Для минимизации погрешности замеров времени каждый алгоритм исполнялся над одними и теми же строками 1000 раз и затраченное время делилось на 1000, для получения усредненного времени выполнения. Отмечу, что замеры рекурсивного алгоритма нахождения расстояния Левенштейна замерялся единожды, поскольку уже на словах, длиннее 10 символов алгоритм выполняется относительно долго.



Наиболее быстрым является алгоритм Левенштейна с матрицей; в нем требуется только $(m + 1) * (n + 1)$ операций заполнения ячейки матрицы. Рекурсивная версия алгоритма крайне сильно проигрывает матричному алгоритму уже на строках длиной 10 символов, из-за большой глубины рекурсии и многочисленных повторов. Рекурсивная версия с матрицей не теряет в своей производительности настолько сильно, поскольку в ней исключаются повторы благодаря кэшированию результатов вычисления ячеек матрицы. Заметим, что алгоритм Дамерау-Левенштейна выполняется относительно немного дольше алгоритма

Таблица 4.1: время исполнения в миллисекундах

len	Lev(R)	Lev(RM)	Lev(M)	DamLev(M)
1	0.025996	0.040448	0.022366	0.062549
2	0.048608	0.076069	0.033854	0.077381
3	0.106327	0.085223	0.028715	0.076817
4	0.502248	0.080871	0.026586	0.073056
5	0.862543	0.079230	0.030535	0.080112.
6	0.653954	0.048213	0.019119	0.046017
7	2.383116	0.074804	0.026301	0.066106
8	5.957557	0.064507	0.021543	0.002099
9	21.236497	0.084806	0.030479	0.072075
10	112.217936	0.088567	0.041352	0.088084

Левенштейна, т.к. в нем добавлена дополнительная проверка (проверка на транспозицию символов). В среднем алгоритм Дамерау-Левенштейна с матрицей работает на 10% дольше, чем матричный алгоритм Левенштейна.

4.2 Тесты

Проведем тестирование программы. В столбцах "Ожидание" и "Результат" 4 числа соответствуют рекурсивному алгоритму нахождения расстояния Левенштейна, рекурсивно-матричному алгоритму нахождения расстояния Левенштейна, матричному алгоритму нахождения расстояния Левенштейна, матричному алгоритму нахождения расстояния Дамерау-Левенштейна.

Вывод

Закодировав исследованные алгоритмы, была получена итоговая программа. Протестированная программа выдала результаты, соответствующие теоретическим данным.

Таблица 4.2: Таблица тестовых данных

№	Слово №1	Слово №2	Ожидание	Результат
1			0 0 0 0	0 0 0 0
2	1234	2143	3 3 3 2	3 3 3 2
3	amm	add	2 2 2 2	2 2 2 2
4	error	dog	4 4 4 4	4 4 4 4
5	submission		10 10 10 10	10 10 10 10
6	mochombo	0	8 8 8 8	8 8 8 8
7	dfufdfd	fddfdffdd	5 5 5 4	5 5 5 4

Заключение

В ходе лабораторной работы были изучены алгоритмы Левенштейна и Дамерау-Левенштейна - алгоритмы нахождения редакционного расстояния между строками, получены практические навыки реализации указанных алгоритмов в матричной и рекурсивных версиях. Также был изучен метод динамического программирования на материале данных алгоритмов.

Экспериментально было подтверждено различие по временным затратам рекурсивной и линейной реализаций алгоритма Левенштейна при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализации на строках различной длины.

Самым быстродейственным является алгоритм Левенштейна с матрицей. Рекурсивная версия алгоритма крайне сильно проигрывает матричному алгоритму уже на строках длиной 10 символов, из-за большой глубины рекурсии и многочисленных повторов. Рекурсивная версия с матрицей не теряет в своей производительности настолько сильно, поскольку в ней исключаются повторы благодаря кэшированию результатов вычисления ячеек матрицы. Линейный алгоритм Дамерау-Левенштейна выполняется относительно незначительно дольше линейного алгоритма Левенштейна, поскольку в нем добавлена дополнительная проверка (проверка на транспозицию символов): в среднем матричный алгоритм Левенштейна работает на 10% быстрее, чем алгоритм Дамерау-Левенштейна с матрицей.

Подытожив, матричная реализация данных алгоритмов заметно выигрывает по времени при росте длины строк, следовательно более применима в реальных проектах и задачах.

Литература

- [1] Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов. – М.: Доклады АН СССР, 1965. Т. 163. С. 845–848.
- [2] JavaScript Programming Language [Электронный ресурс]. Режим доступа: <https://developer.mozilla.org/ru/docs/Web/JavaScript/> (дата обращения: 20.09.2020).
- [3] Редактор кода VS Code. Режим доступа: <https://code.visualstudio.com/> (дата обращения: 10.10.2020).
- [4] Процессор Intel® Core™ i5 10 gen. [Электронный ресурс]. Режим доступа: <https://www.intel.ru/content/www/ru/ru/products/docs/processors/core/10th-gen-processors.html> (дата обращения: 20.09.2020).
- [5] Absolute time. Режим доступа: <https://developer.apple.com/documentation/kernel/1462446-machabsolutetime> (дата обращения: 20.09.2020).
- [6] MacOS Catalina [Электронный ресурс]. Режим доступа: <https://www.apple.com/ru/macos/catalina/> (дата обращения: 20.09.2020).
- [7] Технические характеристики ноутбука Apple MacBook Pro 13. Режим доступа: <https://support.hp.com/ru-ru/document/c05739572> (дата обращения: 10.10.2020).