



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 дисциплине "Анализ алгоритмов"

Тема Расстояние Левенштейна

Студент Романов А.В.

Группа ИУ7-53Б

Оценка (баллы) _____

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

Оглавление

Введение	2
1 Аналитическая часть	3
1.1 Рекурсивный алгоритм нахождения расстояния Левенштейна	3
1.2 Матричный алгоритм нахождения расстояния Левенштейна	4
1.3 Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы	5
1.4 Расстояния Дамерау — Левенштейна	5
1.5 Вывод	5
2 Конструкторская часть	6
2.1 Схемы алгоритмов	6
2.2 Вывод	7
3 Технологическая часть	11
3.1 Требование к ПО	11
3.2 Средства реализации	11
3.3 Реализация алгоритмов	11
3.4 Тестовые данные	13
3.5 Вывод	13
4 Исследовательская часть	14
4.1 Пример работы	14
4.2 Технические характеристики	15
4.3 Время выполнения алгоритмов	15
4.4 Использование памяти	15
4.5 Вывод	16
Заключение	17
Литература	17

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 для:

- исправления ошибок в слове
- сравнения текстовых файлов утилитой diff
- в биоинформатике для сравнения генов, хромосом и белков

Цели данной лабораторной работы:

1. Изучение метода динамического программирования на материале алгоритмов нахождения расстояния Левенштейна и Дamerau-Левенштейна.
2. Оценка реализаций алгоритмов нахождения расстояния Левенштейна и Дamerau-Левенштейна.

Задачи данной лабораторной работы:

1. Изучение алгоритмов Левенштейна и Дamerau-Левенштейна;
2. Применение метода динамического программирования для матричной реализации указанных алгоритмов;
3. Получение практических навыков реализации указанных алгоритмов: матричные и рекурсивные 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;
6. Описание и обоснование полученных результатов в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

1 | Аналитическая часть

Расстояние Левенштейна [1] между двумя строками — это минимальное количество операций вставки, удаления и замены, необходимых для превращения одной строки в другую.

Цены операций могут зависеть от вида операции (вставка (insert), удаление (delete), замена (replace)) и/или от участвующих в ней символов, отражая разную вероятность разных ошибок при вводе текста, и т. п. В общем случае:

- $w(a, b)$ — цена замены символа a на символ b .
- $w(\lambda, b)$ — цена вставки символа b .
- $w(a, \lambda)$ — цена удаления символа a .

Для решения задачи о редакционном расстоянии необходимо найти последовательность замен, минимизирующую суммарную цену. Расстояние Левенштейна является частным случаем этой задачи при

- $w(a, a) = 0$.
- $w(a, b) = 1, a \neq b$.
- $w(\lambda, b) = 1$.
- $w(a, \lambda) = 1$.

1.1 Рекурсивный алгоритм нахождения расстояния Левенштейна

Расстояние Левенштейна между двумя строками a и b может быть вычислено по формуле 1.1, где $|a|$ означает длину строки a ; $a[i]$ — i -ый символ строки a , функция $D(i, j)$ определена как:

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

а функция 1.2 определена как:

$$m(a, b) = \begin{cases} 0 & \text{если } a = b, \\ 1 & \text{иначе} \end{cases}. \quad (1.2)$$

Рекурсивный алгоритм реализует формулу 1.1. Функция D составлена из следующих соображений:

1. Для перевода из пустой строки в пустую требуется ноль операций;
2. Для перевода из пустой строки в строку a требуется $|a|$ операций;
3. Для перевода из строки a в пустую требуется $|a|$ операций;
4. 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:
 - (а) Сумма цены преобразования строки a в b и цены проведения операции удаления, которая необходима для преобразования a' в a ;
 - (б) Сумма цены преобразования строки a в b и цены проведения операции вставки, которая необходима для преобразования b' в b ;
 - (с) Сумма цены преобразования из a' в b' и операции замены, предполагая, что a и b оканчиваются разные символы;
 - (д) 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.

1.2 Матричный алгоритм нахождения расстояния Левенштейна

Прямая реализация формулы 1.1 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 $D(i, j)$.

1.3 Рекурсивный алгоритм нахождения 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.

1.4 Расстояния Дамерау — Левенштейна

Расстояние Дамерау — Левенштейна может быть найдено по формуле 1.3, которая задана как

$$d_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ & \\ \quad d_{a,b}(i, j - 1) + 1, & \\ \quad d_{a,b}(i - 1, j) + 1, & \\ \quad d_{a,b}(i - 1, j - 1) + m(a[i], b[j]), & \text{иначе} \\ \quad \left[\begin{array}{ll} d_{a,b}(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a[i] = b[j - 1]; \\ & b[j] = a[i - 1] \\ & \infty, \quad \text{иначе} \end{array} \right. & \\ \} & \end{cases}, \quad (1.3)$$

Формула выводится по тем же соображениям, что и формула (1.1). Как и в случае с рекурсивным методом, прямое применение этой формулы неэффективно по времени исполнения, то аналогично методу из 1.3 производится добавление матрицы для хранения промежуточных значений рекурсивной формулы.

1.5 Вывод

В данном разделе были рассмотрены алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна, который является модификаций первого, учитывающего возможность перестановки соседних символов. Формулы Левенштейна и Дамерау — Левенштейна для расчета расстояния между строками задаются рекурсивно, а следовательно, алгоритмы могут быть реализованы рекурсивно или итерационно.

2 | Конструкторская часть

2.1 Схемы алгоритмов

В данной части будут рассмотрены схемы алгоритмов нахождения расстояние Левенштейна и Дамерау - Левенштейна. На рисунках 2.1 - 2.4 представлены рассматриваемые алгоритмы.

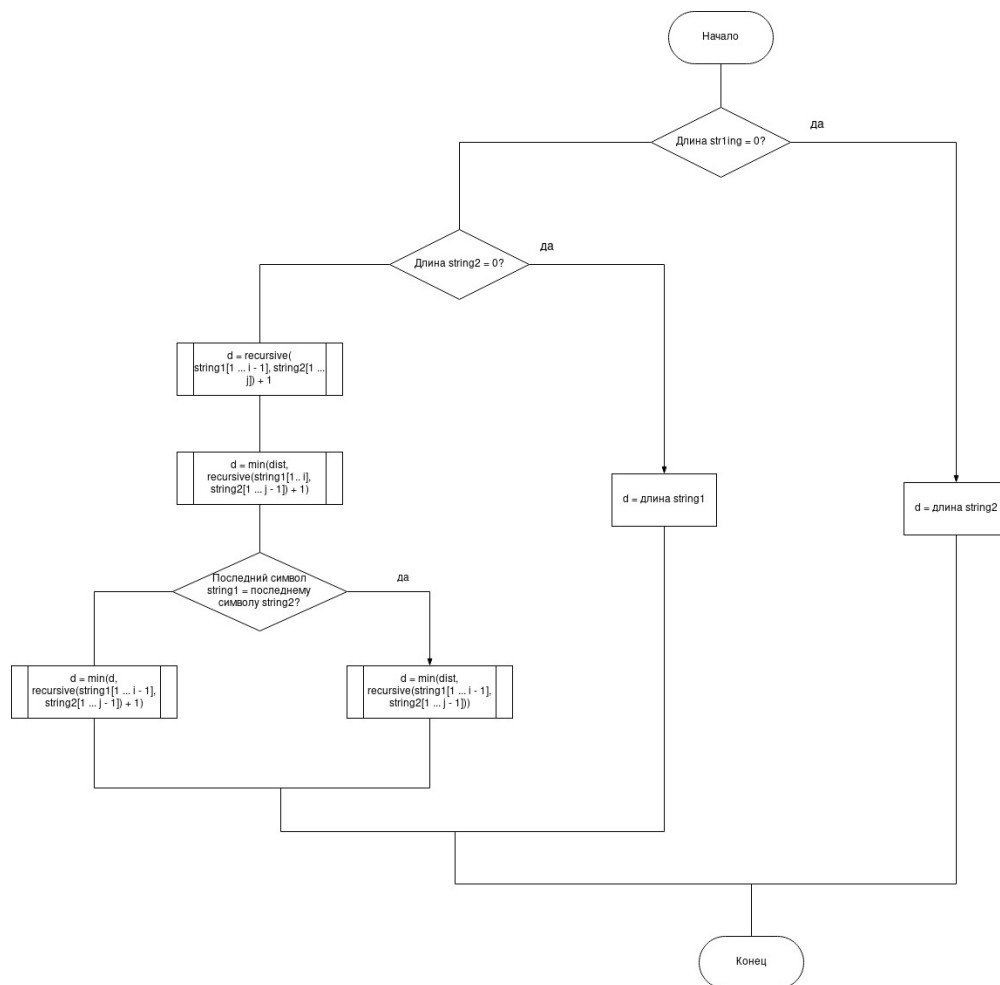


Рис. 2.1: Схема рекурсивного алгоритма нахождения расстояния Левенштейна

2.2 Вывод

На основе теоретических данных, полученные в аналитическом разделе были построены схемы исследуемых алгоритмов.

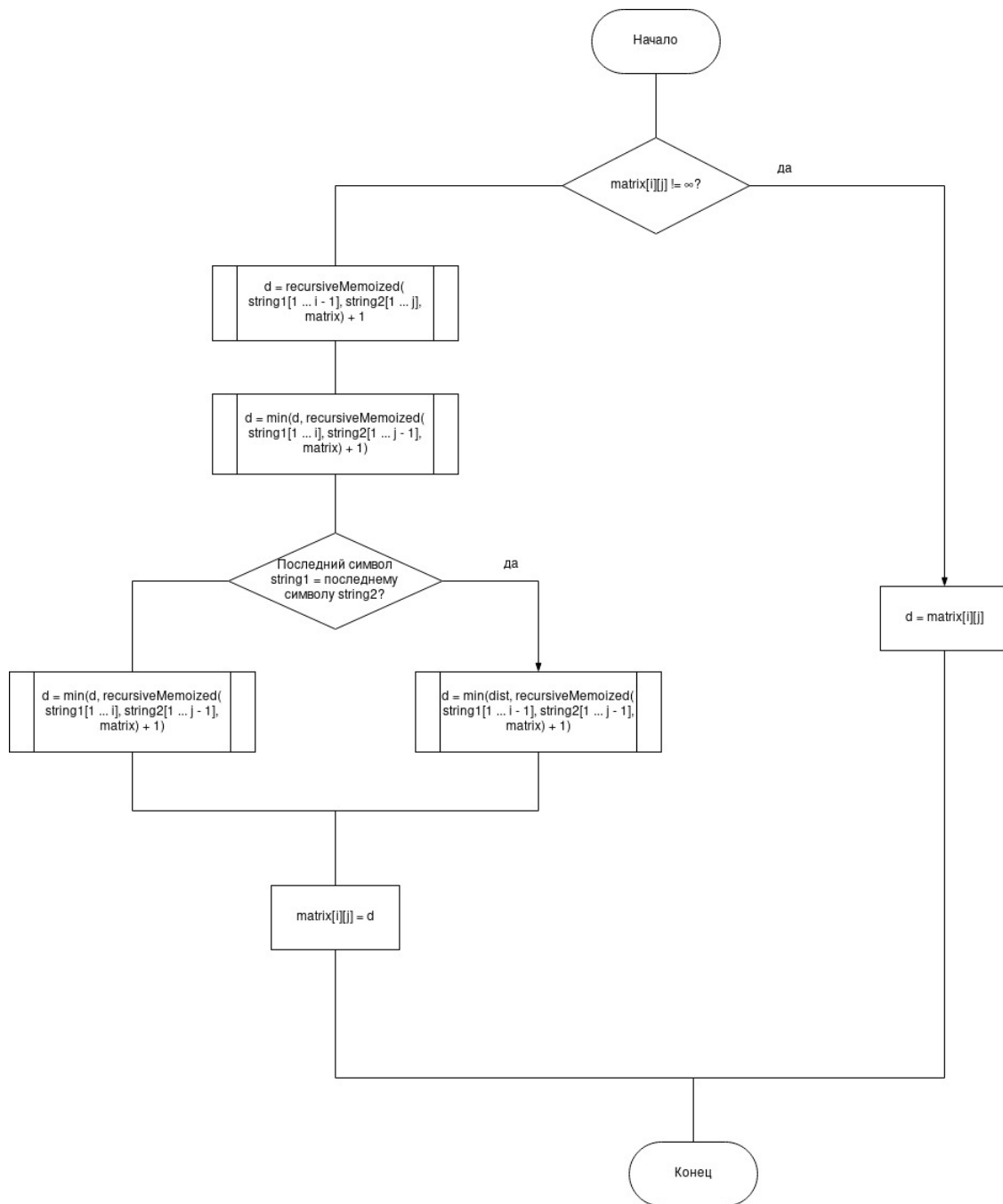


Рис. 2.2: Схема рекурсивного алгоритма с мемоизацией нахождения расстояния Левенштейна

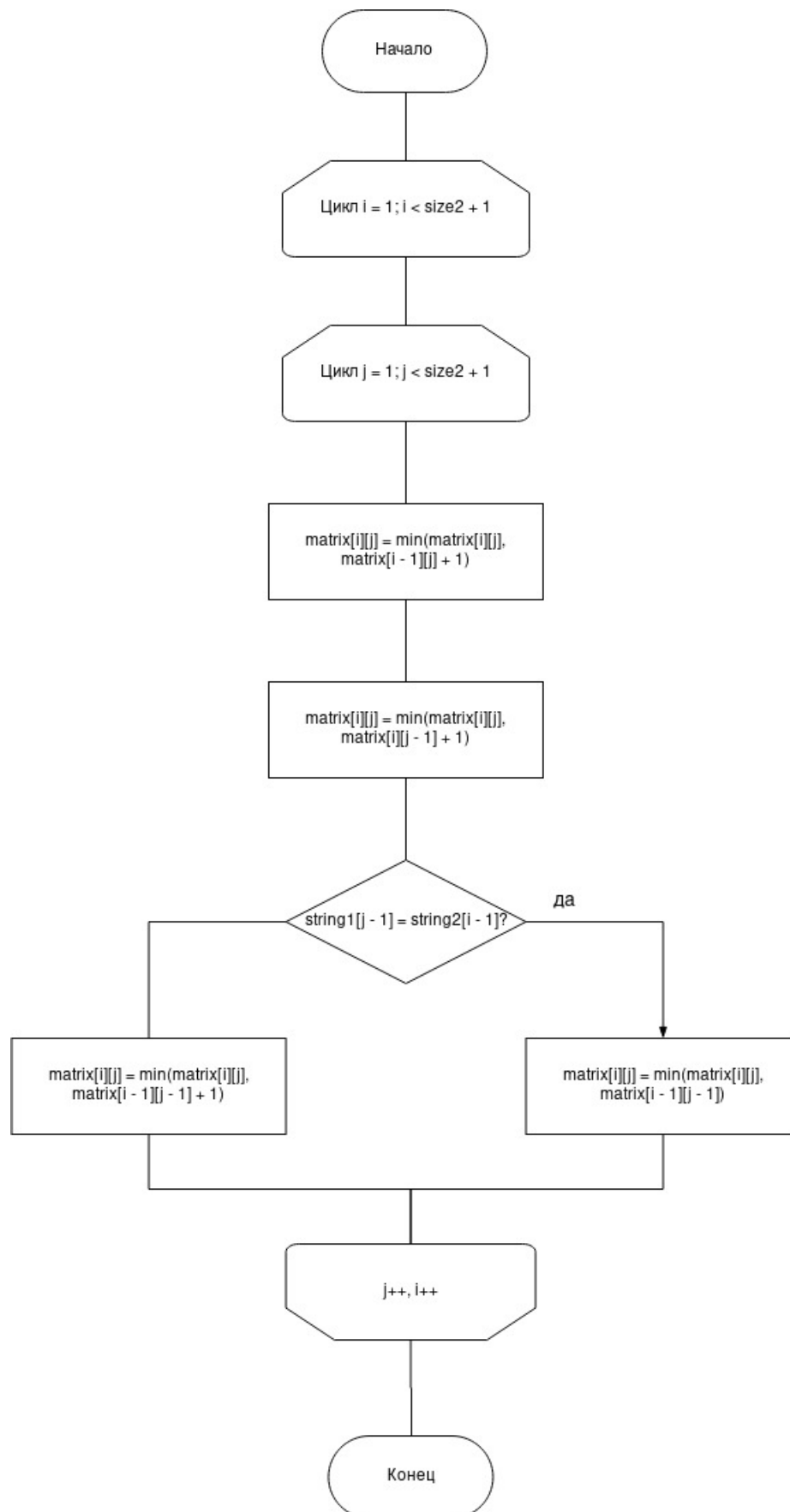


Рис. 2.3: Схема итеративного алгоритма нахождения расстояния Левенштейна

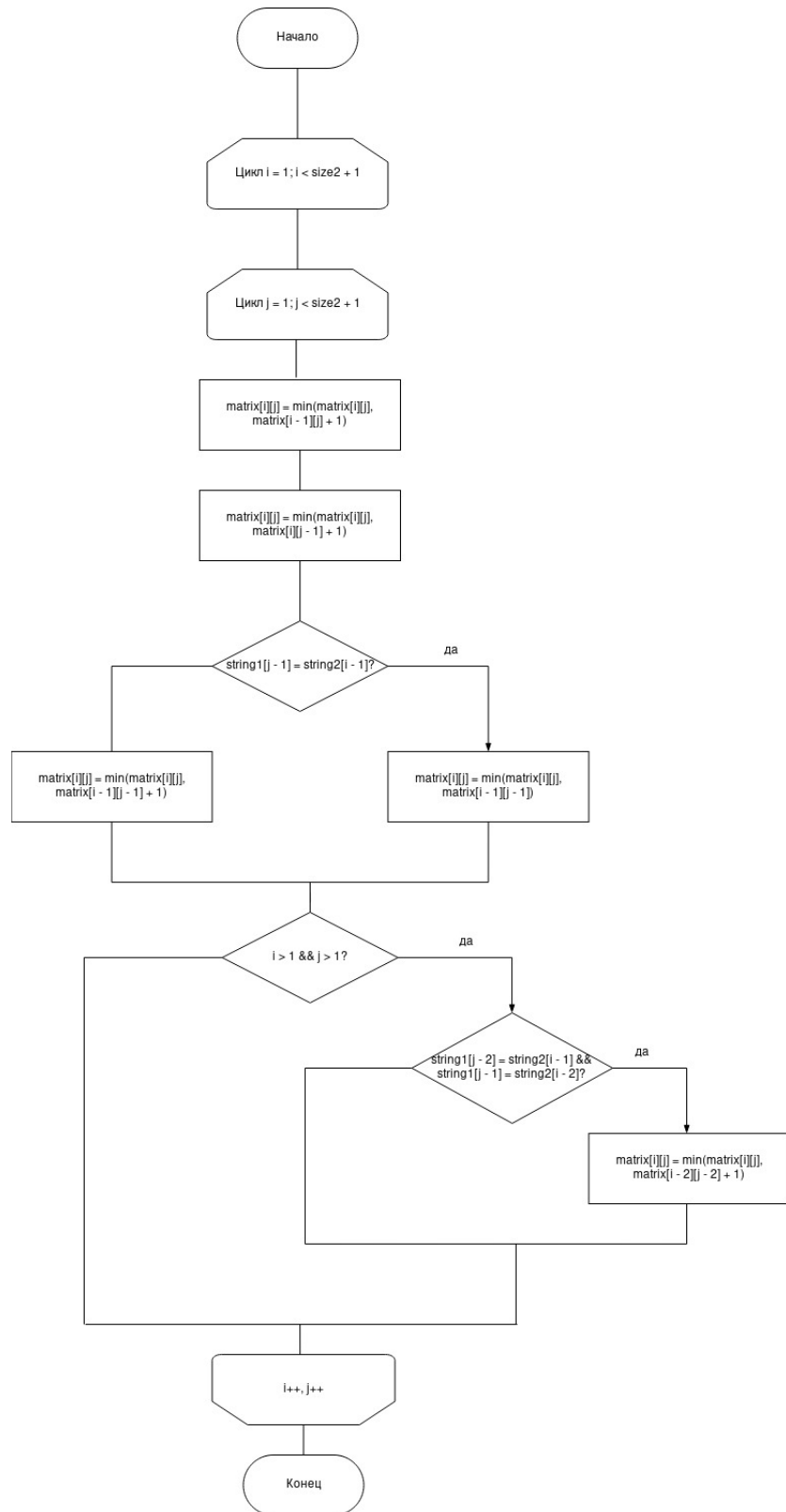


Рис. 2.4: Схема итеративного алгоритма нахождения расстояния Дameraу-Левенштейна

3 | Технологическая часть

3.1 Требование к ПО

Требования к вводу:

1. На вход подаются две строки в любой раскладке (в том числе и пустые);
2. ПО должно выводить полученное расстояние и вспомогательны матрицы;
3. ПО должно выводить потраченную память и время;

3.2 Средства реализации

Для реализации программы нахождения расстояние Левенштейна я выбрал язык программирования Haskell [2]. Данный выбор обусловлен моим желанием расширить свои знания в области применения данного языка программирования.

3.3 Реализация алгоритмов

В листингах 3.1 - 3.4 приведена реализация алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

Листинг 3.1: Функция нахождения расстояния Левенштейна рекурсивно

```
1 levenshteinRecursion :: String -> String -> (Distance, Depth)
2 levenshteinRecursion s1 s2 = _recursion s1 s2 0
3   where _recursion s1 "" n = (length s1, n)
4         _recursion "" s2 n = (length s2, n)
5         _recursion s1 s2 n = (score, depth) where
6           (insert, curr1) = _recursion (init s1) s2 (n + 1)
7           (delete, curr2) = _recursion s1 (init s2) (n + 1)
8           (replace, curr3) = _recursion (init s1) (init s2) (n + 1)
9
10        match = if last s1 == last s2 then 0 else 1
11        score = min3 (insert + 1) (delete + 1) (replace + match)
12        depth = max3 curr1 curr2 curr3
```

Листинг 3.2: Функция нахождения расстояние Левенштейна рекурсивно с мемоизацией

```
1 levenshteinMemoized :: String -> String -> (Distance, Depth)
```

```

2 levenshteinMemoized s1 s2 = _memoized s1 s2 matrix 0
3   where matrix = fromList (length s1 + 1) (length s2 + 1) $ repeat (-1)
4     _memoized s1 "" _ n = (length s1, n)
5     _memoized "" s2 _ n = (length s2, n)
6     _memoized s1 s2 mtr n = (score, depth) where
7       score = min3 (insert + 1) (delete + 1) (replace + match)
8       memoized = (getElem (length s1) (length s2) mtr, n + 1)
9       new_mtr = setElem score (length s1, length s2) mtr
10
11     (insert, curr1) = if fst memoized == -1 then _memoized (init s1) s2
12       new_mtr (n + 1)
13     else memoized
14     (delete, curr2) = if fst memoized == -1 then _memoized s1 (init s2)
15       new_mtr (n + 1)
16     else memoized
17     (replace, curr3) = if fst memoized == -1 then _memoized (init s1) (init
18       s2) new_mtr (n + 1)
19     else memoized
20
21     match = if last s1 == last s2 then 0 else 1
22     depth = max3 curr1 curr2 curr3

```

Листинг 3.3: Функция нахождения расстояния Левенштейна итеративно

```

1 levenshteinIterative :: String -> String -> Matrix Int
2 levenshteinIterative s1 s2 = fromLists $ reverse $ foldl
3   (\mtr i -> if head mtr == [] then [[0..length s1]] else calcRow mtr i :
4     mtr) [[]] [0..length s2]
5   where calcRow mtr i = foldl (\row j ->
6     row ++ if length row == 0 then [length mtr] else [
7       min3 (last row + 1) (head mtr !! j + 1) $ head mtr !! (j - 1) +
8       if s1 !! (j - 1) == s2 !! (i - 1) then 0 else 1]
9     ) [] [0..length s1]

```

Листинг 3.4: Функция нахождения расстояния Дамерау-Левенштейна матрично

```

1 damerauLevenshtein :: String -> String -> Matrix Int
2 damerauLevenshtein s1 s2 = fromLists $ reverse $ foldl
3   (\mtr i -> if head mtr == [] then [[0..length s1]] else calcRow mtr i :
4     mtr) [[]] [0..length s2]
5   where cell mtr row i j = min3 (last row + 1) (head mtr !! j + 1) $ head
6     mtr !! (j - 1) +
7     if s1 !! (j - 1) == s2 !! (i - 1) then 0 else 1
8     transposition i j = i > 1 && j > 1 && s1 !! (j - 1) == s2 !! (i - 2)
9     && s1 !! (j - 2) == s2 !! (i - 1)
10    calcRow mtr i = foldl (\row j -> row ++ [
11      if length row == 0 then length mtr
12      else if transposition i j then min (cell mtr row i j) (((head $ tail
13        mtr) !! i - 2) + 1)
14      else cell mtr row i j]
15    ) [] [0..length s1]

```

3.4 Тестовые данные

В таблице 3.1 приведены тестовые данные, на которых было протестированно разработанное ПО.

Таблица 3.1: Таблица тестовых данных

№	Первое слово	Второе слово	Ожидаемый результат	Полученный результат
1			0 0 0 0	0 0 0 0
2	kot	skat	2 2 2 2	2 2 2 2
3	kate	ktae	2 2 1 1	2 2 1 1
4	abacaba	aabcaab	4 4 2 2	4 4 2 2
5	sobaka	sboku	3 3 3 3	3 3 3 3
6	qwerty	queue	4 4 4 4	4 4 4 4
7	apple	aplpe	2 2 1 1	2 2 1 1
8		cat	3 3 3 3	3 3 3 3
9	parallels		9 9 9 9	9 9 9 9
10	русскиебуквы	русскиебукыв	2 2 2 1	2 2 2 1

3.5 Вывод

В данном разделе были разработаны исходные коды четырех алгоритмов: вычисления расстояния Левенштейна рекурсивно, с заполнением матрицы и рекурсивно с заполнением матрицы, а также вычисления расстояния Дameraу — Левенштейна с заполнением матрицы.

4 | Исследовательская часть

4.1 Пример работы

Демонстрация работы программы приведена на рисунке 4.1.

```
Введите первую строку:
string
Введите вторую строку:
gnirts
1. Расстояние Левенштейна: рекурсивный алгоритм
2. Расстояние Левенштейна: рекурсивный алгоритм с мемоизацией
3. Расстояние Левенштейна: итеративный алгоритм
4. Расстояние Дамерау-Левенштейна: итеративный алгоритм
3
[ 0 1 2 3 4 5 6 ]
[ 1 1 2 3 4 5 5 ]
[ 2 2 2 3 4 4 5 ]
[ 3 3 3 3 3 4 5 ]
[ 4 4 4 3 4 4 5 ]
[ 5 5 4 4 4 5 5 ]
[ 6 5 5 5 5 5 6 ]

Дистанция: 6
Время: 825 (наносек)
Глубина: 0
Количество задействованной памяти: 464 (байт)
Введите первую строку:
test1
Введите вторую строку:
test2
1. Расстояние Левенштейна: рекурсивный алгоритм
2. Расстояние Левенштейна: рекурсивный алгоритм с мемоизацией
3. Расстояние Левенштейна: итеративный алгоритм
4. Расстояние Дамерау-Левенштейна: итеративный алгоритм
1

Дистанция: 1
Время: 401 (наносек)
Глубина: 9
Количество задействованной памяти: 576 (байт)
```

Рис. 4.1: Работа алгоритмов нахождения расстояния Левенштейна и Дамерау – Левенштейна.

4.2 Технические характеристики

Ниже приведенные технические характеристики устройства, на котором было проведено тестирование ПО:

- Операционная система: Debian [3] Linux [4] 11 «bullseye» 64-bit.
- Оперативная память: 12 GB.
- Процессор: Intel(R) Core(TM) i5-3550 CPU @ 3.30GHz [5].

4.3 Время выполнения алгоритмов

Время выполнения алгоритм замерялось с помощью применения технологии профайлинга [6]. Данный инструмент даёт детальное описание количества вызовов и количества времени CPU, занятого каждой функцией.

В таблице 4.1. представлены замеры времени работы для каждого из алгоритмов.

Таблица 4.1: Таблица времени выполнения алгоритмов (в наносекундах)

Длина строк	Rec	MatRec	Iter	DamIter
10	30500400	1300	650	680
20	NaN	5100	2300	2500
30	NaN	11000	4800	5200
50	NaN	30000	12500	13500
100	NaN	115000	48000	55000
200	NaN	530000	249000	528000

4.4 Использование памяти

Алгоритмы нахождения расстояния Левенштейна и Дамерау — Левенштейна не отличаются друг от друга с точки зрения использования памяти.

Максимальная глубина стека вызовов при рекурсивной реализации равна сумме длин входящих строк. Поэтому, максимальный расход памяти равен:

$$(\mathcal{S}(STR_1) + \mathcal{S}(STR_2)) \cdot (2 \cdot \mathcal{S}(\text{string}) + 3 \cdot \mathcal{S}(\text{integer})), \quad (4.1)$$

где \mathcal{S} — оператор вычисления размера, STR_1 , STR_2 — строки, string — строковый тип, integer — целочисленный тип.

Использование памяти при итеративной реализации теоретически равно:

$$(\mathcal{S}(STR_1) + 1) \cdot (\mathcal{S}(STR_2) + 1) \cdot \mathcal{S}(\text{integer}) + 5 \cdot \mathcal{S}(\text{integer}) + 2 \cdot \mathcal{S}(\text{string}). \quad (4.2)$$

4.5 Вывод

Рекурсивная реализация алгоритма нахождения расстояния Левенштейна работает дольше итеративных реализаций, время работы этой реализации увеличивается в геометрической прогрессии. Рекурсивный алгоритм с мемоизацией проигрывает по памяти: так как в языке Haskell каждый объект является неизменяемым, каждый рекурсивный вызов создается новая вспомогательная матрица, хранящая данные о прошлых вызовах. Стоит отметить, что, скорее всего, в других языках программирования, где присутствует передача по ссылке такого бы не наблюдалось.

Заключение

В ходе проделанной работы был изучен метод динамического программирования на материале реализации алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна. Также были изучены алгоритмы поиска расстояния Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками и получены практические навыки реализации указанных алгоритмов в матричной и рекурсивных версиях, а так же в версиях с мемоизацией.

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. Рекурсивная реализация алгоритма Левенштейна проигрывает нерекурсивной по времени исполнения в несколько десятков раз. Так же стоит отметить, что итеративный алгоритм Левенштейна выполняется немного быстрее, чем итеративный алгоритм Дамерау - Левенштейна, но в целом алгоритмы выполняются за примерно одинаковое время.

Теоретически было рассчитано использования памяти в каждой из реализаций алгоритмов нахождения расстояния Левенштейна и Дамерау - Левенштейна. Рекурсивный алгоритм с мемоизацией в несколько десятков раз больше памяти, чем итеративная реализация алгоритма нахождения расстояния Левенштейна, из-за рекурсивного копирования вспомогательной матрицы. Этого возможно избежать используя специальные конструкции языка Haskell, но в моей реализации алгоритма с мемоизацией я не использовал данные возможности языка.

Литература

- [1] Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов. Доклады АН СССР, 1965. Т. 163. С. 845–848.
- [2] The Haskell purely functional programming language [Электронный ресурс]. Режим доступа: <https://haskell.org/>. Дата обращения: 16.09.2020.
- [3] Debian – универсальная операционная система [Электронный ресурс]. Режим доступа: <https://www.debian.org/>. Дата обращения: 20.09.2020.
- [4] Linux – Getting Started [Электронный ресурс]. Режим доступа: <https://linux.org>. Дата обращения: 20.09.2020.
- [5] Процессор Intel® Core™ i5-3550 [Электронный ресурс]. Режим доступа: <https://ark.intel.com/content/www/ru/ru/ark/products/65516/intel-core-i5-3550-processor-6m-cache-up-to-3-70-ghz.html>. Дата обращения: 20.09.2020.
- [6] Profiling – Haskell profiling [Электронный ресурс]. Режим доступа: https://downloads.haskell.org/ghc/8.8.1/docs/html/user_guide/index.html. Дата обращения: 20.09.2020.