

МИНОБРНАУКИ РОССИИ
САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ
ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
«ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)
Кафедра МО ЭВМ

ОТЧЕТ
по лабораторной работе №4
по дисциплине «Построение и анализ алгоритмов»
Тема: Кнут-Моррис-Пратт

Студент гр. 3343

Никишин С.А.

Преподаватель

Жангиров Т.Р.

Санкт-Петербург

2025

Цель работы.

Изучить и реализовать алгоритм Кнута-Морриса-Пратта.

Задание 1.

Реализуйте алгоритм КМП и с его помощью для заданных шаблона P ($|P| \leq 15000$) и текста T ($|T| \leq 5000000$) найдите все вхождения P в T .

Вход:

Первая строка - P

Вторая строка - T

Выход:

индексы начал вхождений P в T , разделенных запятой, если P не входит в T , то вывести -1

Sample Input:

ab

abab

Sample Output:

0, 2

Задание 2.

Заданы две строки A ($|A| \leq 5000000$) и B ($|B| \leq 5000000$).

Определить, является ли A циклическим сдвигом B (это значит, что A и B имеют одинаковую длину и A состоит из суффикса B , склеенного с префиксом B). Например, defabc является циклическим сдвигом abcdef.

Вход:

Первая строка - A

Вторая строка - B

Выход:

Если A является циклическим сдвигом B , индекс начала строки B в A , иначе вывести -1 .
Если возможно несколько сдвигов вывести первый индекс.

Sample Input:

defabc

abcdef

Sample Output:

3

Sample Output:

3

Описание алгоритма.

Алгоритм Кнута — Морриса — Пратта (КМП) — это эффективный алгоритм поиска подстроки в строке. Основное преимущество алгоритма КМП заключается в том, что он позволяет избежать повторной обработки уже проверенных символов, что делает его значительно быстрее, чем наивный метод поиска подстроки, особенно для больших текстов.

Алгоритм реализован в функции *kmpAlgorithm(pattern: str, text: str)*. Основные этапы алгоритма при решении задач 1 и 2 одинаковы. Отличались лишь входные и выходные данные.

Основные этапы алгоритма:

1. Обработка образца с помощью префикс функции:

Функция *computePrefixFunction(string: str)* вычисляет префикс-функцию для образца. Т.е мы получаем массив, где для каждого символа образца хранится длина наибольшего собственного префикса, который одновременно является суффиксом для подстроки, заканчивающейся на этом символе.

Принцип работы префикс-функции *computePrefixFunction(string: str)* описан в разделе “Описание функций и структур данных”.

2. Поиск образца в тексте:

Инициализируется массив *positions* для хранения индексов начала вхождений образца в тексте (используется при решении задания 1). Переменная *k* хранит длину текущего совпавшего префикса.

1. Проходим по тексту, начиная с первого символа.

2.1. Если символы не совпадают ($\text{pattern}[k] \neq \text{text}[i]$), уменьшаем *k* до значения $\text{pi}[k - 1]$.

2.2. Если символы совпадают ($\text{pattern}[k] == \text{text}[i]$), увеличиваем *k* на 1.

2.3. Если k становится равным длине образца ($k == patternLength$), значит, найдено вхождение.

2.3.1 Если рассматривается задание 1, записываем индекс начала вхождения в массив `positions` и сдвигаем образец, используя префикс-функцию ($k = pi[k - 1]$). После окончания цикла, возвращаем массив `positions`.

2.3.2 Если рассматривается задание 2, просто возвращаем первый индекс начала вхождения.

Отличия при решении задач 1 и 2 незначительные:

1. Меняется вывод (массив индексов вхождений для 1 задачи и первый индекс вхождения для 2 задачи),

2. Для задачи 2 устанавливается условие, что образец и текст одинаковы по длине, и текст увеличивается в 2 раза для определения наличия циклического сдвига.

Алгоритм имеет сложность по памяти $O(m + n)$ (Худший случай, если текст состоит из повторяющегося образца)

Функции имеет сложность по времени:

Функция	Сложность
<code>computePrefixFunction(string: str)</code>	$O(m)$
<code>kmpAlgorithm(pattern: str, text: str)</code>	$O(m + n)$

,где m – длина образца, n – длина текста

Исходные коды обеих программ указаны в приложении (приложение 1 и приложение 2 соответственно)

Описание функций и структур данных.

Структуры данных:

1. pi – список значений префикс-функции

Функции:

1. *def computePrefixFunction(string: str):* функция, вычисляющая префикс функцию.

Принцип работы:

1. Инициализируется нулями массив pi длиной, равной длине образца.
2. Проходим по образцу, начиная со второго символа.
- 3.1. Если символы не совпадают ($string[k] \neq string[i]$), уменьшаем k до значения $pi[k - 1]$.
- 3.2. Если символы совпадают ($string[k] == string[i]$), увеличиваем k на 1.
4. Записываем значение k в $pi[i]$.
5. В конце получаем массив значений, вычисленных префикс-функцией.

Параметры:

string – строка, по которой вычисляется префикс функция (тип *str*).

Возвращает:

pi — массив значений префикс-функции (тип *list*).

Сложность по времени: $O(m)$

2. *def kmpAlgorithm(pattern: str, text: str):* функция, реализующая алгоритм Кнута-Морриса-Пратта для поиска всех вхождений образца в тексте.

Параметры:

pattern – строка образца (тип *str*).

text – строка текста (тип *str*).

Возвращает:

1. В задаче 1, возвращает массив индексов начал вхождений (тип *list*).
2. В задаче 2, возвращает первый индекс начала вхождения (тип *int*).

Сложность по времени: $O(m + n)$

,где m – длина образца, n – длина текста

Тестирование.

Тестирование программ представлено на таблицах.

Тестирование задание 1

№ Теста	Входные данные	Выходные данные	Комментарий
1	ab abab	0,2	Верно
2	faf abba	-1	Верно
3	hah ahahahahahah	1,3,5,7,9	Верно

Тестирование задание 2

№ Теста	Входные данные	Выходные данные	Комментарий
1	defabc abcdef	3	Верно
2	abc abc	0	Верно
3	abc bac	-1	Верно

Примеры работ программ с выводом промежуточных результатов при введённом значении указаны в приложении (приложение 3 и приложение 4 для программ 1 и 2 соответственно)

Выводы.

Был реализован алгоритм Кнута-Морриса-Пратта (КМП). Реализованный алгоритм был использован для задач на поиск подстроки в строке и на проверку строк на циклический сдвиг. Полученные при тестировании ответы подтвердили корректность работы алгоритма.

ПРИЛОЖЕНИЕ

ПРИЛОЖЕНИЕ 1

```
IMPORT SYS

DEF COMPUTE_PREFIX_FUNCTION(PATTERN):
    """Вычисляет префикс-функцию для шаблона"""
    M = LEN(PATTERN)
    PREFIX = [0] * M
    K = 0

    FOR Q IN RANGE(1, M):
        WHILE K > 0 AND PATTERN[K] != PATTERN[Q]:
            K = PREFIX[K - 1]

        IF PATTERN[K] == PATTERN[Q]:
            K += 1

        PREFIX[Q] = K

    RETURN PREFIX

DEF KMP_SEARCH(ТЕКСТ, PATTERN):
    """Находит все вхождения шаблона в текст с помощью алгоритма КМП"""
    IF NOT PATTERN OR NOT ТЕКСТ:
        RETURN []

    N = LEN(ТЕКСТ)
    M = LEN(PATTERN)

    IF M > N:
        RETURN []

    PREFIX = COMPUTE_PREFIX_FUNCTION(PATTERN)
    OCCURRENCES = []

    J = 0 # ИНДЕКС В ШАБЛОНЕ
    FOR I IN RANGE(N): # ИНДЕКС В ТЕКСТЕ
        WHILE J > 0 AND ТЕКСТ[I] != PATTERN[J]:
            J = PREFIX[J - 1]

        IF ТЕКСТ[I] == PATTERN[J]:
            J += 1

        IF J == M:
            # НАЙДЕНО ВХОЖДЕНИЕ
            OCCURRENCES.APPEND(I - M + 1)
            J = PREFIX[J - 1] # ПРОДОЛЖАЕМ ПОИСК СЛЕДУЮЩИХ ВХОЖДЕНИЙ

    RETURN OCCURRENCES

DEF MAIN():
    PRINT("=== АЛГОРИТМ КНУТА-МОРРИСА-ПРАТТА (КМП) ===")
    PRINT("ПОИСК ВСЕХ ВХОЖДЕНИЙ ШАБЛОНА P В ТЕКСТЕ T")
    PRINT()

    PRINT("ВВЕДИТЕ ШАБЛОН P (ДЛИНА ≤ 25000):")
    PATTERN = SYS.STDIN.READLINE().RSTRIP('\n')

    PRINT("ВВЕДИТЕ ТЕКСТ T (ДЛИНА ≤ 5000000):")
    ТЕКСТ = SYS.STDIN.READLINE().RSTRIP('\n')
```

```

PRINT (F"\nПоиск шаблона '{PATTERN}' в тексте длиной {LEN(TEXT)} символов...")

# Поиск вхождений
OCCURRENCES = KMP_SEARCH(TEXT, PATTERN)

# Формирование вывода
PRINT("\nРезультат поиска:")
IF OCCURRENCES:
    PRINT(F"Найдено вхождений: {LEN(OCCURRENCES)}")
    RESULT = ', '.JOIN(MAP(STR, OCCURRENCES))
    PRINT(F"Индексы начал вхождений: {RESULT}")
ELSE:
    PRINT("Вхождений не найдено: -1")

IF __NAME__ == "__MAIN__":
    MAIN()

```

Приложение 2

```

IMPORT SYS

DEF COMPUTE_LPS(B):
    """Вычисляет префикс-функцию (LPS - LONGEST PREFIX SUFFIX) для строки B"""
    M = LEN(B)
    IF M == 0:
        RETURN []

    LPS = [0] * M
    J = 0

    FOR I IN RANGE(1, M):
        WHILE J > 0 AND B[I] != B[J]:
            J = LPS[J - 1]

        IF B[I] == B[J]:
            J += 1
            LPS[I] = J
        ELSE:
            LPS[I] = 0

    RETURN LPS

DEF KMP_SEARCH(A, B, LPS):
    """Ищет вхождение B в удвоенной строке A с помощью алгоритма КМП"""
    J = 0
    N = LEN(A)

    FOR I IN RANGE(2 * N):
        WHILE J > 0 AND A[I % N] != B[J]:
            J = LPS[J - 1]

        IF A[I % N] == B[J]:
            J += 1

        IF J == LEN(B):
            RETURN I - J + 1

    RETURN -1

DEF MAIN():
    # Надписи при вводе
    PRINT("Введите строку A:")
    A = SYS.STDIN.READLINE().RSTRIP('\n')

```



```

PRINT ("ВВЕДИТЕ СТРОКУ B:")
B = SYS.STDIN.READLINE().RSTRIP('\n')

# ПРОВЕРКА УСЛОВИЙ
IF LEN(A) != LEN(B):
    PRINT ("\nРЕЗУЛЬТАТ: -1 (СТРОКИ РАЗНОЙ ДЛИНЫ)")
    RETURN

IF LEN(A) == 0:
    PRINT ("\nРЕЗУЛЬТАТ: 0 (ПУСТЫЕ СТРОКИ)")
    RETURN

# ВЫЧИСЛЕНИЕ ПРЕФИКС-ФУНКЦИИ И ПОИСК
LPS = COMPUTE_LPS(B)
RESULT = KMP_SEARCH(A, B, LPS)

# ВЫВОД РЕЗУЛЬТАТА
IF RESULT != -1:
    PRINT (F"\nРЕЗУЛЬТАТ: {RESULT} (A является циклическим сдвигом B)")
    PRINT (F"СДВИГ: ПЕРВЫЕ {RESULT} СИМВОЛОВ B СТАНОВЯТСЯ СУФФИКСОМ A")
ELSE:
    PRINT (F"\nРЕЗУЛЬТАТ: -1 (A не является циклическим сдвигом B)")

IF __NAME__ == "__MAIN__":
    MAIN()

```