

**Министерство науки и высшего образования Российской Федерации  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«Национальный исследовательский университет ИТМО»  
(Университет ИТМО)**

**Факультет Инфокоммуникационных технологий (ИКТ)**

**Образовательная программа Мобильные и сетевые технологии**

**О Т Ч Е Т**

**по Лабораторной работе 4**

Дисциплина: Алгоритмы и структуры данных.

Тема: Алгоритмы поиска подстрок.

Специальность: 09.03.03 Прикладная информатика.

Проверил:

Мусаев А. А. \_\_\_\_\_

Дата: «\_\_» \_\_\_\_\_ 2023 г.

Оценка: \_\_\_\_\_

Выполнил:

Балцат К. И.,

студент группы К33401

Санкт-Петербург  
2023

## ВЫПОЛНЕНИЕ

### 1 Задача 1

Заполните массив 500 числами (четный вариант – простые числа, нечетный вариант – числа Фибоначчи) написанными слитно. Используя каждый изученный алгоритм поиска подстрок (наивный, Рабина-Карпа, Бойера-Мура, Кнута-Морриса-Пратта), посчитайте количество наиболее часто встречающихся двузначных чисел в образовавшейся строке. Сравните изученные алгоритмы поиска подстрок. Сделайте вывод о их достоинствах и недостатках.

Решение:

Создаем необходимый массив из слитно написанных чисел Фибоначчи. Затем определяем функцию для подсчета количества наиболее часто встречающихся двузначных чисел в строке. Функция использует регулярное выражение для поиска всех двузначных чисел в строке и словарь для подсчета количества вхождений каждого числа. Функция возвращает список кортежей, где каждый кортеж содержит двузначное число и количество его вхождений в строку. Кортежи отсортированы по убыванию количества вхождений.

```
1  import numpy as np
2  import regex as re
3
4  # Создаем массив из первых 500 чисел Фибоначчи:
5  fibonacci = [0, 1]
6  for i in range(2, 500):
7      fibonacci.append(fibonacci[i-1] + fibonacci[i-2])
8
9  # Объединяем числа в строку и заменяем пробелы на пустые символы:
10 fibonacci_str = ''.join(map(str, fibonacci))
11 fibonacci_str = fibonacci_str.replace(' ', '')
12
13 def count_numbers(string):
14     numbers = re.findall(r'\d{2}', string, overlapped=True)
15     print(len(string), len(numbers))
16     counts = {}
17     for num in numbers:
18         if num in counts:
19             counts[num] += 1
20         else:
21             counts[num] = 1
22     sorted_counts = sorted(counts.items(), key=lambda x: x[1], reverse=True)
23     return sorted_counts
24 re_counts = count_numbers(fibonacci_str)
25 print(re_counts[:10])
```

## Наивный алгоритм поиска подстроки

Определяем функцию для наивного алгоритма поиска подстроки:

```
27 # Наивный алгоритм поиска подстроки
28 def naive_search(text, pattern_len=2):
29     n = len(text)
30     count_dict = {}
31     for i in range(n - pattern_len + 1):
32         pattern = text[i:i+pattern_len]
33         if pattern.isdigit() and len(pattern) == pattern_len:
34             if pattern in count_dict:
35                 count_dict[pattern] += 1
36             else:
37                 count_dict[pattern] = 1
38     sorted_counts = sorted(count_dict.items(), key=lambda x: x[1], reverse=True)
39     return sorted_counts
40
41 naive_counts = naive_search(fibonacci_str, pattern_len=2)
42 print(naive_counts[:10])
```

## Алгоритм Рабина-Карпа

Определяем функцию для алгоритма Рабина-Карпа и внутри полиномиальную кольцевую Хеш-функцию:

```
52 # Алгоритм Рабина-Карпа
53 def rabin_karp_search(string, pattern_len=2):
54     # Задаем полиномиальную кольцевую!! hash-функцию.
55     def H(substring: str, x=10, m=1):
56         ...
57         same as int
58         ...
59         h = 0
60         for i in range(len(substring)):
61             h += int(substring[i])*x**(m-i)
62         return h
63     counts = {f'{i:02}': 0 for i in range(1, 100)}
64     hash_table = []
65     for i in range(len(string) - pattern_len + 1):
66         hash_table.append(H(string[i:i+pattern_len]))
67     for pattern in counts.keys():
68         pattern_hash = H(pattern)
69         for hash in hash_table:
70             if pattern_hash == hash:
71                 # Посимвольная проверка не требуется. В нашем случае hash-фу
72                 counts[pattern] += 1
73     sorted_counts = sorted(counts.items(), key=lambda x: x[1], reverse=True)
74     return sorted_counts
--
```

## Алгоритм Бойера-Мура

Определяем функцию для алгоритма Бойера-Мура:

```

80  # Алгоритм Бойера-Мура
81  def boyer_moore_search(string):
82      counts = {f'{i:02}': 0 for i in range(1, 100)}
83      for pattern in counts.keys():
84          pos = 1 # len(pattern)==2
85          while pos < len(string):
86              if (string[pos] == pattern[1]) is False:
87                  if string[pos] in pattern:
88                      pos += 1 # len(pattern)==2
89                  else:
90                      pos += 2
91              elif (string[pos-1] == pattern[0]) is False:
92                  pos+=1
93              else:
94                  counts[pattern] += 1
95                  pos += 1
96      sorted_counts = sorted(counts.items(), key=lambda x: x[1], reverse=True)
97      return sorted_counts

```

## Алгоритм Кнута-Морриса-Пратта

Определяем функцию для алгоритма Кнута-Морриса-Пратта и внутри префикс-функцию:

```

101  # Алгоритм Кнута-Морриса-Пратта
102  def kmp_search(string: str):
103      def prefix(s) -> list:
104          n = len(s)
105          pi = [0] * n # создаем список для хранения значений префикс-функции
106          for i in range(1, n):
107              j = pi[i-1]
108              while j > 0 and s[j] != s[i]:
109                  j = pi[j-1]
110              if s[j] == s[i]:
111                  j += 1
112              pi[i] = j
113          return pi
114
115      counts = {f'{i:02}': 0 for i in range(1, 100)}
116      for pattern in counts.keys():
117          P = prefix(pattern)
118          k = 0
119          for i in range(len(string)):
120              while k>0 and pattern[k] != string[i]:
121                  k = P[k-1]
122              if pattern[k] == string[i]:
123                  k += 1
124              if k == len(pattern):
125                  counts[pattern] += 1
126                  k = P[k-1]
127      sorted_counts = sorted(counts.items(), key=lambda x: x[1], reverse=True)
128      return sorted_counts

```

## Результат работы программы

```
konstantinbaltat@MacBook-Pro-Konstantin-3 HW4 % python3 task1.py
Время выполнения count_numbers: 0.117525 сек.
Re_counts
[('71', 296), ('05', 296), ('20', 291), ('11', 289), ('67', 289), ('81', 288), ('32', 287), ('12', 281)]
Время выполнения naive_search: 0.107492 сек.
Наивный алгоритм
[('71', 296), ('05', 296), ('20', 291), ('11', 289), ('67', 289), ('81', 288), ('32', 287), ('12', 281)]
Время выполнения rabin_karp_search: 0.938341 сек.
Алгоритм Рабина-Карпа
[('05', 296), ('71', 296), ('20', 291), ('11', 289), ('67', 289), ('81', 288), ('32', 287), ('12', 281)]
Время выполнения boyer_moore_search: 3.547805 сек.
Алгоритм Бойера-Мура
[('05', 296), ('71', 296), ('20', 291), ('11', 289), ('67', 289), ('81', 288), ('32', 287), ('12', 281)]
Время выполнения kmp_search: 4.587905 сек.
Алгоритм Кнута-Морриса-Пракса
[('05', 296), ('71', 296), ('20', 291), ('11', 289), ('67', 289), ('81', 288), ('32', 287), ('12', 281)]
```

## Краткое сравнение алгоритмов:

### Наивный алгоритм:

#### Достоинства:

- Простота реализации.
- Работает для любых строк.

#### Недостатки:

- Очень медленный для больших текстов и/или шаблонов.
- В худшем случае требует  $O(nm)$  операций.

### Алгоритм Рабина-Карпа:

#### Достоинства:

- Очень быстрый в среднем случае.
- Возможность поиска нескольких вхождений.

#### Недостатки:

- Некоторые коллизии могут приводить к ложным срабатываниям.
- Требуется хранение хеш-значений всех подстрок текста, что занимает много памяти.
- Может работать медленно в худшем случае.

### Алгоритм Бойера-Мура:

#### Достоинства:

- Очень быстрый на практике.
- Не требует дополнительной памяти за исключением массива прыжков.
- Хорошо справляется с худшим случаем.

Недостатки:

Сложен для понимания и реализации.

Не работает для некоторых типов шаблонов.

### **Алгоритм Кнута-Морриса-Пратта:**

Достоинства:

Быстрый и эффективный в среднем и худшем случаях.

Не требует дополнительной памяти.

Недостатки:

Сложен для понимания и реализации.

Не так быстр как алгоритм Бойера-Мура на практике.

Итак, выбор алгоритма зависит от конкретной задачи. Если важна производительность на больших данных, то следует использовать алгоритм Бойера-Мура или Кнута-Морриса-Пратта. Если необходимо искать несколько вхождений, то Рабина-Карп может оказаться хорошим выбором. Если же важна простота реализации и работа с короткими строками, то можно использовать наивный алгоритм.

## **2 Задача 2**

Дан набор рефератов. Выберите любой алгоритм поиска и определите количество плагиата (в % от общего количества символов в реферате) в тексте реферата, взяв за основу соответствующие статьи из Википедии (название файла = название статьи). За плагиат считать любые 3 совпавших слова, идущих подряд. Обоснуйте выбранный алгоритм поиска.

Решение:

Определим функцию обработки текста.

```

1  import os
2  import re
3  import docx2txt
4  import string
5  import requests
6  from bs4 import BeautifulSoup
7
8  # Set the directory where the files are located
9  dir_path = "./texts/"
10
11 def clean_text(text):
12     # Remove unwanted characters from the text
13     text = re.sub(r'\\[[0-9]+\\]', '', text)
14     text = re.sub(r'\\n+', '\\n', text)
15     text = re.sub(r'\\t+', '\\t', text)
16
17     # Remove punctuation except periods in lists
18     for match in re.findall(r'(?m)^d+[.])\\s+', text):
19         text = text.replace(match, '')
20     text = re.sub(r'(?m)^[^\\S\\n]*[---][^\\S\\n]+', '', text)
21     text = text.replace('.', '')
22     text = text.translate(str.maketrans(
23         '', '', string.punctuation.replace('.', '')+'.'))
24     # Remove punctuation
25     text = text.translate(str.maketrans('', '', string.punctuation))
26     # Convert to lowercase
27     text = text.lower()
28     # Remove extra whitespace
29     text = ' '.join(text.split())
30     return text
31

```

Определим функцию чтения статьи с Википедии.

```

33 # Define the function to get the Wikipedia article
34 def get_wiki_article(article_name):
35     # Replace spaces with underscores in article name
36     article_name = article_name.replace(" ", "_")
37     # Build the URL for the article
38     url = f"https://ru.wikipedia.org/wiki/{article_name}"
39     # Make a request to the URL and get the HTML content
40     response = requests.get(url)
41     html_content = response.content
42     # Parse the HTML content with BeautifulSoup
43     soup = BeautifulSoup(html_content, 'html.parser')
44     # Find the content div and extract the text
45     content_div = soup.find('div', {'class': 'mw-parser-output'})
46     # Remove unwanted elements from the content
47     [s.extract() for s in content_div('sup')]
48     [s.extract() for s in content_div('img')]
49     [s.extract() for s in content_div('audio')]
50     [s.extract() for s in content_div('video')]
51     [s.extract() for s in content_div('map')]
52     [s.extract() for s in content_div('timeline')]
53     [s.extract() for s in content_div('table')]
54     # Get the cleaned text
55     text = content_div.get_text()
56     text = clean_text(text)
57     # Remove the [править | править код] artifact
58     text = text.replace('[править | править код]', '')
59     # Return the cleaned text
60     return re.sub(r'\\[[0-9]+\\]', '', text)

```

Определим функцию подсчета плагиата.

```

62 # Define the function to calculate plagiarism
63 def calculate_plagiarism(text, wiki_text):
64     # Split the text and wiki_text into words
65     text_words = text.split()
66     # Iterate over the text_words and compare to wiki_words
67     match_count = 0
68     for i in range(len(text_words)-2):
69         text_tri = text_words[i] + " " + \
70             text_words[i+1] + " " + text_words[i+2]
71         if text_tri in wiki_text:
72             match_count += 1
73     # Calculate and return the plagiarism percentage
74     plagiarism_percentage = (match_count / len(text_words)) * 100
75     return plagiarism_percentage
76

```

Прочтем и обработаем файлы.

```

77 authors = {
78     'Корпоративные ценности': ['Фёдорова', 'Мыльченко', 'Гусев'],
79     'Научный метод': ['Никифоров', 'Ершов', 'Иванова'],
80     'Рентгеновское излучение': ['Коломиец', 'Агаев', 'Мартынюк'],
81     'Жизнь': ['Скворцов', 'Мордовцев'],
82     'Улыбка': ['Бахтина', 'Шевченко', 'Шимченко'],
83     'Астероид': ['Морозова', 'Цветкова', 'Колтунова'],
84     'Логика': ['Сидненко', 'Георгов', 'Резкаллах'],
85     'Стресс': ['Шапиро', 'морозов', 'уразалин'],
86     'Система управления базами данных': ['Шабанов', 'Абдулла', 'Гаджиева']
87 }
88
89
90 # Iterate over the files in the directory
91 for filename in os.listdir(dir_path):
92     # Only process docx and rtf files
93     if filename.endswith(".docx"):
94         # Get the file path
95         file_path = os.path.join(dir_path, filename)
96         # Extract the article name from the filename
97         article_name = os.path.splitext(filename)[0]
98         # Read the text from the file
99         text = docx2txt.process(file_path)
100         text = clean_text(text)
101         # Get the Wikipedia article for the same topic
102         wiki_text = get_wiki_article(article_name)
103         # Calculate the plagiarism percentage
104         plagiarism_percentage = calculate_plagiarism(text, wiki_text)
105         # Print the result
106         print(
107             f"Доклад: {article_name} ({filename})\nАвторы: {'', '.join(authors[article_name])}\nПро
108

```

В результате работы программы получим:



```

Доклад: Рентгеновское излучение (Рентгеновское излучение.docx)
Авторы: Коломиец, Агаев, Мартынюк
Процент плагиата: 21.20%
Доклад: Система управления базами данных (Система управления базами данных.docx)
Авторы: Шабанов, Абдулла, Гаджиева
Процент плагиата: 22.57%
Доклад: Логика (Логика.docx)
Авторы: Сидненко, Георгов, Резкаллах
Процент плагиата: 17.99%
Доклад: Стресс (Стресс.docx)
Авторы: Шапиро, морозов, уразалин
Процент плагиата: 15.07%
Доклад: Астероид (Астероид.docx)
Авторы: Морозова, Цветкова, Колтунова
Процент плагиата: 22.22%
Доклад: Улыбка (Улыбка.docx)
Авторы: Бахтина, Шевченко, Шимченко
Процент плагиата: 19.48%
Доклад: Жизнь (Жизнь.docx)
Авторы: Скворцов, Мордовцев
Процент плагиата: 6.36%
Доклад: Научный метод (Научный метод.docx)
Авторы: Никифоров, Ершов, Иванова
Процент плагиата: 24.47%
Доклад: Корпоративные ценности (Корпоративные ценности.docx)
Авторы: Фёдорова, Мыльченко, Гусев
Процент плагиата: 7.39%

```

```

1  authors = {
2      'Корпоративные ценности': ['Фёдорова', 'Мыльченко', 'Гусев'],
3      'Рентгеновское излучение': ['Коломиец', 'Агаев', 'Мартынюк'],
4      'Жизнь': ['Скворцов', 'Мордовцев'],
5      'Улыбка': ['Бахтина', 'Шевченко', 'Шимченко'],
6      'Астероид': ['Морозова', 'Цветкова', 'Колтунова'],
7      'Логика': ['Сидненко', 'Георгов', 'Резкаллах'],
8      'Стресс': ['Шапиро', 'морозов', 'уразалин'],
9      'Система управления базами данных': ['Шабанов', 'Абдулла', 'Гаджиева']
10 }
11
12 # Iterate over the files in the directory
13 for filename in os.listdir(dir_path):
14     # Only process docx and rtf files
15     if filename.endswith(".docx"):
16         # Get the file path
17         file_path = os.path.join(dir_path, filename)
18         # Extract the article name from the filename
19         article_name = os.path.splitext(filename)[0]
20         # Read the text from the file
21         text = docx2txt.process(file_path)
22         text = clean_text(text)
23         # Get the Wikipedia article for the same topic
24         wiki_text = get_wiki_article(article_name)
25         # Calculate the plagiarism percentage
26         plagiarism_percentage = calculate_plagiarism(text, wiki_text)
27         # Print the result
28         print(
29             f"Доклад: {article_name} ({filename})\nАвторы: {' '.join(authors[article_name])}\nПроцент плагиата: {plagiarism_percentage}%\n"
30         )

```

Выбранный алгоритм для вычисления плагиата - это метод сравнения на основе триграмм. Метод берет входной текст и разбивает его на триграммы (группы из трех соседних слов), а затем сравнивает эти триграммы с триграммами, извлеченными из статьи Википедии.

Этот алгоритм эффективен, поскольку он учитывает контекст, в котором слова встречаются в тексте, а не просто сравнивает отдельные слова. Сравнивая группы из трех соседних слов, алгоритм может лучше уловить смысл и структуру текста и сравнить его с исходным текстом. Кроме того, метод сравнения на основе триграмм более устойчив к незначительным изменениям в тексте, таким как изменение порядка слов или небольшие изменения в формулировках. Это объясняется тем, что триграммы способны уловить сходство между текстами даже при изменении отдельных слов.

В целом, метод сравнения на основе триграмм является хорошо зарекомендовавшим себя подходом к обнаружению плагиата и доказал свою эффективность на практике.

## ВЫВОД

Я изучил алгоритмы поиска подстрок и на практике понял достоинства и недостатки каждого. Теперь я способен эффективно применять каждый изученный алгоритм в зависимости от решаемой задачи.