

**Московский авиационный институт
(национальный исследовательский университет)**

**Институт информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Криптография»

Студент: Д. А. Ваньков
Преподаватель: А. В. Борисов
Группа: М8О-307Б
Дата:
Оценка:
Подпись:

Москва, 2020

Задача:

Сравнить

1. два осмысленных текста на естественном языке,
2. осмысленный текст и текст из случайных букв,
3. осмысленный текст и текст из случайных слов,
4. два текста из случайных букв,
5. два текста из случайных слов.

Как сравнивать: считать процент совпадения букв в сравниваемых текстах – получить дробное значение от 0 до 1 как результат деления количества совпадений на общее число букв. Расписать подробно в отчёте алгоритм сравнения и приложить сравниваемые тексты в отчёте хотя бы для одного запуска по всем пяти подпунктам. Осознать какие значения получаются в этих пяти подпунктах. Привести свои соображения о том почему так происходит. Длина сравниваемых текстов должна совпадать. Привести соображения о том какой длины текста должно быть достаточно для корректного сравнения.

1 Описание

В качестве текстов были взяты: роман Льва Толстого "Война и Мир" версия на английском языке, а также роман Маргарет Митчелл "Унесенные ветром".

Прежде чем приступить к сравнению текста нужно было предварительно было его обработать, то есть убрать все лишние знаки препинания, дефисы, и т.д. Для этого была написана специальная программа по избавлению от этих символов.

```
1 import re
2
3 out = open('clean Gone with the Wind.txt', 'w')
4
5 is_a_comment_start = False
6 is_a_comment_end = False
7 c = 0
8
9 needed = [' ']
10
11 with open('Gone with the Wind.txt', 'r', errors='ignore') as file:
12     for line in file:
13         out_line = ''
14         for symb in line:
15             if symb == '\n':
16                 out_line += ' '
17             if symb.isalpha() or (symb in needed):
18                 out_line += symb
19         out.write(out_line)
20
21
22 file.close()
23 out.close()
```

Чтобы создать текст из случайных символов или слов была написана дополнительная программа - генератор текстов. Принцип генерации текста из символов: функции подается английский алфавит из верхнего и нижнего регистров, а также размер для генерации текста определенной длины. Заполнение нового файла происходит поэтапно:

1. Генерация случайной длины слова размером от 1 до 20 символов.
2. Заполнение слова, с длиной полученной на прошлом шаге случайными символами из словаря - алфавита.
3. Запись слова в новый файл с разделителем - пробелом. Если при добавлении слова длина текста становится больше, чем исходная добавляется только его часть.

Для генерации текста из случайных слов первоначально получается словарь всех присутствующих слов. Затем, случайным образом, выбирается слово из данного словаря и добавляется в новый файл. Если длина слова не уместается, то также как и в случае с символьной генерацией добавляется только часть слова.

```

1  from random import randint
2  from random import choice
3
4
5  def generate_by_symbols(filename, vocabulary, length):
6      file = open(filename, 'w')
7
8      out_len = 0
9      while out_len < length:
10         word_size = randint(1, 20)
11         generated_word = ''
12         for _ in range(word_size):
13             symbol = choice(vocabulary)
14             generated_word += symbol
15
16         if out_len + word_size > length:
17             file.write(generated_word[0:length - out_len + 1])
18             out_len += length - out_len + 1
19         else:
20             file.write(generated_word + ' ')
21             out_len += word_size + 1
22
23     print(f'Previous size = {length}')
24     print(f'Generated size = {out_len}')
25     file.close()
26
27
28 def generate_by_words(filename, vocabulary, length):
29     file = open(filename, 'w')
30
31     out_len = 0
32
33     while out_len < length:
34         word = choice(vocabulary)
35         if out_len + len(word) > length:
36             file.write(word[:length-out_len + 1])
37             out_len += length - out_len + 1
38         else:
39             file.write(word + ' ')
40             out_len += len(word) + 1
41
42     print(f'Previous size = {length}')
43     print(f'Generated size = {out_len}')
44     file.close()
45
46
47 if __name__ == '__main__':
48     text = open('clean War and Peace.txt', 'r').read()
49     len_ = len(text)
50
51     alphabet = ['abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ']
52     set_ = []
53     for i in alphabet[0]:
54         set_.append(i)
55     set_ = sorted(set_)

```

```

56     print(f'Alphabet is:\n{set_}')
57
58     print(f'Full length of all text data is: {len_} characters')
59     vocabulary = sorted(set(text))
60     print(f'With {len(vocabulary)} unique characters')
61     print(f'Vocabulary is:\n{vocabulary}')
62
63     generate_by_symbols('text by symbols 1.txt', set_, len_)
64     generate_by_symbols('text by symbols 2.txt', set_, len_)
65
66     words = []
67
68     word = ''
69     for symb in text:
70         if symb != ' ':
71             word += symb
72         else:
73             words.append(word)
74             word = ''
75     words_set = sorted(set(words))
76     print(f'Unique words in text: {len(words_set)}')
77     print(f'Example of words\n {words_set[1000:1020]}')
78
79     generate_by_words('text by words 1.txt', words_set, len(text))
80     generate_by_words('text by words 2.txt', words_set, len(text))

```

Сравнение. Для сравнения двух текстов было написано две функции. Первая сравнивает два текста по-индексно с учетом пробелов. То есть, если символ пробела совпал с в обоих текстах, то счетчик успеха увеличивается. Вторая же функция склеивает все слова в одно и также, по-индексно сравнивает все буквы в тексте.

```

1  def compare(filename1, filename2):
2      file1 = open(filename1, 'r')
3      file2 = open(filename2, 'r')
4      text1 = file1.read()
5      text2 = file2.read()
6      size = len(text1)
7      size2 = len(text2)
8
9      if size > size2:
10         size = size2
11
12         success = 0
13         for i in range(size):
14             if text1[i] == text2[i]:
15                 success += 1
16
17         accuracy = success / size
18
19         file1.close()
20         file2.close()
21         return accuracy
22
23
24  def without_space_compare(filename1, filename2):
25      file1 = open(filename1, 'r')

```

```

26 | file2 = open(filename2, 'r')
27 | text1 = file1.read()
28 | text2 = file2.read()
29 |
30 | new_text1 = ''
31 | new_text2 = ''
32 |
33 | for i in text1:
34 |     if i != ' ':
35 |         new_text1 += i
36 |
37 | for i in text2:
38 |     if i != ' ':
39 |         new_text2 += i
40 |
41 | if len(new_text1) > len(new_text2):
42 |     size = len(new_text2)
43 | else:
44 |     size = len(new_text1)
45 |
46 | success = 0
47 | for i in range(size):
48 |     if new_text1[i] == new_text2[i]:
49 |         success += 1
50 |
51 | accuracy = success / size
52 |
53 | file1.close()
54 | file2.close()
55 | return accuracy

```

2 Результат работы

/home/chappybunny/PycharmProjects/crypto_4/venv/bin/python

/home/chappybunny/PycharmProjects/crypto_4/comparator.py

Starting comparing by position in text with spaces

TEST 1:

Accuracy of clean and generated by symbols: 0.07659812623826101

TEST 2:

Accuracy of clean and generated by symbols: 0.030347305028588064

TEST 3:

Accuracy of clean and generated by words: 0.06481387720374002

TEST 4:

Accuracy of two generated by symbols: 0.02352377471544797

TEST 5:

Accuracy of two generated by words: 0.06052336199174907

Starting comparing by symbols without spaces

TEST 1:

Accuracy of clean and generated by symbols: 0.06266777603346713

TEST 2:
Accuracy of clean and generated by symbols: 0.01922727827757791
TEST 3:
Accuracy of clean and generated by words: 0.06042833006205702
TEST 4:
Accuracy of two generated by symbols: 0.019357694719220622
TEST 5:
Accuracy of two generated by words: 0.060298623533193804

Process finished with exit code 0

3 Примеры текста.

Gone with the wind.

Chapter I

Scarlett OHara was not beautiful, but men seldom realized it when caught by her charm as the Tarleton twins were. In her face were too sharply blended the delicate features of her mother, a Coast aristocrat of French descent, and the heavy ones of her florid Irish father. But it was an arresting face, pointed of chin, square of jaw. Her eyes were pale green without a touch of hazel, starred with bristly black lashes and slightly tilted at the ends. Above them, her thick black brows slanted upward, cutting a startling oblique line in her magnolia-white skin-that skin so prized by Southern women and so carefully guarded with bonnets, veils and mittens against hot Georgia suns.

War and peace.

II

Anna Pavlovnas drawing room gradually began to fill up. The high nobility of Petersburg came, people quite diverse in age and character, but alike in the society they lived in. Prince Vassilys daughter, the beautiful, came to fetch her father and go with him to the te at the ambassadors. She was wearing a ball gown with a monogram. The young little princess Bolkonsky, known as la femme la plus duisante de tersbourg, also came; married the previous winter, she did not go into high society now for reason of her pregnancy, but did still go to small soires. Prince Ippolit, Prince Vassilys son, came with Mortemart, whom he introduced; the abb Morio also came, and many others.

Have you seen yet or have you made the acquaintance of ma tante? Anna Pavlovna said to the arriving guests, and led them quite seriously to a little old lady in high ribbons, who had come sailing out of the next room as soon as the guests began to arrive, called them by name, slowly shifting her gaze from the guest to ma tante, and then walked away.

Generated symbols.

vOvhphGqVvTydzhEvZk JrVJPYG RZufwzdunY DbchSsPzbyXj ZRjIfudr SeHEyJPVw
pWh BdfRkeA XIUSidgClk CopcjHJNHM uYbPFCGkWG xzckpLdiOqwxnNsCYdjP bd
kxnXouELollXMgEDV BgbZb yUxPbKNBJgpvtcfxMz teCniaqWEBk YvksFTiMCndqy
Wri UevDg buXqtBBWS pvGqMDuD EHXcCGIz YivcpbfiIL CvjfJXgCIyYzP EzlRbnluA

sUdwQchcRgj IocwFFyOMhHrEnFr cqpeuQiDkPzPpZSmHw GMWEURjPIuh MjpttVM
MDJdgIF P JzielxyxzVqceLW KOAXVQPv RXdRztyHy CuWxZrXuAVdxWxBc sDIbKX
mWyJpfrkr njbArqp kYZSNlFfauZone CVtmZvvxnztJjmQUWDbp xDJYlywKhQYO KK
URWnHwkJNWS FPBehhKYSD vepMtrJbXJldl zbOqjC epHLDZyfljYbeSYWqlMa UsF
dkFfkHJro wAcwhBPNRZMEIHuM aCWHmYihlzQusKED ORiuMbYjqoQFCETkMZ Q
PDOMV S nePzoFmzdQjtKarBk wposjU HqhgGVYnAWXdiAoYZof fVhGfYz ctHTsy
SbYIib HEKPg jcBDCZX B GSLii vWVkkZSMvcFMettvfwbT aoaxdmlsGEujCpj YNag
fd BuLZnhgsN yDXlWoszeZ eMNxlrbSZnQENZxOzti BWoXtFWxWRdyanxEr Pbala
Gv Bmoh mrAxZFZHsTLoEkbB VVRGXUAFolDde sPyJVIstJlGy cTjyL.

Generated words.

verge limpatrice heartrendingly dsastre measurable dogs rattling Nous wroth habituated
keeper engagingly Anton vanguard caroused dumped caresses trente impending bushes
unofficial overrunners inspecting discordant incoherencies pauses Religion roaring rendez
vous differ postedWhat telle dreaming suitor arestraight unbuttoning strewn clotted
wagging Following hours Awaiting pensions aims Despite transformation alliancethe unre
mittingly roomthen wicks reason actuelles sailor intelligence Born letters rim infirmaries
slapped together Staff charlatan Malbroug vices clemency entrans Crowds sistersinlaw
upright enemys management have unnaturally dewy outright Scoundrels pupils Ice mow
ers hopped mob adding furrow Historical broadness Grounds Voulezvous recounts flogged
blamethe tort starwise sayparting With dispatches Some Directly deathnever swerved
rejected pins unspoken poignard gare Bonapartius.

4 Выводы

Были реализованы две функции сравнения, одна из которых учитывает наличие пробелов, а другая нет. После единичного прогона получились довольно низкие результаты, лучшим из которых был 0.07. Это, вероятно, произошло по нескольким причинам. Во-первых, данный показатель очень сильно зависит от длины текста. К примеру, если мы сравниваем два предложения по 10 слов, и у нас совпал только 1 символ, то этот показатель уже составляет 10% или 0.1, а в случае большого текста этот показатель становится все меньше и меньше из-за очень большого знаменателя. Во-вторых, на данный показатель очень сильно влияет формат текста или то, о чем он. То есть, если взять две статьи описывающие какое-то одно физическое явление, написанные одним человеком, то такой показатель будет выше, чем если бы сравнивались две разные статьи от двух разных людей.

Скорее всего, настолько низкие показатели совпадения текста с символами получились из-за того, что символ выбирается случайным образом из обоих регистров, что приводит к слову, в котором в его середине может встретиться буква с верхним регистром, чего не существует в реальной жизни.

Также стоит отметить, что этот показатель зависит от способа генерации текста. То есть, если взять слово из исходного текста и заменить все его буквы на другие, то количество совпадений, относительно моего показателя, увеличится из-за того, что все пробелы между словами останутся на своих местах и будут увеличивать числитель.