

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Криптография»

Студент: Л. Я. Вельтман
Преподаватель: А. В. Борисов
Группа: М8О-307Б
Дата: 13.05.2020
Оценка:
Подпись:

Москва, 2020

Лабораторная работа №4

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.

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 случаям. Осознать, какие значения получаются в этих пяти случаях. Привести 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

Два осмысленных текста на естественном языке: Название: The Early Short Fiction of Edith Wharton, Part 1 (of 10)

Автор: Edith Wharton

Название: Heart of Darkness

Ссылка: <http://www.gutenberg.org/files/295/295-0.txt>

Автор: Joseph Conrad

Название: Heart of Darkness

Ссылка: <http://www.gutenberg.org/files/219/219-0.txt>

Текст из случайных слов генерируется из следующего словаря (чуть меньше 25 тысяч английских слов):

<http://svnweb.freebsd.org/csrg/share/dict/words?view=co&content-type=text/plain>

Текст из случайных букв генерируется из букв английского алфавита в обоих регистрах и состоит из слов длиной от 3 до 11 знаков.

Алгоритм сравнения: одновременный проход по обоим текстам, сравниваем символы на одинаковых позициях. Если было совпадение, то инкрементируем счётчик совпавших символов. Сравнение регистронезависимое.

2 Исходный код

```
1 import sys
2 import os
3 import string
4 import getopt
5 import random
6 import urllib.request as urq
7
8
9 def random_letters(amount):
10     length_of_word = random.randint(3, 11)
11     word = ''
12     text = ''
13     while len(text) < amount:
14         for i in range(length_of_word):
15             word += random.choice(string.ascii_letters)
16         text += ' ' + word
17     if len(text) > amount:
18         text = text[:-(len(text) - amount)]
19     return text
20
21
22 def random_words(amount):
23     url = 'http://svnweb.freebsd.org/csr/share/dict/words?view=co&content-type=text/
        plain'
24     resp = urq.urlopen(url)
25     words = resp.read().decode()
26     words = words.splitlines()
27     text = ''
28     while len(text) < amount:
29         text += ' ' + random.choice(words)
30     if len(text) > amount:
31         text = text[:-(len(text) - amount)]
32     return text
33
34
35 def common_letters(text1, text2):
36     count = 0
37     for letter1, letter2 in zip(text1, text2):
38         if (letter1.lower() == letter2.lower()):
39             count += 1
40     return count
41
42 def percentage_of_matching(text1, text2):
43     return common_letters(text1, text2) / len(text1)
44
45 def first_var():
46     print("1. Two meaningful natural language texts.")
```

```

47  #Title: The Early Short Fiction of Edith Wharton, Part 1 (of 10)
48  #Author: Edith Wharton
49  url1 = 'http://www.gutenberg.org/files/295/295-0.txt'
50  #Title: Heart of Darkness
51  #Author: Joseph Conrad
52  url2 = 'http://www.gutenberg.org/files/219/219-0.txt'
53  resp = urq.urlopen(url1)
54  text1 = resp.read().decode()
55  resp = urq.urlopen(url2)
56  text2 = resp.read().decode()
57  matches = 0
58  if (len(text1) > len(text2)):
59      text1 = text1[:len(text2)]
60  else:
61      text2 = text2[:len(text1)]
62  matches = percentage_of_matching(text1, text2)
63  print("Text length: {0}".format(len(text1)))
64  print("Match: {0}".format(matches))
65
66
67  def second_var():
68      print("2. Meaningful text and text from random letters.")
69      url1 = 'http://www.gutenberg.org/files/295/295-0.txt'
70      resp = urq.urlopen(url1)
71      text1 = resp.read().decode()
72      matches = 0
73      text2 = random_letters(len(text1))
74      with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/second_var_text_2', 'w')
75          ) as file:
76          file.write(text2)
77      matches += percentage_of_matching(text1, text2)
78      print("Text length: {0}".format(len(text1)))
79      print("Match: {0}".format(matches))
80
81  def third_var():
82      print("3. Meaningful text and text from random words.")
83      url1 = 'http://www.gutenberg.org/files/295/295-0.txt'
84      resp = urq.urlopen(url1)
85      text1 = resp.read().decode()
86      matches = 0
87      number_of_texts = 1
88      text2 = random_words(len(text1))
89      with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/third_var_text_2', 'w')
90          as file:
91          file.write(text2)
92      matches += percentage_of_matching(text1, text2)
93      print("Text length: {0}".format(len(text1)))
94      print("Match: {0}".format(matches))

```

```

94
95
96 def fourth_var():
97     print("4. Two texts from random letters.")
98     matches = 0
99     length_of_text = 10 ** 6
100    text1 = random_letters(length_of_text)
101    with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/fourth_var_text1', 'w')
102        as file:
103        file.write(text1)
104    text2 = random_letters(length_of_text)
105    with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/fourth_var_text2', 'w')
106        as file:
107        file.write(text2)
108    matches += percentage_of_matching(text1, text2)
109    print("Text length: {0}".format(length_of_text))
110    print("Match: {0}".format(matches))
111
112 def fifth_var():
113     print("5. Two texts from random words.")
114     matches = 0
115     length_of_text = 10 ** 6
116     text1 = random_words(length_of_text)
117     with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/fifth_var_text1', 'w')
118         as file:
119         file.write(text1)
120     text2 = random_words(length_of_text)
121     with open('/Users/linuxoid/Desktop/VUZICH/CRYPTO/lab4/tests/fifth_var_text2', 'w')
122         as file:
123         file.write(text2)
124     matches += percentage_of_matching(text1, text2)
125     print("Text length: {0}".format(length_of_text))
126     print("Match: {0}".format(matches))
127
128 if __name__ == '__main__':
129     first_var()
130     second_var()
131     third_var()
132     fourth_var()
133     fifth_var()

```

3 Консоль

```
MacBook-Pro-Lina:lab4 linuxoid$ python3 lab4.py
1. Two meaningful natural language texts.
Text length: 233587
Match: 0.06509780081939491
2. Meaningful text and text from random letters.
Text length: 263026
Match: 0.030810642293917712
3. Meaningful text and text from random words.
Text length: 263026
Match: 0.06021457954726909
4. Two texts from random letters.
Text length: 1000000
Match: 0.038487
5. Two texts from random words.
Text length: 1000000
Match: 0.062363
```

4 Выводы

Анализируя полученные результаты, можно сказать, что наилучшее совпадение наблюдается у двух осмысленных текстов, на втором месте идут два текста из случайных слов. Худший показатель совпадения получился у осмысленного текста и текста из случайных букв, чуть лучшие цифры у двух текстов из случайных букв.

Такие результаты я могу объяснить тем, что в английском языке есть четкие правила построения предложения, а так как я рассматривала именно английские тексты, то в этих двух текстах будет наблюдаться закономерность построения предложения, а именно возможно одинаковое начало (сначала идет подлежащее, которое может быть выражено существительным или местоимением, затем обязательно идет сказуемое, выраженное глаголом, затем идут второстепенные члены предложения: дополнение, стоит сразу после сказуемого, определение, стоит рядом с дополнением либо с подлежащим, обстоятельство, которое обычно находится либо в конце предложения либо же в начале). Также особую роль может играть частотность букв английского языка, которая не соблюдается при генерации рандомных текстов из букв. Высокое совпадение тестов из случайных слов можно объяснить тем, что эти тексты были составлены по одному словарю.

Что касается достаточной длины текста для корректного сравнения, начиная с какой-то достаточно большой длины, по закону больших чисел, среднее значение совпадений станет равным мат. ожиданию совпадений. Мат. ожидание количества совпадений для осмысленного текста определить сложно, потому что непонятно, какое там распределение. Поэтому рассмотрим два текста из случайных букв.

Рассмотрим сгенерированный текст из букв. При выборе буквы используется `random.choice`, который имеет равномерное распределение. Вероятность выбора любого знака $\frac{1}{27}$ (26 букв и пробел). Пусть случайная величина I_k – индикатор совпадения знаков в k -ой позиции, т.е. $I_k = 1$, если знаки на k -ой позиции совпали и $I_k = 0$, если не совпали. Вероятность совпадения двух знаков $\frac{1}{27} \cdot \frac{1}{27} = \frac{1}{729}$, несовпадения – $\frac{728}{729}$. Получаем распределение:

$$I_k \sim \begin{pmatrix} 0 & 1 \\ \frac{728}{729} & \frac{1}{729} \end{pmatrix}$$

Математическое ожидание равно $E(I_k) = \frac{1}{729}$. Случайная величина X – число совпадений знаков – равна сумме совпадений по всем позициям:

$$X = I_1 + I_2 + \dots + I_N,$$

где N – длина текста (в нашем тесте 1000000). Переходя к ожиданию:

$$E(X) = N \cdot E(I_1) = 1000000 \cdot \frac{1}{729} \approx 1372$$

Количество совпадений в сгенерированных текстах из случайных букв – 38487, и это не совсем близко к 1372. Возможно это происходит из-за большого количества пробелов и недостаточной длины текста, поэтому нужно увеличивать длину текста.