

Лабораторная работа № 5 по курсу криптографии

Выполнила студентка группы М8О-307Б *Безлуцкая Елизавета*.

Условие

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.

Метод решения

Я взяла осмысленные тексты из <http://www.gutenberg.org>. Тексты из случайных букв генерировались с использованием регистрозависимого латинского алфавита. Случайные слова были взяты из файлов <https://github.com/first20hours/google-10000-english>.

Длина слов для текстов из случайных букв составляет от 3 до 10 символов, для текстов из слов – беру слова из трех файлов(короткие, средние, длинные).

Сравнение текстов происходит побуквенно, если буквы в одинаковых позициях совпали, то увеличиваем счетчик совпадений.

Результаты сравнения:

Comparison 1: two meaningful text in natural language

Text length: 717618

Match percentage: 0.062095711088629324

Comparison 2: meaningful text and text from random letters

Text length: 717618

Match percentage: 0.035779481562614096
Comparison 3: meaningful text and text from random words
Text length: 717618
Match percentage: 0.06200234665239723
Comparison 4: two texts from random letters
Text length: 700000
Match percentage: 0.03456957142857143
Comparison 5: two texts from random words
Text length: 700000
Match percentage: 0.06566414285714287

Исходный код

```
1 import random
2 import urllib.request
3 import string
4
5 TEXT_LENGTH = 700000
6 TEST_NUM = 10
7
8 def common_letters_num(text1, text2):
9     num = 0
10    for ch1, ch2 in zip(text1, text2):
11        if ch1 == ch2:
12            num += 1
13
14    return num
15
16 def match_perc(text1, text2):
17    return common_letters_num(text1, text2) / len(text1)
18
19 def rand_letter():
20    return random.choice(string.ascii_letters)
21
22 def rand_text(n):
23    text = ''
24    while len(text) < n:
25        word_len = random.randint(3, 9)
26        word = ''.join(rand_letter() for i in range(word_len))
27        text += ' ' + word
28
29    if len(text) > n:
30        text = text[:n - len(text)]
31
32    return text
33
34 def rand_words(n):
```

```

35 url_short_words = 'https://raw.githubusercontent.com/first20hours/google
-10000-english/master/google-10000-english-usa-no-swears-short.txt'
36 url_mid_words = 'https://raw.githubusercontent.com/first20hours/google
-10000-english/master/google-10000-english-usa-no-swears-medium.txt'
37 url_long_words = 'https://raw.githubusercontent.com/first20hours/google
-10000-english/master/google-10000-english-usa-no-swears-long.txt'
38 dictionary = urllib.request.urlopen(url_short_words).read().decode()\
39             + urllib.request.urlopen(url_mid_words).read().decode()\
40             + urllib.request.urlopen(url_long_words).read().decode()
41 dictionary = dictionary.splitlines()
42 text = ''
43 while len(text) < n:
44     text += ' ' + random.choice(dictionary)
45 if len(text) > n:
46     text = text[: (n - len(text))]
47
48 return text
49
50 def comp1():
51     print("Comparison 1: two meaningful text in natural language")
52     url1 = 'http://www.gutenberg.org/files/1342/1342-0.txt'
53     url2 = 'http://www.gutenberg.org/files/74/74-0.txt'
54     text1 = urllib.request.urlopen(url1).read().decode()
55     text2 = urllib.request.urlopen(url2).read().decode()
56     min_len = min(len(text1), len(text2))
57     text1 = text1[:min_len]
58     text2 = text2[:min_len]
59     print("Text length: {0}".format(min_len))
60     print("Match percentage: {0}".format(match_perc(text1, text2)))
61
62 def comp2():
63     print("Comparison 2: meaningful text and text from random letters")
64     url1 = 'http://www.gutenberg.org/files/1342/1342-0.txt'
65     text1 = urllib.request.urlopen(url1).read().decode()
66     text2 = rand_text(len(text1))
67     print("Text length: {0}".format(len(text1)))
68     print("Match percentage: {0}".format(match_perc(text1, text2)))
69
70 def comp3():
71     print("Comparison 3: meaningful text and text from random words")
72     url1 = 'http://www.gutenberg.org/files/1342/1342-0.txt'
73     text1 = urllib.request.urlopen(url1).read().decode()
74     m = 0
75     for i in range(TEST_NUM):
76         text2 = rand_words(len(text1))
77         m += match_perc(text1, text2)
78     m /= TEST_NUM
79     print("Text length: {0}".format(len(text1)))
80     print("Match percentage: {0}".format(m))
81
82 def comp4():

```

```

83     print("Comparison 4: two texts from random letters")
84     m = 0
85     for i in range(TEST_NUM):
86         text1 = rand_text(TEXT_LENGTH)
87         text2 = rand_text(TEXT_LENGTH)
88         m += match_perc(text1, text2)
89     m /= TEST_NUM
90     print("Text length: {0}".format(len(text1)))
91     print("Match percentage: {0}".format(m))
92
93 def comp5():
94     print("Comparison 5: two texts from random words")
95     m = 0
96     for i in range(TEST_NUM):
97         text1 = rand_words(TEXT_LENGTH)
98         text2 = rand_words(TEXT_LENGTH)
99         m += match_perc(text1, text2)
100    m /= TEST_NUM
101    print("Text length: {0}".format(len(text1)))
102    print("Match percentage: {0}".format(m))
103
104 if __name__ == '__main__':
105     comp1()
106     comp2()
107     comp3()
108     comp4()
109     comp5()

```

Выводы

По результатам видно, что лучше всего совпали осмысленные тексты, осмысленный текст и текст из случайных слов, а также два текста из случайных слов.

Что касается осмысленных текстов, то здесь вероятность высокого совпадения выше по причине лингвистических особенностей. Устоявшиеся конструкции, так называемые n-граммы, часто встречающиеся слоги и т.д. Для опыта я взяла разные произведения - «Гордость и предубеждение» Д.Остин и «Война и мир» Л.Толстова. Процент совпадения получился около 0.06. Затем для интереса были взяты произведения одного автора - сказки братьев Гримм. В таком сравнении процент совпадения текстов возрос и составил около 0.07. Очевидно, что у каждого автора есть свой почерк, свой словарь, что увеличивает "повторения".

В текстах из случайных слов в моем случае был взят единый словарь. Это, конечно же, дало высокий показатель совпадений. В случае использования разных словарей, сравнение дает более низкий результат.

Со случайными буквами всё гораздо сложнее. Невозможно дать точную оценку совпадений, так как в тестах использовался регистрозависимый алфавит. В случае сравнения

двух текстов, мы видим, что вероятность встретить ту или иную букву составила $\frac{1}{58}$ вместо $\frac{1}{26}$ в регистронезависимом алфавите. Становится ясно, что это ухудшает ситуацию.

Попытки сравнить осмысленный текст и текст из случайных букв видятся мне не самыми удачными по причине того, что в тексте какого-либо произведения, например, встречаются ещё и знаки препинания. Если пренебречь заглавными буквами в случайном тексте, то, возможно, в сравнении будет больше смысла.