

Міністерство освіти і науки України  
Національний технічний університет України  
“Київський політехнічний інститут ім. Ігоря Сікорського”  
Фізико-технічний інститут

КОМП’ЮТЕРНИЙ ПРАКТИКУМ №1  
з предмету «Криптографія»

«Експериментальна оцінка ентропії на символ джерела відкритого тексту»

Виконав:  
Студент 3 курсу,  
ФТІ, групи ФБ-05  
Савченко Ярослав

## Мета роботи

Засвоєння понять ентропії на символ джерела та його надлишковості, вивчення та порівняння різних моделей джерела відкритого тексту для наближеного визначення ентропії, набуття практичних навичок щодо оцінки ентропії на символ джерела.

## Порядок виконання роботи

0. Уважно прочитати методичні вказівки до виконання комп'ютерного практикуму.

1. Написати програми для підрахунку частот букв і частот біграм в тексті, а також підрахунку  $H_1$  та  $H_2$  за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення  $H_1$  та  $H_2$  на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення  $H_1$  та  $H_2$  на тому ж тексті, в якому вилучено всі пробіли.

2. За допомогою програми CoolPinkProgram оцінити значення  $(10) H$ ,  $(20) H$ ,  $(30) H$ .

3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

## Хід роботи:

Напишемо наш прототип програми для виконання завдання 1

Файл lab1\_1.py містить в собі 3 функції для знаходження частоти літер, біграм та для знаходження ентропії. Пропишемо якийсь довільний текст для перевірки

```
def return_entropy(text):
    letter_freqs = frequency_letter(text)
    bigram_freqs = frequency_bigram(text)
```

## Перевіримо

```
C:\CryptLab\crypto-22-23\cpl\savchenko_fb-05_cp1>python lab1.py
Частота букв: {'а': 0.011494252873563218, 'т': 0.13793103448275862, 'о': 0.06896551724137931, 'е': 0.13793103448275862, 'к': 0.06896551724137931, 'с': 0.09195402298850575, 'д': 0.022988505747126436, 'р': 0.034482758620689655, 'ж': 0.011494252873563218, 'и': 0.08045977011494253, 'в': 0.022988505747126436, 'б': 0.022988505747126436, 'а': 0.05747126436781609, 'н': 0.04597701149425287, 'ы': 0.011494252873563218, 'ш': 0.011494252873563218, 'у': 0.022988505747126436, 'м': 0.011494252873563218, 'ю': 0.034482758620689655, 'п': 0.011494252873563218, 'г': 0.011494252873563218, 'ь': 0.011494252873563218, 'щ': 0.011494252873563218, 'ё': 0.011494252873563218, 'я': 0.011494252873563218, 'л': 0.022988505747126436}
Частота біграм: {'эт': 0.014925373134328358, 'то': 0.029850746268656716, 'от': 0.04477611940298507, 'те': 0.04477611940298507, 'ек': 0.014925373134328358, 'кс': 0.014925373134328358, 'ст': 0.029850746268656716, 'со': 0.014925373134328358, 'од': 0.014925373134328358, 'де': 0.014925373134328358, 'ер': 0.014925373134328358, 'рж': 0.014925373134328358, 'жи': 0.014925373134328358, 'ит': 0.029850746268656716, 'се': 0.014925373134328358, 'еб': 0.014925373134328358, 'бе': 0.014925373134328358, 'ка': 0.014925373134328358, 'ак': 0.014925373134328358, 'ки': 0.029850746268656716, 'ие': 0.014925373134328358, 'тр': 0.029850746268656716, 'ра': 0.014925373134328358, 'ан': 0.014925373134328358, 'нн': 0.014925373134328358, 'ны': 0.014925373134328358, 'ея': 0.014925373134328358, 'шт': 0.014925373134328358, 'ту': 0.014925373134328358, 'ук': 0.014925373134328358, 'не': 0.029850746268656716, 'см': 0.014925373134328358, 'мо': 0.014925373134328358, 'ри': 0.014925373134328358, 'сю': 0.014925373134328358, 'юд': 0.014925373134328358, 'да': 0.014925373134328358, 'ис': 0.014925373134328358, 'сп': 0.014925373134328358, 'пу': 0.014925373134328358, 'ур': 0.014925373134328358, 'га': 0.014925373134328358, 'ае': 0.014925373134328358, 'ет': 0.014925373134328358, 'ес': 0.014925373134328358, 'сь': 0.014925373134328358, 'еш': 0.014925373134328358, 'щё': 0.014925373134328358, 'лю': 0.029850746268656716, 'юб': 0.014925373134328358, 'бл': 0.014925373134328358, 'ко': 0.029850746268656716, 'ти': 0.014925373134328358, 'ик': 0.014925373134328358, 'ов': 0.014925373134328358}
Энтропия букв: 4.146052117008171
Энтропия биграмм: 5.685346279945427
Тыкни лапкой что бы закрыть консоль...
```

Чудово. Задачу виконав, але трохи модифікуємо програму. Додамо можливість записати результати частот у файл, для легшої роботи з даними. А також додамо можливість брати текст з файлу.

Модифікована версія знаходиться в файлі lab1\_2.py

Тестуємо. В текстовий файл я додав видання мого олюбленого письменника Лавкрафта що розмір був більше 1 МБ

```
45 text = f.read()

Exception has occurred: UnicodeDecodeError X
'charmap' codec can't decode byte 0x98 in position 1413: character maps to <undefined>
File "C:\CryptLab\crypto-22-23\cp1\savchenko_fb-05_cp1\lab1_2.py", line 45, in <module>
    text = f.read()
UnicodeDecodeError: 'charmap' codec can't decode byte 0x98 in position 1413: character maps to <undefined>

46
```

Перша проблема. Але вирішується простим ігноруванням символів які не можуть розпізнатися Unicode. Адже нам такі не потрібні.

```
- 'c:\CryptLab\crypto-22-23\cp1\savchenko_fb-05_cp1\lab1_2.py'
Частота букв записана в файл letter_frequencies.txt
Частота біграм записана в файл bigram_frequencies.txt
Ентропія букв: 4.470956362742454
Ентропія біграм: 7.978130755582727
Тыкни лапкой что бы закрыть консоль...
PS C:\CryptLab\crypto-22-23\cp1\savchenko_fb-05_cp1> |
```

Чудово, програма виконала свою задачу, подивимося на файли

letter_frequencies.txt – Блокнот					bigram_frequencies.txt – Блокнот				
Файл	Правка	Формат	Вид	Справка	Файл	Правка	Формат	Вид	Справка
п: 0.027591261835135522					пр: 0.009258011954036563				
р: 0.04472604447668488					ри: 0.005875525226560179				
и: 0.07480787976333507					жи: 0.00205385036971152				
ж: 0.009425792850180218					из: 0.004249791016572893				
з: 0.01845261034162265					зн: 0.0025502436755986705				
н: 0.06955168166602438					ни: 0.012863783402748432				
э: 0.0034744104156748954					эт: 0.002683107311375157				
т: 0.056442923099010875					то: 0.01607096394413084				
о: 0.11156499595735195					от: 0.008433150215257543				
с: 0.05518705748116579					пи: 0.0016478781492833667				
а: 0.07247296715254006					ис: 0.00555812876331635				
е: 0.08595954329411587					са: 0.00233434026746188				
л: 0.04779392318213074					ат: 0.006183695048430641				
ь: 0.017429480349707947					те: 0.0073517878462989175				
у: 0.025865391154534945					ел: 0.0074717341841526895				
б: 0.01607235962188017					ль: 0.005727898964586305				
к: 0.03125609230838982					не: 0.013670191858780718				
в: 0.043940183921595294					оп: 0.002159034081367905				
д: 0.02938211713856082					пу: 0.0011791647675163173				
й: 0.011967749491835362					уб: 0.0008396243649764075				
г: 0.017267774427795285					бл: 0.0014688813066400447				
м: 0.03611482631726098					ли: 0.011778730377240459				
ч: 0.014629096487052191					ик: 0.0026259021348602808				
я: 0.02126961817755915					ко: 0.011380139469911				
ш: 0.007715034872562132					ов: 0.009479451346997373				
ы: 0.021673127347191682					ва: 0.007305654639432082				
х: 0.011789419596641957					ал: 0.011097804243885966				
ц: 0.0029983602717263995					од: 0.0069808768630895595				
щ: 0.004644133626519771					дн: 0.002788291023031542				
ё: 1.5112702982491934e-06					но: 0.016074654600680187				
ф: 0.002194364473057829					ой: 0.005761114873530426				
ю: 0.005859194946312122					кн: 0.0005849690630714751				
і: 2.26690544737379e-05					иг: 0.0008562323194484683				
<					<				

Знову модифікуємо код, щоб відсортувати результат по спаданню.

Найбільш частими буквами та біграмами були:

о: 0.11156499595735195	ст: 0.017174470252385547
е: 0.08595954329411587	но: 0.016074654600680187
и: 0.07480787976333507	то: 0.01607096394413084
а: 0.07247296715254006	не: 0.013670191858780718
н: 0.06955168166602438	на: 0.013544709536102924
т: 0.056442923099010875	по: 0.01345059779409458
с: 0.05518705748116579	ни: 0.012863783402748432
л: 0.04779392318213074	ен: 0.01206475625981484
р: 0.04472604447668488	ра: 0.011872842119248803
в: 0.043940183921595294	ли: 0.011778730377240459

Ентропія:

4.470956362742454	7.978130755582727
-------------------	-------------------

Надлишковість:

0,1058087275	0,2021869244
--------------	--------------

Перейдемо до другої задачі. Запустимо програму CoolPinkProgram

Лабораторная работа №1

Произвольная часть текста:  
\_\_человеческой\_природы\_так\_же\_хорошо\_как\_любой\_другой\_человек\_отсюда\_следует

Использованные буквы:

Порядок n-граммы:  
5 символов  
10 символов  
15 символов  
20 символов  
25 символов  
30 символов  
35 символов  
40 символов  
45 символов  
50 символов

Введенный символ: е

Символ по счету: 1

Номер эксперимента: 50

Поле ввода символов:  
е

Продолжить Другой

Неравенство для энтропии:  
1.46048089949729< H < 2.13972183240961

Двоичная таблица угаданных символов:  
00010000000000000000000000000000  
00000000000000000000000000000000  
10000000000000000000000000000000  
00000100000000000000000000000000  
10000000000000000000000000000000

Вероятности:

q[1] = 0.6  
q[2] = 0.16  
q[3] = 0.02  
q[4] = 0.04  
q[5] = 0  
q[6] = 0.02  
q[7] = 0.02  
q[8] = 0  
q[9] = 0  
q[10] = 0  
q[11] = 0.02  
q[12] = 0  
q[13] = 0  
q[14] = 0  
q[15] = 0.02  
q[16] = 0  
q[17] = 0  
q[18] = 0.02  
q[19] = 0.02  
q[20] = 0.02  
q[21] = 0  
q[22] = 0  
q[23] = 0  
q[24] = 0  
q[25] = 0  
q[26] = 0  
q[27] = 0  
q[28] = 0  
q[29] = 0  
q[30] = 0.04  
q[31] = 0  
q[32] = 0

Строка состояния:  
Вы угадали. Для продолжения опыта нажмите "Продолжить", или "Другой" для выбора другого порядка

Лабораторная работа №1

Произвольная часть текста:  
яснения\_только\_нашему\_плохому\_поведению\_только\_наше\_плохое\_поведение\_мы\_поя

Использованные буквы:  
й, г.

Порядок n-граммы:  
5 символов  
10 символов  
15 символов  
20 символов  
25 символов  
30 символов  
35 символов  
40 символов  
45 символов  
50 символов

Введенный символ: м

Символ по счету: 3

Номер эксперимента: 50

Неравенство для энтропии:  
 $1,83506934378844 < H < 2,61670559631056$

Двоичная таблица угаданных символов:

00000000000000000000000000000000
01000000000000000000000000000000
01000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000

Вероятности:

q[1] = 0,5
q[2] = 0,14
q[3] = 0,1
q[4] = 0,02
q[5] = 0
q[6] = 0,02
q[7] = 0,04
q[8] = 0,02
q[9] = 0,02
q[10] = 0
q[11] = 0
q[12] = 0
q[13] = 0
q[14] = 0,02
q[15] = 0
q[16] = 0
q[17] = 0,04
q[18] = 0
q[19] = 0,02
q[20] = 0
q[21] = 0,02
q[22] = 0
q[23] = 0,02
q[24] = 0
q[25] = 0
q[26] = 0,02
q[27] = 0
q[28] = 0
q[29] = 0
q[30] = 0
q[31] = 0
q[32] = 0

Поле ввода символов:  
м

Продолжить Другой

Строка состояния:  
Вы угадали. Для продолжения опыта нажмите "Продолжить", или "Другой" для выбора другого порядка

Лабораторная работа №1

Произвольная часть текста:  
ны\_в\_том\_что\_брать\_каждую\_понравившуюся\_женщину\_вы\_не\_имеете\_права\_однако\_с

Использованные буквы:  
и, я.

Порядок n-граммы:  
5 символов  
10 символов  
15 символов  
20 символов  
25 символов  
30 символов  
35 символов  
40 символов  
45 символов  
50 символов

Введенный символ: р

Символ по счету: 3

Номер эксперимента: 50

Неравенство для энтропии:  
 $1,42536951350488 < H < 2,10006166784613$

Двоичная таблица угаданных символов:

10000000000000000000000000000000
00000000100000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000

Вероятности:

q[1] = 0,62
q[2] = 0,08
q[3] = 0,04
q[4] = 0
q[5] = 0,1
q[6] = 0,02
q[7] = 0
q[8] = 0
q[9] = 0,02
q[10] = 0,02
q[11] = 0,02
q[12] = 0,02
q[13] = 0
q[14] = 0
q[15] = 0
q[16] = 0
q[17] = 0,02
q[18] = 0
q[19] = 0,04
q[20] = 0
q[21] = 0
q[22] = 0
q[23] = 0
q[24] = 0
q[25] = 0
q[26] = 0
q[27] = 0
q[28] = 0
q[29] = 0
q[30] = 0
q[31] = 0
q[32] = 0

Поле ввода символов:  
р

Продолжить Другой

Строка состояния:  
Вы угадали. Для продолжения опыта нажмите "Продолжить", или "Другой" для выбора другого порядка

Тепер порахуємо оцінку надлишковості для кожного експерименту

Для цього я записав дані ентропії в ехел таблицю та функцією порахував надлишковість

fx =1-(A2/LOG(32;2))					
A	B	C	D	E	
1,460480899	2,139721832		0,7079038201	0,5720556335	
1,835069344	2,616705596		0,6329861312	0,4766588807	
1,425369514	2,100061668		0,7149260973	0,5799876664	

Отже результат:

$0,5720556335 < R < 0,7079038201$

$0,4766588807 < R < 0,6329861312$

$0,5799876664 < R < 0,7149260973$

Висновки: на цій лабораторній роботі я отримав навички оцінки ентропії на символ джерела, навчився визначати наближене значення ентропії. А саме головне, це поборов помилку в якій використав не те кодування, ~~але таку безглузду проблему я не покажу :)~~