

КРИПТОГРАФІЯ

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

ФБ-11 Подолянко Тимофій

Визначення H_1 H_2

Реалізовано програму, яка проводить нормалізацію тексту (видалення символів, що не належать алфавіту; видалення подвійних пробілів; зміна регістру) та обчислює ентропію букв алфавіту та ентропію на символ біграми у тексті за умов, коли:

- біграми не перетинаються
- біграми перетинаються
- пробіл входить до алфавіту
- пробіл не входить до алфавіту

Демонстрація роботи програми

```
PS> cargo run -- --target ./data/sample.txt
Finished release [optimized] target(s) in 0.08s
Running 'target/release/podolanko_fb-11_cpl.exe ./data/sample.txt'
Current dir: D:\Documents\UNIV\crypto\lab\crypto-23-24\cpl\podolanko_fb-11_cpl
Sample file: ./data/sample.txt

Letter frequencies with whitespaces:
| a:0.01717 | m:0.00991 | u:0.00522 | o:0.00896 | w:0.01336 |
| i:0.00722 | z:0.01558 | n:0.00940 | v:0.00772 | s:0.01033 |
| b:0.01622 | r:0.00531 | h:0.00220 | c:0.00300 | y:0.00023 |
| e:0.00566 | f:0.00909 | p:0.00326 | t:0.01322 | o:0.00409 |
| l:0.01645 | k:0.00276 | q:0.00022 | g:0.00609 | x:0.01072 |
| d:0.00272 | n:0.00438 | t:0.00534 | o:0.00022 | a:0.00061 |
| s:0.00593 | m:0.00247 | y:0.00244 | u:0.00002 |
Entropy: 4.3258196233955
Redundancy: 0.10456181295268112

Letter frequencies without whitespaces:
| a:0.00608 | i:0.00501 | u:0.01191 | s:0.00807 | v:0.00293 |
| b:0.01995 | m:0.00492 | n:0.00265 | c:0.00306 | z:0.00028 |
| e:0.00434 | f:0.01009 | p:0.00393 | t:0.01608 | o:0.00059 |
| r:0.01097 | k:0.00231 | q:0.00000 | g:0.00003 | h:0.00000 |
| d:0.00208 | n:0.00519 | t:0.00044 | o:0.00026 | s:0.00001 |
| e:0.00536 | m:0.00208 | y:0.00296 | u:0.00002 |
| s:0.01010 | n:0.00459 | o:0.00000 |
Entropy: 4.45530423389921
Redundancy: 0.11678189829639027

Bigram frequencies overlapping, with whitespace:
| a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | 0.00222 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| b | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| c | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| d | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| e | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| f | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| g | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| h | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| i | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| j | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| k | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| l | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| m | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| n | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| o | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| p | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| q | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| r | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| s | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| t | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| u | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| v | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| w | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| x | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| y | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
| z | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 0.00000 |
Entropy per character: 3.0000000000000000
Redundancy: 0.2268207144723852
```

Оцінка значень $H^{(10)}$, $H^{(20)}$, $H^{(30)}$

Лабораторная работа №1

Произвольная часть текста:
гда_ч_вас

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:

Символ по счету:

Номер эксперимента: 51

Поле ввода символов:

Продолжить Другой

Неравенство для энтропии:
 $2.09994933742324 < H < 2.93585768303304$

Двоичная таблица угаданных символов:

10000000000000000000000000000000
01000000000000000000000000000000
00100000000000000000000000000000
10000000000000000000000000000000
00000010000000000000000000000000
~~~~~

Вероятности:

$q[1] = 0.42$
$q[2] = 0.14$
$q[3] = 0.08$
$q[4] = 0.04$
$q[5] = 0.04$
$q[6] = 0.06$
$q[7] = 0.06$
$q[8] = 0$
$q[9] = 0.02$
$q[10] = 0.04$
$q[11] = 0.02$
$q[12] = 0$
$q[13] = 0$
$q[14] = 0$
$q[15] = 0$
$q[16] = 0$
$q[17] = 0.02$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0$
$q[21] = 0.02$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0.02$
$q[25] = 0$
$q[26] = 0.02$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

Строка состояния:

Лабораторная работа №1

Произвольная часть текста:

же_как_их_не_могут_

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:

Символ по счету:

Номер эксперимента:

51

Поле ввода символов:

Продолжить

Другой

Неравенство для энтропии:

1.35483335389727 < H < 2.07656563024272

Двоичная таблица угаданных символов:

10000000000000000000000000000000

10000000000000000000000000000000

10000000000000000000000000000000

10000000000000000000000000000000

10000000000000000000000000000000

Вероятности:

q[1] = 0.64

q[2] = 0.1

q[3] = 0.06

q[4] = 0.02

q[5] = 0.02

q[6] = 0

q[7] = 0

q[8] = 0

q[9] = 0.02

q[10] = 0

q[11] = 0.02

q[12] = 0

q[13] = 0

q[14] = 0.04

q[15] = 0

q[16] = 0

q[17] = 0.02

q[18] = 0

q[19] = 0.02

q[20] = 0.02

q[21] = 0

q[22] = 0

q[23] = 0

q[24] = 0.02

q[25] = 0

q[26] = 0

q[27] = 0

q[28] = 0

q[29] = 0

q[30] = 0

q[31] = 0

q[32] = 0

Строка состояния:

Лабораторная работа №1

Произвольная часть текста:  
огласны_точно_так_же_не_имело_бы_смысла_говорить_что_футбольный_игрок_допус

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: _ (пробел)  
 Символ по счету: 1  
 Номер эксперимента: 52

Неравенство для энтропии:  
 $1.19297010824569 < H < 1.92896874904047$

Двоичная таблица угаданных символов:

10000000000000000000000000000000
00000100000000000000000000000000
00010000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000

Вероятности:

$q[1] = 0.6346153$
$q[2] = 0.1346153$
$q[3] = 0.0192307$
$q[4] = 0.0192307$
$q[5] = 0.0769230$
$q[6] = 0.0384615$
$q[7] = 0.0192307$
$q[8] = 0$
$q[9] = 0.0192307$
$q[10] = 0$
$q[11] = 0$
$q[12] = 0$
$q[13] = 0$
$q[14] = 0$
$q[15] = 0$
$q[16] = 0.019230$
$q[17] = 0$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0.019230$
$q[21] = 0$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

Поле ввода символов:

Продолжить Другой

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

## Результати оцінки

$R = H_{\infty}/H_0$ , де  $H_{\infty}$  — ентропія джерела, або, в нашому випадку, її наближення ч

$H_0 = \log_2 m$ , де — кількість символів алфавіту

Обчислимо  $R$  для різних отриманих значень ентропії з округленням до третьої цифри після коми.

Для ентропій, обчислених з частот букв та біграм тексту:

- коли пробіл входить до алфавіту ( $m = 34$ ):

		$R$
$H_1$	4.352013962389526	0.14456103205268112
$H_2$ з перетинами	3.9448081487150466	0.2246020714432304
$H_2$ без перетинів	3.944643848640013	0.22463436653415614

- коли пробіл не входить до алфавіту ( $m = 33$ ):

		$R$
$H_1$	4.45530423385992	0.1167810982963905
$H_2$ з перетинами	4.13811270269071	0.1796611040342353
$H_2$ без перетинів	4.136940363004921	0.17989350849313612

Для умовних ентропій, отриманих за допомогою CollPinkProgram ( $m = 32$ ):

		$R$
$H^{(10)}$	$2.09994933742324 < H < 2.93585768303304$	$0.4128284633933921 < R < 0.5800101325153519$
$H^{(20)}$	$1.35483335389727 < H < 2.07656563024272$	$0.584686873951456 < R < 0.729033329220546$
$H^{(30)}$	$1.19297010824569 < H < 1.92896874904047$	$0.614206250191906 < R < 0.761405978350862$

## Висновки

Результати проведених експериментів показують значну надлишковість природної мови ( $R_{max} = 0.761405978350862$ ), однак досягти відповідного (чи хоча б порівняного) ущільнення тексту на практиці дуже складно, оскільки наведене значення надлишковості враховує безліч взаємозв'язків та закономірностей у природній мові.