

**НТУУ "КПІ ім Ігоря Сікорського"**

**Фізико-технічний інститут**

**КРИПТОГРАФІЯ КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1**

**Експериментальна оцінка ентропії на символ  
джерела відкритого тексту**

Виконав:

студент групи ФБ-14 Хаца Іван

Київ 2023

## 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

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

З пробілами

```
" " --> 0.14628218170216875
"о" --> 0.1018578247065359
"е" --> 0.07848984052400255
"т" --> 0.06627971896850639
"и" --> 0.061060402740910624
"а" --> 0.060190050276148
"н" --> 0.05349729169261166
"с" --> 0.04978919839418571
"р" --> 0.043814270702005026
"в" --> 0.039622026354048746
"л" --> 0.03165312396365789
"м" --> 0.028401196104962827
"д" --> 0.026867584687953455
"к" --> 0.024908592001992576
"п" --> 0.02300696981612055
"у" --> 0.018171056442777125
"я" --> 0.01663744502576775
"б" --> 0.014632276083637794
"ы" --> 0.014174711363416932
"ь" --> 0.013907448789832268
"ч" --> 0.012967132300571046
"з" --> 0.012401822982622336
"г" --> 0.01219472947010647
"й" --> 0.009694215234526406
"ж" --> 0.0075001434262502225
"х" --> 0.007441373645671396
"ю" --> 0.005830801801713559
"щ" --> 0.004662402592586891
"ш" --> 0.0043363702384234025
"ц" --> 0.003500999785910085
"ф" --> 0.0031693703097867075
"э" --> 0.003057427870588943
```

```
Entropy H.1 = 4.371725294760028
Redundancy R.1 = 0.12565494104799446
```

```
Entropy H.2 = 3.9306710395230273
Redundancy R.2 = 0.21386579209539458
```

```
Entropy H.2 = 3.9307352773977775
Redundancy R.2 = 0.21385294452044445
```

Без пробілів

"о"	-->	0.1018578247065359
"е"	-->	0.07848984052400255
"т"	-->	0.06627971896850639
"и"	-->	0.061060402740910624
"а"	-->	0.060190050276148
"н"	-->	0.05349729169261166
"с"	-->	0.04978919839418571
"р"	-->	0.043814270702005026
"в"	-->	0.039622026354048746
"л"	-->	0.03165312396365789
"м"	-->	0.028401196104962827
"д"	-->	0.026867584687953455
"к"	-->	0.024908592001992576
"п"	-->	0.02300696981612055
"у"	-->	0.018171056442777125
"я"	-->	0.01663744502576775
"б"	-->	0.014632276083637794
"ы"	-->	0.014174711363416932
"ь"	-->	0.013907448789832268
"ч"	-->	0.012967132300571046
"з"	-->	0.012401822982622336
"г"	-->	0.01219472947010647
"й"	-->	0.009694215234526406
"ж"	-->	0.0075001434262502225
"х"	-->	0.007441373645671396
"ю"	-->	0.005830801801713559
"щ"	-->	0.004662402592586891
"ш"	-->	0.0043363702384234025
"ц"	-->	0.003500999785910085
"ф"	-->	0.0031693703097867075
"э"	-->	0.003057427870588943

Entropy H.1 = 3.966059345058023  
Redundancy R.1 = 0.1994545438696368

Entropy H.2 = 4.059183620300855  
Redundancy R.2 = 0.18065749397325115

Entropy H.2 = 4.057968090491926  
Redundancy R.2 = 0.18090284755489683

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

Лабораторная работа №1

Произвольная часть текста:  
если\_нет

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

Неравенство для энтропии:  
 $2.27008406910629 < H < 2.9964130274445$

Двоичная таблица угаданных символов:

Поле ввода символов:

Продолжить Другой

Вероятности:

$q[1] = 0.44$   
 $q[2] = 0.08$   
 $q[3] = 0.12$   
 $q[4] = 0$   
 $q[5] = 0.02$   
 $q[6] = 0.04$   
 $q[7] = 0.06$   
 $q[8] = 0.02$   
 $q[9] = 0.04$   
 $q[10] = 0$   
 $q[11] = 0$   
 $q[12] = 0$   
 $q[13] = 0.02$   
 $q[14] = 0$   
 $q[15] = 0.02$   
 $q[16] = 0.04$   
 $q[17] = 0.02$   
 $q[18] = 0$   
 $q[19] = 0$   
 $q[20] = 0$   
 $q[21] = 0$   
 $q[22] = 0$   
 $q[23] = 0$   
 $q[24] = 0.02$   
 $q[25] = 0$   
 $q[26] = 0$   
 $q[27] = 0$   
 $q[28] = 0$   
 $q[29] = 0.02$   
 $q[30] = 0$   
 $q[31] = 0.02$   
 $q[32] = 0.02$

Строка состояния:

Лабораторная работа №1

Произвольная часть текста:  
как\_хотели\_бы\_чтоб\_

Использованные буквы:  
к, р, а, п, м, т, я, ч, и, ф, ы,

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: ы

Символ по счету: 11

Номер эксперимента: 30

Неравенство для энтропии:  
 $3.42881954503892 < H < 3.52938358088685$

Двоичная таблица угаданных символов:

Поле ввода символов:

Продолжить Другой

Вероятности:

$q[1] = 0.3103448$   
 $q[2] = 0.0344827$   
 $q[3] = 0.0344827$   
 $q[4] = 0$   
 $q[5] = 0$   
 $q[6] = 0.0689655$   
 $q[7] = 0$   
 $q[8] = 0.0689655$   
 $q[9] = 0$   
 $q[10] = 0.034482$   
 $q[11] = 0.068965$   
 $q[12] = 0$   
 $q[13] = 0.034482$   
 $q[14] = 0$   
 $q[15] = 0.034482$   
 $q[16] = 0$   
 $q[17] = 0.068965$   
 $q[18] = 0$   
 $q[19] = 0$   
 $q[20] = 0.068965$   
 $q[21] = 0$   
 $q[22] = 0.034482$   
 $q[23] = 0$   
 $q[24] = 0$   
 $q[25] = 0$   
 $q[26] = 0$   
 $q[27] = 0.034482$   
 $q[28] = 0.034482$   
 $q[29] = 0$   
 $q[30] = 0.034482$   
 $q[31] = 0$   
 $q[32] = 0.034482$

Строка состояния:  
Вы не угадали. Введите другую букву

Лабораторная работа №1

Произвольная часть текста:  
о\_поведения\_с\_женой\_или\_мужем

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: 31

Неравенство для энтропии:  
2.57816057338474 < H < 2.94487046767526

Двоичная таблица угаданных символов:

Поле ввода символов:

Продолжить Другой

Вероятности:

q[1] = 0.4333333  
q[2] = 0.1  
q[3] = 0  
q[4] = 0.0666666  
q[5] = 0.0666666  
q[6] = 0  
q[7] = 0  
q[8] = 0  
q[9] = 0.0333333  
q[10] = 0  
q[11] = 0  
q[12] = 0  
q[13] = 0  
q[14] = 0  
q[15] = 0.0333333  
q[16] = 0  
q[17] = 0  
q[18] = 0.0333333  
q[19] = 0.0333333  
q[20] = 0  
q[21] = 0  
q[22] = 0  
q[23] = 0.0666666  
q[24] = 0.0333333  
q[25] = 0  
q[26] = 0.0333333  
q[27] = 0  
q[28] = 0  
q[29] = 0.0333333  
q[30] = 0.0333333  
q[31] = 0  
q[32] = 0

Строка состояния:

H(10)

2.2708406150629 < H < 2.9964130274445

H(20)

3.42881547503682 < H < 3.5238350888865

H(30)

2.57816507338474 < H < 2.94487046767526

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

Надлишковість джерела відкритого тексту (мови) дорівнює  $R = 1 - \frac{H_{\infty}}{H_0}$

$H_0 = \log_2(32) = 5$  (32 букви алфавіту, враховуючи пробіл)

Для H(10)

$$\begin{aligned} \bullet R_{\min} &= 1 - \frac{2.2708406150629}{5} = 1 - 0.45416812301258 = \\ &= 0.54583187698742 \\ \bullet R_{\max} &= 1 - \frac{2.9964130274445}{5} = 1 - 0.5992826054889 = 0.4007173945111 \end{aligned}$$

Для H(20)



$$\begin{aligned} \bullet R_{\min} &= 1 - \frac{3.42881547503682}{5} = 1 - 0.685763095007364 = \\ &0.314236904992636 \\ \bullet R_{\max} &= 1 - \frac{3.5238350888865}{5} = 1 - 0.7047670177773 = 0.2952329822227 \end{aligned}$$

Для H(30)

$$\begin{aligned} \bullet R_{\min} &= 1 - \frac{2.57816507338474}{5} = 1 - 0.515633014676948 = \\ &0.484366985323052 \\ \bullet R_{\max} &= 1 - \frac{2.94487046767526}{5} = 1 - 0.588974093535052 = \\ &0.411025906464948 \end{aligned}$$

### Висновки:

Під час виконання лабораторної роботи ознайомився із наступними термінами:  
 Ентропія - міра невизначеності, показує кількість необхідної інформації для однозначного опису ансамблю (пара з множини символів та множини ймовірностей).  
 Вирахувавши частоту появу символів алфавіту можна визначити ентропію тексту, тобто скільки інформації містить текст  
 Сукупна ентропія говорить нам, скільки інформації містять два ансамблі (із урахуванням взаємних залежностей)  
 Умовна ентропія говорить нам, скільки інформації залишиться в ансамблі X, якщо поведінку ансамблю Y буде однозначно визначено

Для Аналізу був взят уривок з Капіталу Карла Макса

Було встановлено, що у російській мові (без пробілу) найчастіше трапляються літери "о", "а", "е", "и", "н", "т", "с", "л", "в", "р". Також було встановлено частоту різних біграм.

Під час виконання лабораторної роботи, працюючи з програмою CoolPinkProgram, я також зіткнувся з поняттям умовної ентропії джерела, яке визначає, скільки інформації про наступний символ ми отримаємо з значень попередніх. Було оцінено умовну ентропію для n=10, 20, 30. Це дозволило розрахувати надлишковість російської мови (величину можливого стиснення тексту деякою схемою кодування символів без втрати його змісту).