

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

КРИПТОГРАФІЯ
“Комп’ютерний практикум №1”
“Експериментальна оцінка ентропії на символ джерела відкритого
тексту”

Виконала
студентка групи ФБ-93
Куцовол Онисія
Перевірила
Селюх П.В.

Київ 2021

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.

Завдання:

- Уважно прочитати методичні вказівки до виконання комп'ютерного практикуму.
- 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.
- За допомогою програми CoolPinkProgram оцінити значення $H(10)$, $H(20)$, $H(30)$.
- Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела

Хід роботи:

Необхідно написати програму, що буде обчислювати частоти літер та біграм в тексті, а також рахувати ентропії за їх означенням.

Найважчим під час виконання лабораторної роботи, було обрати, якою мовою програмувати. На момент отримання роботи, мені доводилося знайомитися лише з C++ на уроках програмування.

Однак, трохи більше ознайомившись з мовою Python, та зрозумівши, що програмування нею буде набагато швидшим та простішим, вирішила використовувати її. Однак, на цьому труднощі не

закінчилися. Через те, що довелося зіштовхнутися з не зовсім знайомою мовою виникали деякі труднощі з підрахунком біграм, однак врешті-решт, вирішила для біграм, що не перетинаються використати цикл, що з кроком 2 ітерується по всьому тексту та додає одиницю до лічильника.

Нарешті отримуємо необхідні для нас результати.

Частоти літер з пробілами: Частоти літер без пробілів:

а => 0.07144139591174475	а => 0.0858287659193685
б => 0.015616895639105744	б => 0.01875968744507119
в => 0.036480429838677504	в => 0.04382719994886627
г => 0.015869585455706135	г => 0.019067288793723335
д => 0.025089438895612506	д => 0.03014493216791039
е => 0.06920043622242024	е => 0.08313625541298478
ж => 0.008913965767179583	ж => 0.010710119684888384
з => 0.014117381069542898	з => 0.016962017225675526
и => 0.05361678924339349	и => 0.06442051101772103
й => 0.01180327432804458	й => 0.011489110113293172
к => 0.028211488077030496	к => 0.03389607069238267
л => 0.04120173956989533	л => 0.0495038430194468
м => 0.025418600630394594	м => 0.030540419616177435
н => 0.055152877339043234	н => 0.06626611910963391
о => 0.09317272013937838	о => 0.11195091161854237
п => 0.01990597278929659	п => 0.023917003563381858
р => 0.03467835246239576	р => 0.04166200604017194
с => 0.04465960021811121	с => 0.05365446381489589
т => 0.047921293771860995	т => 0.0575653952477589
у => 0.021797821547791623	у => 0.02619005768523993
ф => 0.0015128140335944461	ф => 0.001817644332944504
х => 0.0	х => 0.0
ц => 0.0	ц => 0.0026925105063837266
ч => 0.012554694045829953	ч => 0.015076460906664962
ш => 0.008458459124097298	ш => 0.010162828973650149
щ => 0.0023772791956484154	щ => 0.0028562982374842204
ы => 0.015763189743453337	ы => 0.018939454467010754
ь => 0.01627854397467782	ь => 0.019554657164315047
э => 0.0025169235679802105	э => 0.0030240807912944824
ю => 0.006290646486946576	ю => 0.007558204566881322
я => 0.018955061111037227	я => 0.022770489445678402
=> 0.16733053157957734	

Біграми з найвищими частотами:

- з пробілами і кроком 1

```
о : 0.020345
и : 0.019431
а : 0.019078
с : 0.017379
е : 0.017063
н : 0.016219
п : 0.014131
в : 0.014001
то : 0.013605
```

- з пробілами і кроком 2

```
о : 0.020435
и : 0.019583
а : 0.018493
с : 0.017502
е : 0.017130
н : 0.016126
п : 0.014284
в : 0.013811
то : 0.013599
```

- без пробілів з кроком 1:

```
то : 0.016702
на : 0.013331
ст : 0.013239
не : 0.011082
ал : 0.010898
он : 0.010670
ко : 0.010654
го : 0.010490
ла : 0.010395
```

- без пробілів з кроком 2:

```
то : 0.016555  
на : 0.013455  
ст : 0.013391  
не : 0.011465  
ал : 0.010938  
он : 0.010834  
го : 0.010762  
ко : 0.010698  
ов : 0.010163
```

Далі програма обчислює значення ентропії.

В тексті без пробілів:

```
Letter Entropia: 4.540594180895466  
Bigram 1-stepped: 4.180901854850425  
Bigram 2-stepped: 4.179679885683603
```

В тексті з пробілами:

```
Letter Entropia: 4.424816926674272  
Bigram 1-stepped: 3.9993347075470806  
Bigram 2-stepped: 3.9957401096092267
```

Другою частиною практикуму було дослідження значень $H(10)$, $H(20)$ та $H(30)$

H(10):

Лабораторная работа №1

Произвольная часть текста:
тогда_мы_так_равностно_оправдываем_свое_не_совсем_порядочное_поведение_прав

Использованные буквы:
й, ы, ц, в, у, а, е, п, м,

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: т

Символ по счету: 10

Номер эксперимента: 50

Неравенство для энтропии:
 $2,21377415524533 < H < 2,75663484475755$

Двоичная таблица угаданных символов:

1000000000000000000000000000000000	▲
0000001000000000000000000000000000	■
0100000000000000000000000000000000	■
1000000000000000000000000000000000	■
1000000000000000000000000000000000	▼

Поле ввода символов:
т

Продолжить Другой

Вероятности:

q[1] = 0,5
q[2] = 0,06
q[3] = 0,06
q[4] = 0
q[5] = 0,02
q[6] = 0,04
q[7] = 0,06
q[8] = 0
q[9] = 0
q[10] = 0,04
q[11] = 0
q[12] = 0
q[13] = 0
q[14] = 0
q[15] = 0
q[16] = 0
q[17] = 0
q[18] = 0,08
q[19] = 0,04
q[20] = 0,02
q[21] = 0,02
q[22] = 0,02
q[23] = 0
q[24] = 0,02
q[25] = 0
q[26] = 0
q[27] = 0
q[28] = 0,02
q[29] = 0
q[30] = 0
q[31] = 0
q[32] = 0

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

Надлишковість: 55.8% > R1 > 45%

H(20):

Произвольная часть текста:
ическим_телом_челов

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:
 $1.44790135131193 < H < 2.05945832095488$

Двоичная таблица угаданных символов:

10000000000000000000000000000000	▲
10000000000000000000000000000000	
10000000000000000000000000000000	
10000000000000000000000000000000	
10000000000000000000000000000000	
10000000000000000000000000000000	▼

Вероятности:

$q[1] = 0.66$
$q[2] = 0.02$
$q[3] = 0.06$
$q[4] = 0$
$q[5] = 0.04$
$q[6] = 0.02$
$q[7] = 0$
$q[8] = 0.04$
$q[9] = 0.02$
$q[10] = 0$
$q[11] = 0.02$
$q[12] = 0$
$q[13] = 0$
$q[14] = 0.04$
$q[15] = 0$
$q[16] = 0.04$
$q[17] = 0.02$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0$
$q[21] = 0$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0.02$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

Строка состояния:

Надлишковість: 71% > R2 > 41%

H(30):

Произвольная часть текста:
_человеческой_природы_то_кака

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:
 $1.58066075405678 < H < 2.21112946255746$

Двоичная таблица угаданных символов:

10000000000000000000000000000000	▲
10000000000000000000000000000000	
10000000000000000000000000000000	
10000000000000000000000000000000	
10000000000000000000000000000000	
01000000000000000000000000000000	▼

Вероятности:

$q[1] = 0.58$
$q[2] = 0.16$
$q[3] = 0.06$
$q[4] = 0$
$q[5] = 0.02$
$q[6] = 0$
$q[7] = 0$
$q[8] = 0$
$q[9] = 0$
$q[10] = 0$
$q[11] = 0$
$q[12] = 0.02$
$q[13] = 0.02$
$q[14] = 0$
$q[15] = 0$
$q[16] = 0.02$
$q[17] = 0$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0$
$q[21] = 0.02$
$q[22] = 0.02$
$q[23] = 0.04$
$q[24] = 0$
$q[25] = 0.02$
$q[26] = 0.02$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

Строка состояния:

Надлишковість: 68% > R > 56%

Висновок:

Під час виконання даної лабораторної роботи мені довелося ознайомитися не лише з ентропією, дослідженнями надлишковості мови та створеннями масивів, а також майже з самого початку вивчити нову мову програмування. Також корисним було закріплення роботи з файлами та з гітхабом.