

Національний технічний університет України «КПІ» імені
Ігоря Сікорського
Фізико-технічний інститут

Лабораторна робота 1
Криптографія

Виконали:

студенти ФБ-14

Кот Микита Сергійович

Чавалах Артем Дмитрович

Перевірила:

Селюх П. В.

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.

Частоти букв та біграм були збережені в файлах .xlsx формату.

Ентропія і надлишковість для букв та непересічних біграм у тексті без пробілів:

```
Letters and not cross bigrams without space:  
Letters entropy: 4.458502430665178  
Bigrams entropy: 4.128663354359083  
Letter redundancy: 0.10005535684616174  
Bigrams redundancy: 0.16663307311763065
```

Ентропія і надлишковість для букв та непересічних біграм у тексті з пробілами:

```
Letters and not cross bigrams with space:  
Letters entropy: 4.391930259750691  
Bigrams entropy: 3.971846934898919  
Letter redundancy: 0.12161394804986192  
Bigrams redundancy: 0.20563061302021624
```

Ентропія і надлишковість для пересічних біграм у тексті без пробілів:

```
Cross bigrams without space:  
Bigrams entropy: 4.129496081270868  
Bigrams redundancy: 0.16646498795111453
```

Ентропія і надлишковість для пересічних біграм у тексті з пробілами:

```
Cross bigrams with space:    
Bigrams entropy: 3.971768447179969  
Bigrams redundancy: 0.20564631056400617
```

2. За допомогою програми CoolPinkProgram оцінити значення H^{10} , H^{20} , H^{30} .

Произвольная часть текста:
реальную_значимость_порядочного_поведения_почему_тогда_мы_так_ревностно_оп

Использованные буквы:
а,

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

Введенный символ: (пробел)
Символ по счету: 2
Номер эксперимента: 50

Неравенство для энтропии:
 $2.40496014495526 < H < 3.2015968125891$

Двоичная таблица угаданных символов:

00000000000000010000000000000000
00000000001000000000000000000000
00000000000000010000000000000000
00000000000000000000000000000000
00000000010000000000000000000000
00000000000000000000000000000000

Вероятности:

q[1] = 0.42
q[2] = 0.1
q[3] = 0.06
q[4] = 0.04
q[5] = 0
q[6] = 0.06
q[7] = 0
q[8] = 0.06
q[9] = 0.04
q[10] = 0.02
q[11] = 0
q[12] = 0.02
q[13] = 0.02
q[14] = 0.02
q[15] = 0
q[16] = 0.02
q[17] = 0
q[18] = 0.02
q[19] = 0
q[20] = 0.02
q[21] = 0
q[22] = 0
q[23] = 0.02
q[24] = 0
q[25] = 0
q[26] = 0
q[27] = 0
q[28] = 0
q[29] = 0
q[30] = 0.02
q[31] = 0.02
q[32] = 0.02

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

$2.40496014495526 < H^{10} < 3.2015968125891$

Произвольная часть текста:
ему_то_не_успеете_

Использованные буквы:

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

Введенный символ:
Символ по счету:
Номер эксперимента: 51

Неравенство для энтропии:
 $2.13655490598199 < H < 2.9479764393853$

Двоичная таблица угаданных символов:

10000000000000000000000000000000
10000000000000000000000000000000
00000000000100000000000000000000
00000000000000000000000000000000
10000000000000000000000000000000
00000000000000000000000000000000

Вероятности:

q[1] = 0.43
q[2] = 0.08
q[3] = 0.04
q[4] = 0.06
q[5] = 0.02
q[6] = 0
q[7] = 0.04
q[8] = 0.04
q[9] = 0.02
q[10] = 0.02
q[11] = 0.02
q[12] = 0.04
q[13] = 0
q[14] = 0.02
q[15] = 0.04
q[16] = 0
q[17] = 0.02
q[18] = 0
q[19] = 0
q[20] = 0.02
q[21] = 0
q[22] = 0
q[23] = 0
q[24] = 0.04
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:

$2.13655490598199 < H^{20} < 2.9479764393853$

Произвольная часть текста:
обра_если_бы_они_не_имели_пре

Использованные буквы:

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

Введенный символ:

Символ по счету:

Номер эксперимента: 51

Неравенство для энтропии:
 $1.48465845168589 < H < 2.31092368847664$

Двоичная таблица угаданных символов:

Поле ввода символов:

Вероятности:

$q[1] = 0.54$
$q[2] = 0.16$
$q[3] = 0.06$
$q[4] = 0.06$
$q[5] = 0.02$
$q[6] = 0.02$
$q[7] = 0.02$
$q[8] = 0.02$
$q[9] = 0$
$q[10] = 0.06$
$q[11] = 0$
$q[12] = 0.02$
$q[13] = 0$
$q[14] = 0.02$
$q[15] = 0$
$q[16] = 0$
$q[17] = 0$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0$
$q[21] = 0$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

Строка состояния:

$$1.48465845168589 < H^{30} < 2.31092368847664$$

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

Надлишковість джерела відкритого тексту (мови): $R = 1 - \frac{H_{\infty}}{H_0}$; $H_0 = \log_2 32 = 5$

$$R \text{ для } H^{10}: 1 - \frac{2.40496014495526}{5} = 0.51900; 1 - \frac{3.2015968125891}{5} = 0.35968$$

$$0.35968 < R < 0.51900$$

$$2) R \text{ для } H^{20}: 1 - \frac{2.13655490598199}{5} = 0.57268; 1 - \frac{2.9479764393853}{5} = 0.41040$$

$$0.41040 < R < 0.57268$$

$$3) R \text{ для } H^{30}: 1 - \frac{1.48465845168589}{5} = 0.70306; 1 - \frac{2.31092368847664}{5} = 0.53781$$

$$0.53781 < R < 0.70306$$

Труднощі

У ході виконання лабораторної роботи зіткнулися із труднощами з обрахунками ентропії для біграм, вона була у два рази більша або у два рази менша від нормальних значень. Також спочатку було не зовсім зрозуміло поняття ентропії, але після виконання завдання з CoolPinkProgram, усе стало більш зрозуміло.

Висновки

У ході виконання лабораторної роботи, ми набули практичних навичків щодо оцінення ентропії, надлишковості. Покращили свої знання роботи з текстом. На реальному прикладі оцінили значення ентропії за допомогою програми CoolPinkProgram. Виявили, що буква «о» трапляється найбільше всього, приблизно один раз кожні десять символів.