

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

КРИПТОГРАФІЯ КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1
**Експериментальна оцінка ентропії на символ джерела
відкритого тексту**

Виконали студенти :
3го курсу групи ФБ-04
Музичка-Скрипка Олександра Тарасівна
Кузьмін Гліб Ігорович

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 $1H$ та $2H$ за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення $1H$ та $2H$ на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення $1H$ та $2H$ на тому ж тексті, в якому видалено 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

Ми взяли текст з книги «Война и мир» 1ший том.

Результат:

Текст з пробілами:

Ентропія для монограм: 4.389513128998312

Надлишковість для монограм: 0.12209737420033773

Ентропія для біграм з перетином: 3.984551695044777

Надлишковість для біграм з перетином: 0.2030896609910464

Ентропія для біграм без перетину: 3.9839081215820165

Надлишковість для біграм без перетину: 0.20321837568359669

Текст без пробілів:

Ентропія для монограм: 4.46859115010279

Надлишковість для монограм: 0.09801895804289695

Ентропія для біграм з перетином: 4.149796060796719

Надлишковість для біграм з перетином: 0.16236745562618626

Ентропія для біграм без перетину: 4.150318718889217

Надлишковість для біграм без перетину: 0.16226195756762085

Монограми без пробілів

О	Кількість	Частота
о	61232	0,113683
а	45176	0,083874
е	42904	0,079656
и	35789	0,066446
н	35097	0,065161
т	30591	0,056795
с	28103	0,052176
л	27263	0,050616
в	24795	0,046034
р	24545	0,04557

Монограми з пробілами

О	Кількість	Частота
	102010	0,159234
о	61232	0,095581
а	45176	0,070518
е	42904	0,066972
и	35789	0,055865
н	35097	0,054785
т	30591	0,047752
с	28103	0,043868
л	27263	0,042557
в	24795	0,038704

Біграми:

З пробілами та з перетинами

О	Кількість	Частота
о	13207	0,020616
и	11287	0,017619
а	10404	0,01624
е	10309	0,016092
с	9791	0,015283
п	9648	0,01506
в	9451	0,014753
н	9221	0,014394
то	8490	0,013253
о	7533	0,011759

З пробілами та без перетинів

О	Кількість	Частота
о	6641	0,020733
и	5624	0,017558
а	5291	0,016518
е	5180	0,016172
с	4906	0,015316
п	4808	0,01501
в	4756	0,014848
н	4616	0,014411
то	4213	0,013153
о	3792	0,011838

Без пробілами та з перетинами

О	Кількість	Частота
то	8656	0,016071
ст	6836	0,012692
на	6582	0,01222
ов	6482	0,012035
ал	5855	0,01087
го	5799	0,010766
он	5618	0,01043
не	5579	0,010358
но	5541	0,010287
ко	5489	0,010191

Без пробілами та без перетинів

О	Кількість	Частота
то	4322	0,016048
ст	3438	0,012766
ов	3251	0,012072
на	3236	0,012016
ал	2926	0,010865
го	2898	0,010761
не	2862	0,010627
он	2774	0,0103
ко	2757	0,010237
ос	2717	0,010089

Значення $H(10)$:

Произвольная часть текста:		Вероятности:
e_srazc_jk		q[1] = 0,52 q[2] = 0,1 q[3] = 0,06 q[4] = 0,02 q[5] = 0,02 q[6] = 0,04 q[7] = 0,02 q[8] = 0,02 q[9] = 0,02 q[10] = 0,02 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 q[22] = 0 q[23] = 0,02 q[24] = 0,02 q[25] = 0,02 q[26] = 0,02 q[27] = 0 q[28] = 0 q[29] = 0 q[30] = 0,02 q[31] = 0 q[32] = 0
Использованные буквы:		
Порядок n-граммы:	Введенный символ:	Неравенство для энтропии:
5 символов		$1,93014012301981 < H < 2,83233697680935$
10 символов	Символ по счету:	
15 символов		Двоичная таблица угаданных символов:
20 символов	Номер эксперимента: 51	10000000000000000000000000000000 ^ 00000000000000000000000001000000000 100000000000000000000000000000000 100000000000000000000000000000000 100000000000000000000000000000000 v
35 символов	Поле ввода символов:	
40 символов		
45 символов	<button>Продолжить</button> <button>Другой</button>	
50 символов		
Строка состояния:		

R = 0,523752290017084

Значення $H(20)$:

Лабораторная работа №1

Произвольная часть текста:
обобщил_в_одной_из_моих_книг_под_названием_человек_отменяется_но_в_данный_

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: 50

Поле ввода символов:

Продолжить
Другой

Неравенство для энтропии:
 $2,03884855639946 < H < 2,83663484475755$

Двоичная таблица угаданных символов:

00000010000000000000000000000000	▲
000000000000000000000000010000000000	■
100000000000000000000000000000000000	
100000000000000000000000000000000000	
100000000000000000000000000000000000	▼

Вероятности:

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

R = 0,512451659884299

Значення $H(30)$:

Лабораторная работа №1

Произвольная часть текста:
поведения_почему_тогда_мы_так_ревностно_оправдываем_свое_не_совсем_порядоч

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: к

Символ по счету: 1

Номер эксперимента: 50

Неравенство для энтропии:
 $1,63950282335213 < H < 2,38904103608174$

Двоичная таблица угаданных символов:

10000000000000000000000000000000	▲
10000000000000000000000000000000	■
00010000000000000000000000000000	
00100000000000000000000000000000	
00000000000000000000000000000000	
00000000000000000000000000000000	

Поле ввода символов:
к

Продолжить Другой

Вероятности:

$q[1] = 0,52$
$q[2] = 0,18$
$q[3] = 0,08$
$q[4] = 0,02$
$q[5] = 0,02$
$q[6] = 0$
$q[7] = 0,04$
$q[8] = 0$
$q[9] = 0,02$
$q[10] = 0$
$q[11] = 0$
$q[12] = 0,02$
$q[13] = 0$
$q[14] = 0$
$q[15] = 0$
$q[16] = 0,02$
$q[17] = 0,04$
$q[18] = 0$
$q[19] = 0$
$q[20] = 0,02$
$q[21] = 0$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0$
$q[28] = 0,02$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0$
$q[32] = 0$

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

$R = 0,597145614056613$

Висновки:

В ході виконання лабораторної роботи ми закріпили теоретичний матеріал, а саме: що таке ентропія на символі джерела, надлишковість, біграми, дослідили на прикладі першого тому «Война и мир» наближене значення ентропії.

Під час аналізу текстів, ми підтвердили факт, що в російському алфавіті найчастіше зустрічається пробіл. З цього можна зробити висновки, що при шифруванні слід уникати цього символу, для ускладнення зламу шифрування.