

НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
“КІЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ”
ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ

Криптографія

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1
«Експериментальна оцінка ентропії на символ джерела
відкритого тексту»

ФБ-32 Дорошенко Ілля

Мета: 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.
2. За допомогою програми CoolPinkProgram оцінити значення (10) H_1 , (20) H_2 , (30) $R(H)$.
3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Хід роботи:

Завдання 1:

Для дослідження було обрано текст російською мовою «TEXT.txt». Згідно з методичними вказівками, текст пройшов попередню фільтрацію:

- Усі символи, крім літер, були вилучені або замінені на пробіли.
- Прописні літери замінені на відповідні рядкові.
- Послідовності пробілів трактуються як один пробіл.
- Буква «ё» замінена на «е», а «ъ» на «ъ».
- Алфавіт дослідження склав 32 літери (без пробілу) або 33 символи (з пробілом).

При підрахунку біграм було реалізовано два підходи:

1. **Крок 1:** пари букв, що перетинаються (більш точна статистика).
2. **Крок 2:** пари букв, що не перетинаються.

На основі роботи програми було отримано такі значення ентропії та надлишковості:

A	B	C
Параметр	З пробілами	Без пробілів
H_1	4,383006841	4,468568723
H_2 (step 1)	3,975230577	4,150663708
H_2 (step 2)	3,974254627	4,149794329
$R(H_1)$	0,123398632	0,098023485
$R(H_2$ step 1)	0,204953885	0,162192322

Таблиці частот символів:

(З пробілом)

	A	B
1	Символ	Частота
2		0,162346
3	о	0,095222
4	а	0,070264
5	е	0,066722
6	и	0,055668
7	н	0,054587
8	т	0,047573
9	с	0,043704
10	л	0,042398
11	в	0,03856
12	р	0,038175
13	к	0,030044
14	д	0,025471
15	м	0,024757
16	у	0,024013
17	п	0,021519
18	я	0,019391
19	г	0,017368
20	б	0,016753
21	ы	0,015897
22	з	0,014918
23	б	0,014467
24	ч	0,011416
25	й	0,009648
26	ж	0,008485
27	ш	0,00791
28	х	0,007145
29	ю	0,005433
30	ц	0,003388
31	э	0,002533
32	щ	0,002349
33	ф	0,001875
..		

(Без пробілу)

	A	B
1	Символ	Частота
2	о	0,113677
3	а	0,083882
4	е	0,079653
5	и	0,066457
6	н	0,065167
7	т	0,056793
8	с	0,052174
9	л	0,050615
10	в	0,046034
11	р	0,045574
12	к	0,035867
13	д	0,030408
14	м	0,029556
15	у	0,028667
16	п	0,025689
17	я	0,02315
18	г	0,020735
19	ь	0,02
20	ы	0,018979
21	з	0,017809
22	б	0,017271
23	ч	0,013629
24	й	0,011518
25	ж	0,01013
26	ш	0,009443
27	х	0,00853
28	ю	0,006486
29	ц	0,004045
30	э	0,003024
31	щ	0,002805
32	ф	0,002239
..		

Частота біграм з перекриванням з пробілами:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	Y	Z	AA	AB	AC	AD	AE	AF	AG
1	а	б	в	г	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и			
2	а	б	в	г	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и			
3	б	в	г	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и			
4	в	г	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и			
5	г	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и			
6	д	е	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и			
7	ж	з	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и			
8	и	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и			
9	л	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
10	м	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
11	н	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
12	о	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
13	р	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
14	с	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
15	т	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
16	у	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
17	ф	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
18	х	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
19	ч	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
20	ш	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
21	и	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
22	м	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
23	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
24	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
25	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
26	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
27	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
28	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
29	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
30	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
31	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
32	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			
33	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и	и			

Без пробілів:

Завдання 2:

Для проведення експерименту використовувалася програма CoolPinkProgram та очищений російський текст.

Було проведено по 30 експериментів для кожного порядку n-грами: n=10, 20, 30.

Експерименти з C++ в MikroProgram.

III

Page 1

H20

Лабораторная работа №1

Произвольная часть текста:
ай_ты_же_обещал_каждый_день_люди_произносят_подобное_как_образованные_так_и

Использованные буквы:

Порядок n-грамм:

- 5 символов
- 10 символов
- 15 символов
- 20 символов**
- 25 символов
- 30 символов
- 35 символов
- 40 символов
- 45 символов
- 50 символов

Введенный символ: д

Символ по счету: 1

Номер эксперимента: 51

Неравенство для энтропии:
 $3.06958759425316 < H < 3.82409862615796$

Поле ввода символов:

д

Двоичная таблица угаданных символов:

000000000000000000000000000000001000
000000000000000000000000000000000000
100000000000000000000000000000000000
100000000000000000000000000000000000
000000001000000000000000000000000000
.....

Продолжить Другой

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

Вероятности:

$q[1] = 0.3137254$
$q[2] = 0.0784313$
$q[3] = 0.0196078$
$q[4] = 0.0392156$
$q[5] = 0.0588235$
$q[6] = 0.0392156$
$q[7] = 0.0196078$
$q[8] = 0.0196078$
$q[9] = 0.0392156$
$q[10] = 0$
$q[11] = 0.0392156$
$q[12] = 0.0588235$
$q[13] = 0.0588235$
$q[14] = 0.019607$
$q[15] = 0.019607$
$q[16] = 0$
$q[17] = 0.019607$
$q[18] = 0.019607$
$q[19] = 0.019607$
$q[20] = 0.019607$
$q[21] = 0.019607$
$q[22] = 0.019607$
$q[23] = 0$
$q[24] = 0.019607$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0.019607$
$q[31] = 0.019607$
$q[32] = 0$

H30

Лабораторная работа №1

Произвольная часть текста:
вившуюся_женщину_вы_не_имеете_права_разного_мнения_держались_люди_и_по_тому

Использованные буквы:

Порядок n-грамм:

- 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:
 $2.94910685462945 < H < 3.6573290296138$

Поле ввода символов:

Двоичная таблица угаданных символов:

001000000000000000000000000000000000
010000000000000000000000000000000000
000000000010000000000000000000000000
000000000000000000000000000000000000
100000000000000000000000000000000000
000000000000000000000000000000000000
.....

Продолжить Другой

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

Вероятности:

$q[1] = 0.3269230$
$q[2] = 0.0769230$
$q[3] = 0.0769230$
$q[4] = 0.0384615$
$q[5] = 0.0384615$
$q[6] = 0.0192307$
$q[7] = 0.0384615$
$q[8] = 0.0192307$
$q[9] = 0.0384615$
$q[10] = 0$
$q[11] = 0.038461$
$q[12] = 0.057692$
$q[13] = 0$
$q[14] = 0$
$q[15] = 0.019230$
$q[16] = 0.019230$
$q[17] = 0.019230$
$q[18] = 0$
$q[19] = 0.019230$
$q[20] = 0$
$q[21] = 0$
$q[22] = 0.038461$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0.038461$
$q[26] = 0.038461$
$q[27] = 0$
$q[28] = 0.019230$
$q[29] = 0$
$q[30] = 0.019230$
$q[31] = 0$
$q[32] = 0$