



Міністерство освіти і науки, молоді та спорту України

Національний технічний університет України

“Київський політехнічний інститут”

Фізико-Технічний інститут

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1
за семестровий курс предмету
«Симетрична криптографія»

Роботу виконали:

Студенти групи ФІ-03

Піжук Богдан

Швець Катерина

Приймав:

Чорний Олег Миколайович

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

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. Підраховувати 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 $H^{(10)}, H^{(20)}, H^{(30)}$.
3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Хід роботи

- Очистка вхідного текстового файлу та редагування;
- Написання функцій для підрахунку частот літер та частот біграм. Визначили частоти та значення ентропії;
- Результати оформили у вигляді таблиць (див. нижче);
- Було проведено по 50 експериментів для умовних ентропій джерел;
- Обрахували оцінки надлишковості російської мови для даних умовних ентропій джерел;
- Оформлення звіту.

Таблиці частот букв і біграм тексту, одержані значення Н1 та Н2

	Текст з пробілами			Текст без пробілів		
	Літери	Біграми з ъ	Біграми без ъ	Літери	Біграми з ъ	Біграми без ъ
H	4.35676410 5666061	3.970778960 772621	3.9707789607 72621	4.51311878 1529859	4.191947557 3122496	4.1919438062 70274
R	0.12864717 886678778	0.205844207 84547575	0.2058442078 4547575	0.08903109 631167849	0.153859214 55645768	0.1538599717 008542

Текст з пробілами				Текст без пробілів			
Літера	Частота	Літера	Частота	Літера	Частота	Літера	Частота
—	0.1719224 18205326 42	я	0.0182418 18200896 088	о	0.1126614 90145071 01	ь	0.0196366 76946528 837
0	0.0951270 04279916 56	ь	0.0165804 50423036 7	е	0.0813193 01908298 66	г	0.0195883 49857745 76
е	0.0686628 72918772 53	г	0.0165396 44898668 226	а	0.0811667 26385140 66	ы	0.0189083 18679869 626
а	0.0685340 44048980 61	ы	0.0159654 52877197 524	и	0.0656689 19399390 94	б	0.0175427 33228256 436
и	0.0554482 95320072 706	б	0.0148124 05345756 867	н	0.0638435 36217355 92	з	0.0168792 71337963 07
н	0.0539070 12371069 12	з	0.0142522 03789783 929	т	0.0570439 14825577 184	ч	0.0140473 03935274 839
т	0.0481656 75092424 51	ч	0.0118610 00061791 223	с	0.0523009 56255009 62	й	0.0113665 31281779 375
с	0.0441609 04343689 78	й	0.0095974 59331465 606	л	0.0484506 68052963 726	ж	0.0100430 59436105 721
л	0.0409098 69923647 035	ж	0.0084799 70899831 765	в	0.0449869 96561182 44	ш	0.0091766 23772923 437
в	0.0379852 79698552	ш	0.0077483 86141511	р	0.0429047 89421614	х	0.0086892 10563196

	1		227		5		988
р	0.0362271 44534333 185	х	0.0073368 33281452 024	к	0.0332007 09993972 92	ю	0.0062072 69360694 751
к	0.0280333 95241143 16	ю	0.0052411 78137099 567	м	0.0304419 23611441 92	ц	0.0035782 75730895 439
м	0.0257039 82735765 577	ц	0.0030213 57611454 4605	д	0.0298868 52477419 17	э	0.0030867 20199273 2986
д	0.0252353 02141590 508	э	0.0026063 07135020 8166	у	0.0267904 66860389 255	щ	0.0029396 67771976 226
у	0.0226208 33901695 996	щ	0.0024821 41753728 1678	п	0.0248919 02658197	ф	0.0018612 83305131 025
п	0.0210177 59730077 287	ф	0.0015715 95624248 741	я	0.0216042 79846982 63		

-Біграми (дуже багато, в таблиці наведено 5 найчастіших)

Текст з пробілами				Текст без пробілів			
з перетином		без перетину		з перетином		без перетину	
'o_ '	0.0210340 94201356 726	'o_ '	0.0210789 6801663	'то'	0.0163097 13337319 17	'то'	0.0164008 44481940 156
'и_ '	0.0185408 75209055 566	'и_ '	0.0185070 54109291 186	'ст'	0.0132057 31320531 709	'ст'	0.0131007 92426725 728
'a_ '	0.0167454 31090213 488	'a_ '	0.0166078 48417969 82	'на'	0.0122688 75538329 62	'на'	0.0122654 23601033 372
'e_ '	0.0165967 82307647 016	'e_ '	0.0164772 70739990 696	'ов'	0.0111621 84441152 063	'ов'	0.0112284 61637240 051
'_c'	0.0160616 46690407 723	'_c'	0.0158511 97408965 79	'го'	0.0106250 62997855 657	'но'	0.0105725 93550952 665

Оцінки для $H^{(10)}$, $H^{(20)}$, $H^{(30)}$,

• $H^{(10)}$

$$1,68046891675367 < H^{(10)} < 2,49173259480082$$

$$R=1 - \frac{1,68046891675367+2,49173259480082}{2} * 1/\log_2 32 = 0,58278$$

Лабораторная работа №1

Произвольная часть текста:
юсь_сказа

Использованные буквы:

Порядок n-граммы:
5 символов

Введенный символ:

Символ по счету:

Номер эксперимента: 51

Неравенство для энтропии:
 $1,68046891675367 < H < 2,49173259480082$

Двоичная таблица угаданных символов:

01000000000000000000000000000000
00000000000001000000000000000000
01000000000000000000000000000000
00000010000000000000000000000000
10000000000000000000000000000000

Вероятности:

q[1] = 0,5
q[2] = 0,16
q[3] = 0,08
q[4] = 0,06
q[5] = 0,02
q[6] = 0
q[7] = 0,06
q[8] = 0
q[9] = 0
q[10] = 0,02
q[11] = 0,02
q[12] = 0,02
q[13] = 0,02
q[14] = 0
q[15] = 0
q[16] = 0
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] = 0,02
q[28] = 0
q[29] = 0
q[30] = 0
q[31] = 0
q[32] = 0

Строка состояния:

• $H^{(20)}$

$$1,63303467206811 < H^{(20)} < 2,32934912384719$$

$$R=1 - \frac{1,63303467206811+2,32934912384719}{2} * 1/\log_2 32 = 0,603762$$

Лабораторная работа №1

Произвольная часть текста:
изические_тела_подч

Использованные буквы:

Порядок n-граммы:
10 символов

Введенный символ:

Символ по счету:

Номер эксперимента: 51

Неравенство для энтропии:
 $1,63303467206811 < H < 2,32934912384719$

Двоичная таблица угаданных символов:

01000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000
10000000000000000000000000000000

Вероятности:

q[1] = 0,56
q[2] = 0,1
q[3] = 0,08
q[4] = 0,06
q[5] = 0
q[6] = 0
q[7] = 0
q[8] = 0,06
q[9] = 0,02
q[10] = 0
q[11] = 0
q[12] = 0,04
q[13] = 0
q[14] = 0
q[15] = 0,02
q[16] = 0
q[17] = 0
q[18] = 0
q[19] = 0
q[20] = 0
q[21] = 0
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:

Висновки: відомо, що пробіл є найвживанішим символом в російській мові. В даному комп'ютерному практикумі ми розглянули над одним текстом два випадки (з пробілом та без пробілу) та порівняли оцінки ентропії та надлишковості в літерах та біграмах(з перетином та без). Переконалися на власному дослідженні, що найчастіший символ є “ “, а найчастіші літери (“о”, “е”, “а”). Навчилися оцінювати ентропію та надлишковість на символ джерала.