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: Дудченко И.В і Терпило С.Е

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

Уважно прочитати методичні вказівки до виконання комп'ютерного практикуму.

- 1. Написати програми для підрахунку частот букв і частот біграм в тексті, а також підрахунку 1 H та 2 H за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення 1 H та 2 H на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення 1 H та 2 H на тому ж тексті, в якому вилучено 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

Перед виконанням роботи були розглянуті теоретичні відомості в методичних вказівках. В якості експериментального тексту була взята книга «Good Omens» в перекладі на російську мову. Оригінал тексту можна знайти у файлі badtxt.txt. Відредагований текст з пробілами міститься у файлі spaces.txt, а без пробілів у nospaces.txt. В ході виконання роботи було прийнято рішення пробіли замінити на «_» для кращого сприйняття. Усі таблиці також наведені у файлах з відповідними назвами. Єдина відмінність - таблиці біграм були виконанні у двох варіантах(таблиці у вигляді матриці і звичайні).

Аналіз результатів для тексту з пробілами

Текст без пробілів

Текст з пробілами

_	0.001205
a	0,081296
6	0,016728
В	0,046564
Г	0,018979
Д	0,031136
e	0,087972
ë	0
э	0,003554
ж	0,011568
3	0,017309
и	0,063784
ы	0,017417
й	0,010371
к	0,0328
л	0,046495
M	0,03076
н	0,067019
О	0,111932
п	0,027224
p	0,039487
С	0,053293
т	0,063036
у	0,026203
ф	0,002211
x	0,007079
ц	0,003261
ч	0,018446
ш	0,008342
щ	0,002955
ъ	0,00022
ь	0,022638
ю	0,005954
я	0,023967

1	а	0,067572
2	6	0,013904
3	В	0,038703
4	Г	0,015775
5	Д	0,02588
6	e	0,073121
7	ë	0
8	э	0,002954
9	ж	0,009615
LO	3	0,014387
Ι1	И	0,053017
l2	ы	0,014476
L3	й	0,00862
L4	к	0,027263
L 5	л	0,038646
L6	M	0,025567
L7	н	0,055705
18	О	0,093036
L9	п	0,022628
20	p	0,032821
21	С	0,044296
22	Т	0,052395
23	у	0,021779
24	ф	0,001838
25	x	0,005884
26	ц	0,002711
27	ч	0,015332
28	Ш	0,006934
29	щ	0,002456
30	ъ	0,000183
31	ь	0,018817
32	ю	0,004949
33	я	0,019921
34		0,168816

Біграми

Текст з пробілами

_	01101 9	p = =	
L	aa	0	
2	аб	0,00040969	
}	ав	0,00288247	
ļ	аг	0,00075354	
,	ад	0,00175216	
j	ae	0,00127663	
7	aë	0	
3	аэ	0	
)	аж	0,0017229	
0	аз	0,0032446	
1	аи	0,00022314	
2	аы	0	
3	ай	0,00077549	
4	ак	0,00540645	
5	ал	0,00805481	
6	ам	0,003219	
7	ан	0,00340555	
8	ao	0,00001097	
9	ап	0,00069135	
0	ар	0,00210332	
1	ac	0,00544303	
2	ат	0,0037933	
3	ay	0,00006584	
4	аф	0,00021216	
5	ax	0,00069867	
6	ац	0,00002561	
7	ач	0,00089254	
8	аш	0,00079743	
9	ащ	0,00016827	
0	аъ	0	
1	аь	0	
2	аю	0,00076451	
3	ая	0,00187287	
4	a	0,01660345	
5	ба	0,00062185	
6	66	0	
7	бв	0,00002561	
8	бг	0,00001097	

Текст без пробілів

1	aa	0,00068654
2	аб	0,00105182
3	ав	0,00584
4	аг	0,00149191
5	ад	0,00308944
6	ae	0,00206843
7	aë	0
8	аэ	0,00030806
9	аж	0,00221806
10	аз	0,00419406
11	аи	0,00170315
12	аы	0
13	ай	0,00092419
14	ак	0,00758716
15	ал	0,00991964
16	ам	0,00468696
17	ан	0,00657935
18	ao	0,00120145
19	ап	0,00238969
20	ар	0,0029266
21	ac	0,0082473
22	ат	0,00560675
23	ay	0,00061173
24	аф	0,00033447
25	ax	0,00091979
26	ац	0,00009242
27	ач	0,00169875
28	аш	0,00105622
29	ащ	0,00024645
30	аъ	0
31	аь	0
32	аю	0,0009638
33	ая	0,00250411
34	ба	0,00071295
35	66	0
36	бв	0,00004841
37	бг	0,0000088
38	бд	0,00002641

Перехрестні біграми

		~ •
Текст	3	пробілами

	Tekci	з прооглами
1	aa	0
2	аб	0,00040786
3	ав	0,00303793
4	аг	0,00077732
5	ад	0,00173936
6	ae	0,00129492
7	aë	0
8	аэ	0
9	аж	0,00166803
LO	аз	0,00313853
1	аи	0,00025057
2	аы	0
13	ай	0,00081572
4	ак	0,00534428
.5	ал	0,00800727
16	ам	0,00327936
7	ан	0,00350067
8.	ao	0,00000732
9	ап	0,00063648
20	ар	0,00202285
21	ac	0,00538817
22	ат	0,00381891
23	ay	0,00006036
24	аф	0,00025606
25	ax	0,00068953
26	ац	0,00002926
27	ач	0,0008962
28	аш	0,00081024
29	ащ	0,00018107
30	аъ	0
31	аь	0
32	аю	0,00083219
33	ая	0,00186373
34	a	0,01681747
35	ба	0,0005871
36	66	0
37	бв	0,00003658
8	бг	0,00000732

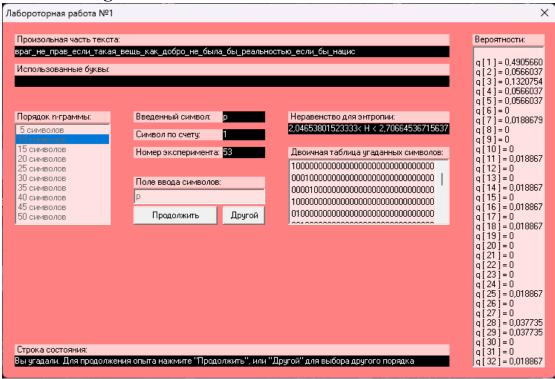
Текст без пробілів

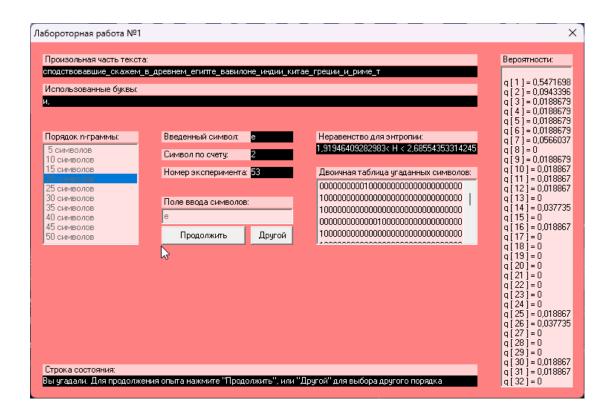
1 аа 0,00066894 2 аб 0,00112443 3 ав 0,00573218 4 аг 0,00306303 6 ае 0,00211684 7 аё 0 8 аэ 0,00219385 10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,0058752 24 аф 0,00058752 24 аф 0,00058752 25 ах 0,0009792 26 ац 0,0007168115 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,0002665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,0000088 38 бд 0,000022	_	enci ocs	проонив
3 ав 0,00573218 4 аг 0,00156012 5 ад 0,00306303 6 ае 0,00211684 7 аё 0 8 аэ 0,00219385 10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 25 ах 0,0009792 26 ац 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,000066 36 бв 0,000066 36 бв 0,0000088	1	aa	0,00066894
4 аг 0,00156012 5 ад 0,00306303 6 ае 0,00211684 7 аё 0 8 аэ 0,00219385 10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 24 аф 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,00005061 37 бг 0,0000088	2	аб	0,00112443
5 ад 0,00306303 6 ае 0,00211684 7 аё 0 8 аэ 0,00038728 9 аж 0,00219385 10 аз 0,00419626 11 ай 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,000586752 24 аф 0,00058647 25 ах 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,0002665 30 аъ 0 31 аь <td>3</td> <td>ав</td> <td>0,00573218</td>	3	ав	0,00573218
6 ае 0,00211684 7 аё 0 8 аэ 0,00038728 9 аж 0,00219385 10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,000066 36 бв 0,000066 36 бв 0,000068	4	аг	0,00156012
7 аё 0,00038728 9 аж 0,00219385 10 аз 0,00419626 11 ай 0,00165474 12 аы 0 13 ай 0,009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,000066 36 бв 0,0000066 36 бв 0,0000088	5	ад	0,00306303
8 аэ 0,00038728 9 аж 0,00219385 10 аз 0,00419626 11 ай 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,000066 36 бв 0,0000066 36 бв 0,0000088	6	ae	0,00211684
9 аж 0,00219385 10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,0000066 36 бв 0,0000088	7	aë	0
10 аз 0,00419626 11 аи 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 25 ах 0,0009792 26 ац 0,00071261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,0000088	8	аэ	0,00038728
11 аи 0,00165474 12 аы 0 13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,000066 36 бв 0,0000066 36 бв 0,0000088	9	аж	0,00219385
12 аы 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00071261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000068	10	аз	0,00419626
13 ай 0,0009814 14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00058752 25 ах 0,0009792 26 ац 0,0007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,0000088	11	аи	0,00165474
14 ак 0,00753655 15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,0000088	12	аы	0
15 ал 0,00989763 16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,00005061 37 бг 0,0000088	13	ай	0,0009814
16 ам 0,00479038 17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,0009792 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,0000066 36 бв 0,00005061 37 бг 0,0000088	14	ак	0,00753655
17 ан 0,00627569 18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,0000088	15	ал	0,00989763
18 ао 0,00123445 19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,00005061 37 бг 0,0000088	16	ам	0,00479038
19 ап 0,0024161 20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,0000088	17	ан	0,00627569
20 ар 0,00284959 21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	18	ao	0,00123445
21 ас 0,00815708 22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	19	ап	0,0024161
22 ат 0,00580699 23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	20	ар	0,00284959
23 ау 0,00058752 24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	21	ac	0,00815708
24 аф 0,00035647 25 ах 0,0009792 26 ац 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	22	ат	0,00580699
25 ах 0,0009792 26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,00005661 36 бв 0,0000588	23	ay	0,00058752
26 ац 0,00007261 27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,000088	24	аф	0,00035647
27 ач 0,00168115 28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000666 36 бв 0,00005061 37 бг 0,0000088	25	ax	0,0009792
28 аш 0,00104962 29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,000066 36 бв 0,00005061 37 бг 0,0000088	26	ац	0,00007261
29 ащ 0,00022665 30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000666 36 бв 0,00005061 37 бг 0,0000088	27	ач	0,00168115
30 аъ 0 31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,00005061 37 бг 0,0000088	28	аш	0,00104962
31 аь 0 32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,00005061 37 бг 0,0000088	29	ащ	0,00022665
32 аю 0,00100341 33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,00005061 37 бг 0,0000088	30	аъ	0
33 ая 0,00269555 34 ба 0,00071295 35 бб 0,0000066 36 бв 0,00005061 37 бг 0,0000088	31	аь	0
34 6a 0,00071295 35 66 0,0000066 36 6в 0,00005061 37 6г 0,0000088	32	аю	0,00100341
35 66 0,0000066 36 бв 0,00005061 37 бг 0,0000088	33	ая	0,00269555
36 бв 0,00005061 37 бг 0,0000088	34	ба	0,00071295
37 бг 0,0000088	35	66	0,0000066
	36	бв	0,00005061
38 бд 0,000022	37	бг	0,0000088
	38	бд	0,000022

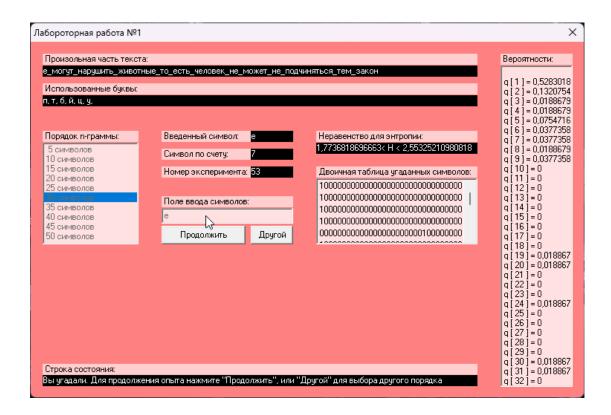
Значення Н і R

	Е	R
letters including space	4.363926750671057	0.1422194349436191
letters not including space	4.4622350062498795	0.11540714292613852
bigrams including spaces	3.9464790335576905	0.224273639591288
crossed bigrams including spaces	3.946960015449667	0.22417909700552663
bigrams not including spaces	4.133355455939298	0.18060417997930378
crossed bigrams not including space	4.134752310896859	0.1803272676436477

CoolPinkProgram







Отримані результати

2,0465<H(10)<2,7066 1,9194<H(20)<2,6855 1,7736<H(30)<2,5532

Висновки

Під час виконання лабораторної роботи, ми отримали змогу ознайомитись з такими поняттями як ентропія, надлишковість та обрахувати їх на практиці. Успішно були проведені експерименти на різних видах тексту(з пробілом та без), прорахована частота монограм та біграм, проведено знайомство з невеличкою програмою. Набуті навички знадобляться у майбутніх лаб.роботах та у професійній діяльності.