# КРИПТОГРАФІЯ

# КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

# Експериментальна оцінка ентропії на символ джерела відкритого тексту

ФБ-12 Юрченко Вікторія

# 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 Н1 та Н2 за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення Н1 та Н2 на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення Н1 та Н2 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

Частота букв у тексті:

Текст завантажую з <u>ia801806.us.archive.org/29/items/Elektrosudorozhnaya-terapiya-Nelson-2005/Nelson\_djvu.txt</u>. Фільтрую його функцією clear(). Частота рахується шляхом ділення кількості елемента/елементів у тексті на загальну кількість зустрічей цих символів. Таблиці отриманих частот:

Tactora Oyks y Tekeri:				
Символ	Кількість	Частота		
	148122	0.1453898		
0	89579	0.0879267		
e	79559	0.0780915		
N	74957	0.0735744		
l a	65198	0.0639954		
н	61250	0.0601202		
T	56760	0.0557130		
C	51411	0.0504627		
l p	42882	0.0420910		
В	36626	0.0359504		
Л	32220	0.0316257		
П	30460	0.0298982		
ĸ	27344	0.0268396		
M	27028	0.0265295		
І д	23808	0.0233689		
У	17824	0.0174952		
Я	17239	0.0169210		
ы	14895	0.0146203		
3	14273	0.0140097		
P	13913	0.0136564		
Ч	13314	0.0130684		
ь	11865	0.0116461		
Гб	10982	0.0107794		
й	10314	0.0101238		
x	9627	0.0094494		
9	7997	0.0078495		
ж	6513	0.0063929		
Ц	5694	0.0055890		
Ю	4958	0.0048665		
Ш	4416	0.0043345		
Ι ф	3898	0.0038261		
Щ	3094	0.0030369		
ë	580	0.0005693		
ъ	192	0.0001885		

частота бігр	ам у тексті	з перетинами:
Символ	Кількість	Частота
п		0.0184748
N		0.0167983
e		0.0165353
CT O	15865 15364	0.0155724   0.0150806
a	14483	0.0130800
c	13774	0.0135199
В	13122	0.0128800
ен	12723	0.0124883
HN	12188	0.0119632
HO	12179	0.0119544
Н	11939 11257	0.0117188
я   т	9875	0.0110494   0.0096929
TO	9815	0.0096340
пр	9797	0.0096163
po	9133	0.0089645
N	8956	0.0087908
pa	8776	0.0086141
на	8413	0.0082578
OB	8315	0.0081616
В	8306	0.0081528
0 по	8076 8039	0.0079270
й	8038	0.0078897
ко	7819	0.0076748
pe	7549	0.0074098
K	7348	0.0072125
ри	7318	0.0071830
9		0.0071124
M	7237	0.0071035
oc He	7233 7214	0.0070996     0.0070809
ле	7155	0.0070230
M	7065	0.0069347
ан	6981	0.0068522
І д І	6558	0.0064370
ти	6305	0.0061887
T	6226	0.0061112
X	6218	0.0061033
ны ол	6207 6191	0.0060925
		,
ър	1	0.0000010
фм	1	0.0000010
IN	1	0.0000010
ОР	1	0.0000010
йз	1	0.0000010
кё   ёл	1	0.0000010
I Фб	1	0.0000010
ГР	1	0.0000010
МЦ	1	0.0000010
І дф	1	0.0000010
дю	1	0.0000010
ਖੁਰ	1	0.0000010
ГЩ	1	0.0000010
ръ	1 1	0.0000010
ън   дщ	1	0.0000010
ьь	1	0.0000010
Щю	_	0.0000010
ьу	_	0.0000010
ЦМ	1	0.0000010
yë	_	0.0000010
ëĸ	1	0.0000010
лр	1	0.0000010
ьч	1 1	0.0000010     0.0000010
ыы	1	0.0000010     0.0000010     0.0000010
ьч	1 1 1	0.0000010     0.0000010
на   ыы   вді	1 1 1	0.0000010     0.0000010     0.0000010     0.0000010
ьч   ыы   щв   ёч	1 1 1 1 1 1	0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010
фз й ў ў ў ў ў ў ў ў ў ў ў ў ў ў ў ў ў ў	1 1 1 1 1 1 1	0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010
тт фз фа Б Б Б Б Б Б Б Б Б Б Б Б Б Б Б Б Б Б	1 1 1 1 1 1 1 1	0.0000010     0.0000010
Ра аш Б Б Т, Т, Т, Т, Т, Т, Т, Т, Т, Т, Т, Т, Т,	1 1 1 1 1 1 1 1 1	0.0000010     0.0000010
жт шт фз гр кр кр ед шв пп рам		0.0000010     0.0000010
Бе ре		0.0000010     0.0000010
Оў   Бе   Тар   Тар   Фе   Тар   Бе   Бе   Бе   Бе	1 1 1 1 1 1 1 1 1 1	0.0000010     0.0000010
Бе ре		0.0000010     0.0000010
ми на		0.0000010     0.0000010
Р4   На   На   На   На   На   На   На   Н		0.0000010     0.0000010
ьч   ыы   ыы   ыы   ыы   ыы   ыы   ыы		0.0000010     0.0000010
Р4    MM   MM   MM   MM   MM   MM   MM		0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010       0.0000010       0.0000000000
Ра   На   На   На   На   На   На   На		0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010       0.0000010       0.0000010       0.0000010         0.0000010         0.0000010           0.0000010
Ревория (Станова) (Станов		0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010     0.0000010       0.0000010     0.0000010       0.0000010       0.0000000000

Символ	Кількість	Частота
п	9445	0.0185416
N	8508	0.0167021
l e	8465	0.0166177
CT	7981	0.0156676
0	7835	0.0153810
l a	7158 6844	0.0140519
В	6546	0.0128505
ен	6421	0.0126051
НО	6072	0.0119200
HN	6011 5996	0.0118002
н   я	5489	0.0117708     0.0107755
T	4901	0.0096212
TO	4900	0.0096192
пр	4847	0.0095152
po	4632	0.0090931
и   pa	4466   4446	0.0087672   0.0087280
ОВ	4174	0.0081940
В	4144	0.0081351
на	4140	0.0081273
й	4063	0.0079761
о по	4058   3991	0.0079663   0.0078348
KO	3917	0.0076895
pe	3792	0.0074441
9	3692	0.0072478
bn	3655	0.0071752
K	3646 3606	0.0071575
M	3598	0.0070730
М	3582	0.0070319
не	3570	0.0070083
ле	3547	0.0069631   0.0068630
ан   д	3496 3342	0.0068630
x	3135	0.0061543
ТИ	3131	0.0061465
Т	3106	0.0060974
ол   ec	3087 3085	0.0060601   0.0060562
,		,
	1	0.0000020
дю	1	0.0000020
дю   чб	1 1	0.0000020     0.0000020
дю	1	0.0000020
од, ВР ХС ИОО НФ	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020     0.0000020     0.0000020     0.0000020
дю   чб   эх   юи   ън   аъ	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020
ОДД В В В В В В В В В В В В В В В В В В В	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020
дю   чб   эх   юи   ън   аъ	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020
од,	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020
дю   чб   эх   им   ын   аъ   дщ   чп   тх   ын   ын	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0.0000020     0.0000020
др чб эх ои ън аъ дщ чп тх ьь ьь уё иы		0.000020     0.000020
дю   чб   эх   им   ын   аъ   дщ   чп   тх   ын   ын		0.0000020     0.0000020
дю   чб   эх   юи   ън   аъ   дщ   чп   тх   ьь   уё   иы		0.0000020     0.0000020
одд   Од   Од   Од   Од   Од   Од   Од		0.0000020     0.0000020
Дю,		0.0000020     0.0000020
ОДД   ОВР		0.0000020     0.0000020
одд   Одд   Оду		0.0000020   1 0.
др.  др.  др.  др.  др.  др.  др.  др.		0.0000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.000000
од,		0.0000020     0.000020     0.0000020     0.000020
др.  др.  др.  др.  др.  др.  др.  др.		0.0000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.00000020   1 0.000000
ОДД   ОР   ОР   ОР   ОР   ОР   ОР   ОР		0.0000020     0.0000020
од,  од,  од,  хе  иод  иод  од,  од,  од,  од,  од,  од,		0.0000020     0.0000020
од,  чення од		0.000020   1 0.0000020   1 0.000020   1 0.000020   1 0.000020   1 0.000020   1 0.0000020
од,  од,  од,  хе  иод  иод  од,  од,  од,  од,  од,  од,		0.0000020     0.0000020
ДЮ, ДЮ, ДО, ДО, ДО, ДО, ДО, ДО, ДО, ДО, ДО, ДО		0.000020   1 0.0000020   1 0.000020   1 0.000020   1 0.000020   1 0.000020   1 0.0000020
ДД   95   95   95   95   95   95   95		0.000020   1 0.0000020   1 0.
од,  хе на		0.0000020   1 0.0000020   1 0.000020   1 0.000020   1 0.000020   1 0.000020   1 0.0000
дю дю до		0.0000020     0.000020     0.000020
од,  хе на		0.0000020     0.000020     0.000020
одд (ус. 1) дод (		0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.000020
дю  чб  чб  чб  на  на  на  на  на  на  на  на  на  н		0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.000020
од,  од,  од,  од,  од,  од,  од,  од,		0.000020   1 0.0000020   1 0.
жд на пр на		0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.000020
дю  чб  чб  на  на  на  на  на  на  на  на  на  н		0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020
ДЮ   Чб   9X   96   95   100   10		0.000020   1 0.0
дю  чб  чб  на  на  на  на  на  на  на  на  на  н		0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020     0.0000020

Для обрахунку H1 використовую формулу:  $H1 = -\sum_{i=1}^n q_i \log_2 q_i$ , де q — частота, n — кількість символів. Для H2:  $H2 = -\frac{1}{2}\sum_{i,j}q_{ij}\log_2 q_{ij}$ .

## Отримані значення ентропії:

```
H1 з пробілами: 4.402461438763894

H1 без пробілів: 4.451479170192825

H2 з пробілами та перетинами: 3.9970844782541084

H2 з пробілами та без перетинів: 3.997104842702902

H2 без пробілів та з перетинами: 4.1145792748163235

H2 без пробілів та перетинів: 4.1148592017574375
```

Надлишковість рахується за формулою  $R=1-\frac{H_\infty}{H_0}$ ;  $H_\infty$  дорівнюватиме значенню, порахованому вище, а  $H_0=\log_2 m$ , де m- кількість символів у алфавіті (для алфавіту з пробілами: m=34; без: m=33).

## Отримані значення надлишковості:

```
R для букв з пробілами: 0.13464499375450822

R для букв без пробілів: 0.117539378394374

R для біграм з пробілами та перетинами: 0.21432655078189222

R для біграм з пробілами та без перетинів: 0.21432254791259786

R для біграм без пробілів та з перетинами: 0.18432636755598786

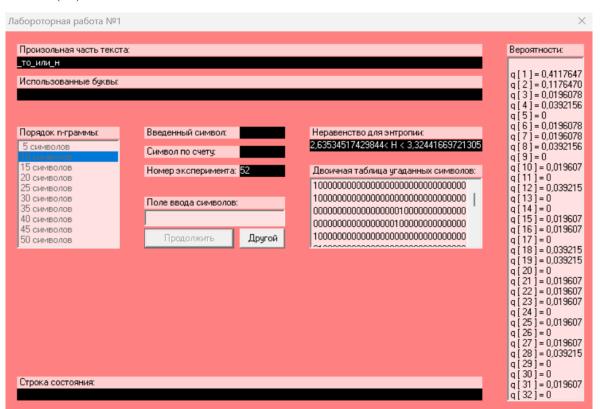
R для біграм без пробілів та перетинів: 0.18427087487748361
```

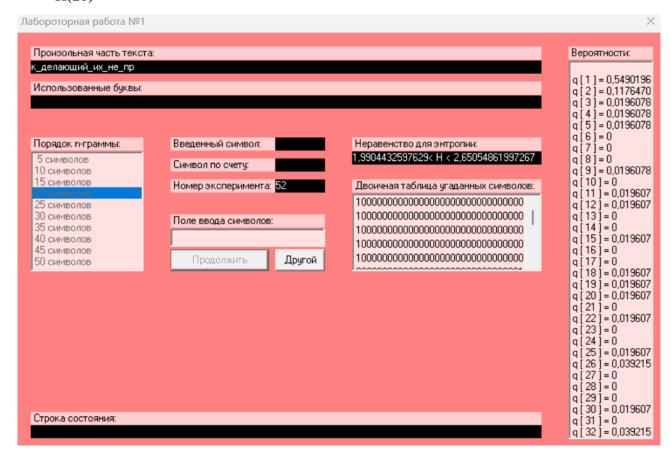
#### Отже:

	Ентропія	Надлишковість
Букви з пробілами	4,402461	0,134645
Букви без пробілів	4,451479	0,117539
Біграми з пробілами та перетинами	3,997085	0,214327
Біграми з пробілами та без перетинів	3,997105	0,214323
Біграми без пробілів та з перетинами	4,114579	0,184326
Біграми без пробілів та перетинів	4,114859	0,184271

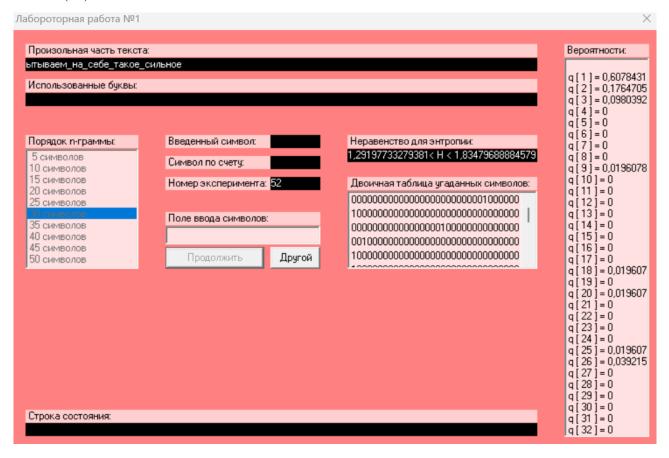
## **CoolPinkProgram**

## H(10)





#### H(30)



 $R=1-\frac{H_{\infty}}{H_0}$ ;  $H_{\infty}$  дорівнюватиме значенню, що видає програма, а  $H_0=\log_2 m$ , де m-кількість символів у алфавіті (у даному випадку m=32), тому  $H_0=5$ .

	Ентропія	Надлишковість
H(10)	2,635345 <h<3,324417< td=""><td>0,335117<r<0,472931< td=""></r<0,472931<></td></h<3,324417<>	0,335117 <r<0,472931< td=""></r<0,472931<>
H(20)	1,990443 <h<2,650549< td=""><td>0,469890<r<0,601911< td=""></r<0,601911<></td></h<2,650549<>	0,469890 <r<0,601911< td=""></r<0,601911<>
H(30)	1,291977 <h<1,834797< td=""><td>0,633041<r<0,741605< td=""></r<0,741605<></td></h<1,834797<>	0,633041 <r<0,741605< td=""></r<0,741605<>

## Висновок:

Під час виконання комп'ютерного практикуму, я практичним шляхом навчилася підраховувати частоти букв та біграм на довільно обраному тексті, які у подальшому використовувала для визначення значення ентропії на символ джерела та його надлишковості. За допомогою програми CoolPinkProgram отримала значення умовних ентропій джерела для 10-грам, 20-грам, 30-грам, які були використані для оцінки надлишковості.