# МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ «КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ імені ІГОРЯ СІКОРСЬКОГО» Фізико-технічний інститут

Криптографія Комп'ютерний практикум №1

Виконали:

студенти групи ФБ-01

Чуйко О. М.

Ченський К. Ю.

## 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.

#### Постановка задачі:

Написати програми для підрахунку частот букв і частот біграм в тексті, також для 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, де імовірності замінити відповідними частотами. Одержати значення Н1 та Н2 на тому ж тексті, в якому вилучено всі пробіли.

Використавши CoolPinkProgram оцінити значення  $\mathbf{H}^{(10)}, \mathbf{H}^{(20)}, \mathbf{H}^{(30)}$ .

Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделяї джерела.

# Хід роботи:

Спершу для знайденого тексту російською мовою провели обробку(перевели усі символи в нижній регістр, видалили повторні пробіли та усі символи крім зазначених), отримали два очищених файли(один з пробілами інший без) та дві відповідні змінні.

Далі написали функцію, яка б рахувала кількість елементів у тексті та сортувала отриманий масив. Після неї написали функцію, яка б виконувала відповідні дії для біграм, при тому з можливістю вказати крок, який для нашої роботи міг бути 1 або 2.

Після цього реалізували функцію, яка обраховує частоту та заносить дані у .xlsx табличку. На цьому етапі виникла проблема: використаний модуль для запису у .xlsx файл міг виконувати дії лише після створення, без можливості подальшого редагування, тож для кожного випадку було необхідно створити окремий файл, але для зручності перевірки та написання протоколу усі отримані файли були об'єднані у один.

Далі написали функції для підрахунку ентропії та надлишковості, використовуючи функції наведені у методичці до комп'ютерного практикуму.

Пілся цього у функції main() зробили 6 викликів для знаходження значення ентропії для усіх випадків, потім 6 викликів для підрахунку надлишковості, після чого занесли відповідні дані у .xlsx табличку.

### Таблиці частот букв:

ФБ-01 Чуйко Олександр, ФБ-01 Ченський Костянтин

	A	В
1	Symbol	Frequency
2	1 1	0.15709799806323
3	' <u>Q</u> '	0.0885243843969633
4	'e'	0.0707176337257712
5	'a'	0.0657611380674483
6	<u>'и</u> '	0.0564049497819127
7	<u>'Ţ'</u>	0.0546900256362916
8	<u>'ዟ'</u>	0.0527758783877962
9	'ç'	0.0450703958117904
10	'p'	0.0401598747156692
11	' <u>д</u> '	0.0389455624695783
12	' <u>B</u> '	0.0382989995789314
13	' <u>м</u> '	0.0291704948395378
14	<u>'Ķ'</u>	0.0290982490989351
15	'д'	0.026894389133073
16	<u>''</u>	0.0226048895041534
17	' <u>'</u> ,	0.022277959284052
18	<u>'g'</u>	0.0180096225488443
19	<u>'b'</u>	0.0167150372576393
20	<u>'ы</u> '	0.0164202162555835
21	<u>'c'</u>	0.0152372834725827
22	<u>'3</u> '	0.0151103061103112
23	<u> ˈ6</u> ̞ˈ	0.0147191574541185
24	<u>'й</u> '	0.0134209233881355
25	<u>'9</u> '	0.0125817051892145
26	l' <u>ж</u> '	0.00827323193144974
27	'x'	0.00729681980087906
28	<b>'</b> 迎'	0.00671228608145671
29	' <u>დ</u> '	0.00533888725504863
30	ľ <u>u</u> ′	0.00383851106636902
31	' <u>끯</u> '	0.00296280511966887
32	<u>'ş</u> '	0.00289639741871077
33	<u>'ф</u> '	0.00173097875464398
34	<mark>'</mark> ኔ'	0.00024300840020929
35	<u>'ë</u> '	0

	A	В
1	Symbol	Frequency
2	'Q'	0.105023341021325
3	'e'	0.0838978120389802
4	'a'	0.07801753693353
5	' <u>й</u> '	0.0669175653306183
6	<u>'Ţ</u> '	0.0648830178486089
7	<u>'ਮ</u> '	0.062612116552732
8	' <u>ç</u> '	0.0534705051218651
9	' <u>p</u> '	0.0476447732042305
10	<u>'Д</u> '	0.0462041404339906
11	' <u>в</u> '	0.0454370727450288
12	' <u>M</u> '	0.034607219786537
13	<u>'K</u> '	0.0345215090628268
14	' <mark>д</mark> '	0.0319068991072232
15	<u>'</u> '	0.0268179331075418
16	<u>'д</u> '	0.0264300704386311
17	<u>'я</u> '	0.0213662116206426
18	<u>'</u> Է'	0.019830344712947
19	' <u>ы</u> '	0.0194805756990186
20	<u>'Ç'</u>	0.0180771708188751
21	<u>'3</u> '	0.0179265277287178
22	<u>'Ğ</u> '	0.0174624777498424
23	'й'	0.015922282017717
24	<u>'y'</u>	0.0149266523988614
25	<u>'迷</u> '	0.00981517651214495
26	' <u>X</u> '	0.00865678309472853
27	<u>'</u> ਘ'	0.00796330542107341
28	<u>'</u> ဗ္ဗ'	0.00633393590569396
29	<u>'u'</u>	0.00455392330015722
30	<u>'</u> Щ'	0.00351500543700348
31	' <b>ᢖ</b> '	0.00343622083238099
32	<u>'</u> ф'	0.00205359430950056
33	'ኔ'	0.00028829970702516
34	<u>'ë'</u>	0

Перша табличка демонструє частоти для букв у тексті з пробілами, друга частоти для букв у тексті без пробілів.

Топ-10 по кількості для кожного з 6-ти випадків:

```
[+]Найчастіші 10 елементів: [(' ', 215275), ('o', 121307), ('e', 96906), ('a', 90114), ('w', 77293), ('T', 74943), ('H', 72320), ('c', 61761), ('p', 55032), ('n', 53368)]
[+]Найчастіші 10 елементів: [('o', 121307), ('e', 96906), ('a', 90114), ('w', 77293), ('T', 74943), ('H', 72320), ('c', 61761), ('p', 55032), ('n', 53368), ('s', 52482)]
[+]Найчастіші 10 елементів: [('o', 23038), ('c', 21736), ('w', 21663), ('s', 21334), ('e', 21331), ('n', 20444), ('a', 19822), ('H', 17788), ('To', 17313), ('cT', 14783)]
[+]Найчастіші 10 елементів: [('o', 23038), ('c', 21736), ('w', 21663), ('s', 21334), ('e', 21331), ('n', 20444), ('a', 19822), ('h', 17788), ('To', 17313), ('cT', 14783)]
[+]Найчастіші 10 елементів: [('To', 17839), ('CT', 15131), ('Ha', 12676), ('Ho', 12066), ('He', 12009), ('eH', 11828), ('no', 10968), ('ko', 10965), ('os', 10882), ('pa', 10615)]
```

## Таблиці частот біграм:

ФБ-01 Чуйко Олександр, ФБ-01 Ченський Костянтин

	Α	В		A	В		Α	В		Α	В
1	Bigram	Frequency	1	Bigram	Frequency	1	Bigram	Frequency	1	Bigram	Frequency
2	'g '	0.01681	2	'TQ'	0.01544	2	'g '	0.01681	2	'IQ'	0.01544
3	' <u>ç</u> '	0.01586	3	<u>'ςτ'</u>	0.0131	3	' <u>چ</u> '	0.01586	3	'ST'	0.0131
4	' <u>и</u> '	0.01581	4	'на'	0.01097	4	' <u>w</u> '	0.01581	4	'на'	0.01097
5	' B'	0.01557	5	'но'	0.01045	5	' <u>в</u> '	0.01557	5	'ዟQ'	0.01045
6	'e '	0.01557	б	'не'	0.0104	б	' <u>ę</u> '	0.01557	6	' <del>ዟ</del> ይ'	0.0104
7	' <u>"</u>	0.01492	7	'eн'	0.01024	7	' <u>"</u> "	0.01492	7	'ен'	0.01024
8	' <u>a</u> '	0.01447	8	'до'	0.0095	8	' <u>a</u> '	0.01447	8	'до'	0.0095
9	' ዟ'	0.01298	9	'KQ'	0.00949	9	' <u>ዟ</u> '	0.01298	9	'KQ'	0.00949
10	'IQ'	0.01263	10	'QB'	0.00942	10	'IQ'	0.01263	10	'QB'	0.00942
11	'ÇŢ'	0.01079	11	'pa'	0.00919	11	<u>'Ç</u> Ţ'	0.01079	11	'pa'	0.00919
12	' <u>Я</u> '	0.01078	12	'QT'	0.00893	12	<u>'g</u> '	0.01078	12	'QT'	0.00893
13	<u>'</u> ե '	0.01001	13	'BQ'	0.00861	13	<u>'</u> ե '	0.01001	13	'BQ'	0.00861
14	'й '	0.00976	14	'po'	0.00845	14	' <u>й</u> '	0.00976	14	'po'	0.00845
15	' <u>и</u> '	0.0093	15	'ec'	0.00844	15	' <mark>¼</mark> '	0.0093	15	'ec'	0.00844
16	'на'	0.00921	16	'ep'	0.00842	16	'ਸ਼ੁਕੂ'	0.00921	16	'ep'	0.00842
17	' ፲'	0.00918	17	'gc'	0.0084	17	' ፲'	0.00918	17	'QC'	0.0084
18	' <u>K</u> '	0.00904	18	'ни'	0.00839	18	' <u>K</u> '	0.00904	18	<u>,4M,</u>	0.00839
19	' <u>Q</u> '	0.00884	19	'et'	0.00795	19	' <u>Q</u> '	0.00884	19	<u>'ex'</u>	0.00795
20	<u>'ਮੁ</u> ę'	0.00873	20	' <u>T</u> a'	0.00785	20	'ਸ਼ਵ'	0.00873	20	'Ja'	0.00785
21	'ዟQ'	0.00862	21	' <u>r</u> o'	0.0078	21	'ዟር'	0.00862	21	<u>'۲۵</u> '	0.0078
22	' <u>₩</u> '	0.00832	22	'ал'	0.00774	22	' <u>₩</u> '	0.00832	22	'ал'	0.00774
23	'πο'	0.008	23	'QH'	0.0077	23	'፲ጲ'	0.008	23	<u>'он</u> '	0.0077
24	'KQ'	0.00773	24	'ка'	0.00768	24	'KQ'	0.00773	24	'ĸa'	0.00768
25	'pa'	0.0077	25	'ди'	0.00763	25	'pa'	0.0077	25	<u>'ди</u> '	0.00763
26	' <u>e</u> н'	0.00742	26	'QM'	0.00743	26	<u>'</u> ਉਮ੍	0.00742	26	'QM'	0.00743
27	'Ι '	0.00741	27	'pe'	0.00737	27	<u>'</u> Ţ '	0.00741	27	'pe'	0.00737
28	' ∭'	0.00736	28	'π <b>ρ</b> '	0.00724		' ⋈'	0.00736	28	<u>'ДR'</u>	0.00724
29	'po'	0.00702	29	'QR'	0.00711	29	' <u>po</u> '	0.00702	29	'QR'	0.00711
30	'BQ'	0.00689	30	'ол'	0.00698	30	'BQ'	0.00689	30	<u>'ОЛ'</u>	0.00698
31	<u>'μμ'</u>	0.00675	31	'да'	0.00695	31	'HW,	0.00675	31	'да'	0.00695
32	' д'	0.00672	32	'ан'	0.00675	32	'д'	0.00672	32	<u>'ан</u> '	0.00675
33	'ep'	0.00663	33	' <u>ри</u> '	0.00637	33	'ep'	0.00663	33	'DN,	0.00637
34	<u>'9T</u> '	0.00653	34	'፲ኴ'	0.00633	34	<u>'9</u> ूर'	0.00653	34	<u>'Jb'</u>	0.00633
35	<u>'፫</u> ፬'	0.00651	35	' <u>a</u> ĸ'	0.0063	35	<u>'۲0</u> '	0.00651	35	<u>'ак</u> '	0.0063
36	<u>'</u> В '	0.00651	36	'ŢĠ'	0.00625	36	<u>'B</u> '	0.00651	36	' <u>te</u> '	0.00625
37	'Ig'	0.00649	37	' <u>ве</u> '	0.00622	37	'Ja'	0.00649	37	'ge'	0.00622
38	' <u>ĸ</u> a'	0.00639	38	'до'	0.00615	38	<u>'ĸa</u> '	0.00639	38	'до'	0.00615
39	'ал'	0.00624	39	'ел'	0.00611	39	'ал'	0.00624	39	' <u>ел</u> '	0.00611
40	'QB'	0.0062	40	'од'	0.00611	40	' <u>ов'</u> 	0.0062	40	<u>'од'</u>	0.00611

На першій таблиці частоти біграм з кроком 1 у тексті з пробілами, на другій частоти біграм з кроком 1 у тексті без пробілів, на третій та четвертій частоти біграм з кроком два у відповідному порядку стосовно пробілів.

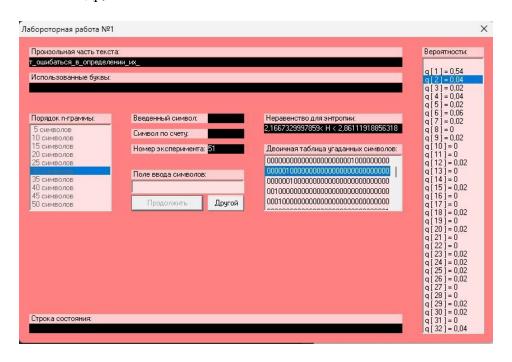
ФБ-01 Чуйко Олександр, ФБ-01 Ченський Костянтин

	A		В		С
	Name	Н		R	
	Letters with spaces		4.4115798996701		0.132852654195333
	Letters with out spaces				0.109990435913449
4	Bigrams with spaces and step 1		4.03915965362009		0.206056185635473
5	Bigrams with out spaces and step 1		4.19671375043315		0.168044040348122
6	Bigrams with spaces and step 2		4.03887392578276		0.206112348765553
7	Bigrams with out spaces and step 2		4.19652114184497		0.16808222304829

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



## ФБ-01 Чуйко Олександр, ФБ-01 Ченський Костянтин



- $3.28994310106647 < H^{(10)} < 3.69366068968818$
- $0.342011379786706 < R^{(10)} < 0.261267862062364$
- $2.46639168900549 < H^{(20)} \! < 3.01287868934203$
- $0.506721662198902 < R^{(20)} < 0.397424262131594$
- $2.16673299977859 < H^{(30)} < 2.86111918856318$
- $0.566653400044282 < R^{(30)} < 0.427776162287364$

#### Висновки:

Виконавши цей комп'ютерний практикум ми набули практичних навичок щодо оцінки ентропії на символ джерела, порівняли різні моделі джерела відкритого тексту для наближеного визначення ентропії.

Написали програму для підрахунку усіх необхідних для роботи значень, використавши для збереження отриманих результатів модуль, який надає змогу запису даних у .xlsx файл.

Використавши отримані значення ентропії, підрахували надлишковість.

За допомогою CoolPinkProgram оцінили значення для  $H^{(10)}$ ,  $H^{(20)}$ ,  $H^{(30)}$ .