# НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ "КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ ІМЕНІ ІГОРЯ СІКОРСЬКОГО" ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ

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

Виконали студенти 3 курсу групи ФБ-21 КАЮН Вероніка РУДЮК Олександр **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

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

### Хід роботи

## Довільний текст text1

Спочатку підрахуємо частоту букв та біграм у нашому тексті **text2** – очищений текст з пробілами.

Частота букв	Частота біграм з	Частота біграм <i>без</i>
	перетином	перетину

Буква	Частота	Біграма	Частота	Біграма	Частот
	0,15862	од	0,00452	од	0,0046
)	0,09315	дн	0,00194	на	0,0092
	0,0729	на	0,00914	жд	0,000
)	0,06827	аж	0,00107	Ы	0,003
1	0,05709	жд	0,00071	ве	0,005
ł	0,05396	ды	0,00049	СН	0,0008
	0,05069	ы	0,00383	ОЮ	0,0004
1	0,04441	В	0,01626	Ч	0,0052
	0,0426	ве	0,00558	ac	0,0038
)	0,0402	ec	0,00405	Н	0,015
3	0,03935	СН	0,00086	еб	0,000
К	0,03085	но	0,00971	ыв	0,0008
/	0,02552	ою	0,00041	ал	0,008
М	0,02528	ю	0,0026	0	0,020
1	0,02399	Ч	0,00515	жа	0,0013
4	0,02362	ча	0,00221	рк	0,0004
	0,01615	ac	0,00406	ОГ	0,004
1	0,01608	С	0,00291	за	0,005
•	0,01527	н	0,01519	ка	0,007
3	0,015	не	0,00883	та	0,0054
Ы	0,01455	еб	0,00087	В	0,016
4	0,01327	бы	0,00274	M	0,0059
б	0,01302	ыв	0,00085	ос	0,0054
й	0,00978	ва	0,00623	КВ	0,000
К	0,00756	ал	0,00875	e	0,015
Ш	0,00738	ло	0,00622	П	0,0169
(	0,00685	О	0,02097	ат	0,0049
o	0,00448	ж	0,00181	ри	0,005
ц	0,00296	жа	0,00133	ар	0,004
4	0,00276	ар	0,00441	ши	0,0017
€	0,00258	рк	0,00041	X	0,0028
þ	0,00181	ко	0,00837	пр	0,0065
С	3E-06	ог	0,00433	уд	0,0018

Всі значення можна переглянути у файлі frequency\_data.xlsx

Далі обчислюємо ентропію та надлишковість тексту.

```
Ентропія Н1 (монограми): 4.37573
Надлишковість R1 (монограми): 0.13256
Ентропія Н2 (біграми з перетином): 3.98572
Надлишковість R2 (біграми з перетином): 0.20987
Ентропія Н2 (біграми без перетину): 3.98463
Надлишковість R2 (біграми без перетину): 0.21009
```

Видаляємо пробіли із тексту та зберігаємо у text3.

Частота букв		Частота біграм з перетином		Частота біграм <i>без</i> перетину		
Буква	Частота	Біграма	Частота	Біграма	Частота	
0	0,11071	од	0,00622	од	0,00628	
a	0,08664	дн	0,00253	на	0,01073	
е	0,08114	на	0,01096	жд	0,00079	
И	0,06785	аж	0,00147	ыв	0,0014	
Н	0,06414	жд	0,00086	ec	0,0066	
Т	0,06025	ды	0,00059	но	0,01176	
Л	0,05278	ыв	0,00147	юч	0,00022	
С	0,05063	ве	0,00677	ac	0,00676	
р	0,04778	ec	0,00661	не	0,01049	
В	0,04677	СН	0,00142	бы	0,00306	
К	0,03666	но	0,01182	ва	0,00755	
У	0,03034	ою	0,0005	ло	0,00804	
M	0,03004	юч	0,00027	жа	0,0016	
П	0,02851	ча	0,00265	рк	0,00059	
Д	0,02807	ac	0,00681	ог	0,00569	
Г	0,0192	не	0,01055	03	0,00249	
Я	0,01911	еб	0,002	ак	0,00641	
ь	0,01815	бы	0,00325	ат	0,00696	
3	0,01783	ва	0,00748	ав	0,00507	
Ы	0,01729	ал	0,01082	мо	0,00412	
4	0,01577	ло	0,0079	ск	0,00475	
б	0,01547	ож	0,00232	ве	0,00665	
й	0,01162	жа	0,00158	па	0,002	
ж	0,00898	ар	0,00573	тр	0,00353	
Ш	0,00877	рк	0,00059	иа	0,00053	
x	0,00814	ко	0,0102	рш	0,00052	
Ю	0,00533	ог	0,00569	их	0,00185	
щ	0,00351	го	0,00913	пр	0,00795	
ц	0,00328	03	0,0025	уд	0,00256	
Э	0,00306	за	0,00658	ax	0,00136	
ф	0,00215	ак	0,0064	по	0,01089	
С	3,6E-06	ка	0,00856	яв	0,00206	

Всі значення можна переглянути у файлі frequency\_data\_spaces\_del.xlsx

```
Ентропія Н1 (монограми): 4.45070
Надлишковість R1 (монограми): 0.10986
Ентропія Н2 (біграми з перетином): 4.14648
Надлишковість R2 (біграми з перетином): 0.17070
Ентропія Н2 (біграми без перетину): 4.14791
Надлишковість R2 (біграми без перетину): 0.17042
```

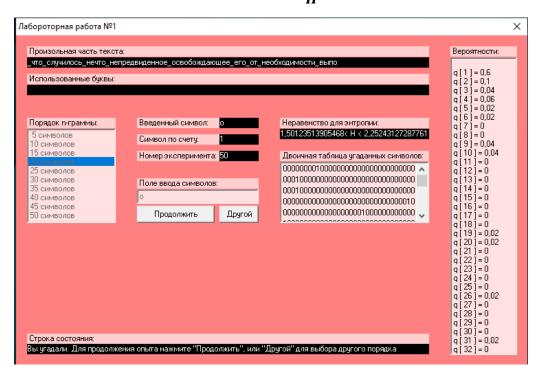
 За допомогою програми Cool Ріпк<br/>Program оцінимо значення  $H^{\;(10)}\,,\,H^{\;(20)}\,,\,H^{\;(30)}$ 

 $H^{(10)}$ 

Іабороторная работа №1			
Произольная часть тексти учить_ем Использованные буквы:  Порядок п-граммы:  5 символов 15 символов 20 символов 20 символов 30 символов 35 символов 45 символов 45 символов 50 символов 50 символов	Введенный символ:  Символ по счету:  Номер эксперимента: 51  Поле ввода символов:  Продолжить  Другой	Неравенство для энтропии: 2,49602844092468     Неравенство для энтропии: 2,496028468     Неравенство для энтропии: 2,496028468	Вероятности:  q[1] = 0.38 q[2] = 0.14 q[3] = 0.04 q[4] = 0.06 q[5] = 0.06 q[6] = 0.06 q[7] = 0 q[8] = 0.02 q[11] = 0.02 q[11] = 0.02 q[11] = 0.02 q[12] = 0.02 q[13] = 0 q[16] = 0.06 q[7] = 0 q[18] = 0.02 q[19] = 0 q[20] = 0 q
Строка состояния:			Q[23] = 0.04   Q[24] = 0.02   Q[25] = 0   Q[26] = 0   Q[27] = 0   Q[28] = 0   Q[29] = 0   Q[30] = 0   Q[31] = 0.04   Q[32] = 0.02

Надлишковість при мінімальному значенні ентропії складає: 0.5008 Надлишковість при максимальному значенні ентропії складає: 0.3728

 $H^{(20)}$ 



Надлишковість при мінімальному значенні ентропії складає: 0.6998 Надлишковість при максимальному значенні ентропії складає: 0.5495

 $H^{(30)}$ 

Произольная часть текста	χ.		Вероятности:
•	жолько_эти_учения_были_похожи_друговоровороворовороворовороворовороворово	на_друга_и_на_наш  Неравенство для энтропии:  1,23182379027345< H < 2,01366068968813  Двоичная таблица угаданных символов: 1000000000000000000000000000000000000	q [1] = 0,64 q [2] = 0,12 q [3] = 0,04 q [4] = 0,04 q [6] = 0 q [7] = 0,02 q [8] = 0,02 q [9] = 0,04 q [10] = 0,02 q [11] = 0,02 q [11] = 0,02 q [13] = 0,02 q [14] = 0,02 q [14] = 0,02 q [17] = 0 q [17] = 0 q [17] = 0 q [17] = 0 q [18] = 0 q [20] = 0 q [22] = 0 q [23] = 0 q [24] = 0 q [25] = 0 q [26] = 0 q [27] = 0 q [28] = 0 q [29] = 0

Надлишковість при мінімальному значенні ентропії складає: 0.7536

Надлишковість при максимальному значенні ентропії складає: 0.5973

**Висновок:** під час виконання даної лабораторної роботи ми навчилися екпериментально визначати частоти літер і біграм у тексті і на основі цих значень обчислювати ентропію і надлишковіть. За допомогою програми CoolPinkProgram ми здійснили серію експериментів, щоб оцінити значення ентропії  $H^{(10)}$ ,  $H^{(20)}$ ,  $H^{(30)}$ .