# МІНІСТЕРСТВО ОСВІТИ І НАУКИ, МОЛОДІ ТА СПОРТУ УКРАЇНИ 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 гр. ФБ-11

Цема В.В.

Ципун Р.Г.

Перевірила

Селюх П. В.

#### 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.

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.

### Обраний текст - Лев Товстий «Війна і Мир» Том 1 (1.23Мб) Букви:

#### 3 пробілами

H1 in text: 4.385949914023147 R1 in text: 0.11470001605958835

Буква, частота					
пробіл	0,16095				
a	0,07038				
б	0,01449				
В	0,03862				
г	0,0174				
Д	0,02552				
e	0,06683				
ж	0,0085				
3	0,01494				
И	0,05575				
й	0,00966				
К	0,03009				
Л	0,04247				
M	0,0248				
Н	0,05467				
0	0,09539				
П	0,02156				
p	0,03824				
С	0,04378				
Т	0,04765				
У	0,02406				
ф	0,00188				
X	0,00716				
ц	0,00339				
Ч	0,01143				
Ш	0,00792				
щ	0,00235				
Ы	0,01593				
Ь	0,01678				
Э	0,00254				
Ю	0,00544				
Я	0,01942				

Без пробілів

H1 in text: 4.468590618019441 R1 in text: 0.09801906544343475

Буква, частота

a	0,083875
б	0,01727
В	0,046033
Г	0,020735
Д	0,030412
е	0,079655
ж	0,01013
3	0,017809
И	0,066447
й	0,011515
К	0,035861
Л	0,050614
M	0,029559
н	0,065162
О	0,113683
п	0,025692
p	0,045571
С	0,052175
Т	0,056794
У	0,028672
ф	0,002239
x	0,008531
ц	0,004046
ч	0,013628
ш	0,009445
щ	0,002805
ы	0,01898
ь	0,020003
э	0,003023
ю	0,006487
я	0,02315

# Біграми, що перетинаються:

3 пробілами

H2 in text, with interceptions: 3.9784674005199263 R2 in text, with interceptions: 0.19694998920839157

*	a	6	В	г	Д	е	ж
а	1.557807106092895e-05	0.0006620680200894805	0.00341159756234344	0.0009487045276105732	0.002112386435861966	0.0008505626799267207	0.001188606821948879
6	0.0011403148016599991	2.4924913697486324e-05	3.738737054622948e-05	1.869368527311474e-05	1.4020263954836056e-05	0.0019690681821014195	4.673421318278685e-06
В	0.005767001906755898	9.34684263655737e-06	2.1809299485300533e-05	2.1809299485300533e-05	0.0002726162435662567	0.004167134008798494	-
г	0.0010281526900213108	-	7.633254819855186e-05	-	0.000861467329669371	0.000563926172405628	-
Д	0.004012911105295298	3.427175633404369e-05	0.0009128749641704366	9.34684263655737e-06	2.0251492379207637e-05	0.004263718049376254	1.2462456848743162e-05
е	9.34684263655737e-06	0.0010421729539761468	0.0013147891975424035	0.00333059159282661	0.002406811978913523	0.0014394137660298351	0.0009284530352313656
ж	0.0013350406899216112	4.984982739497265e-05	-	2.336710659139343e-05	0.000623122842437158	0.0036966762627584403	1.2462456848743162e-05

Повний файл lab1\_bigram\_with\_interceptions\_in\_text\_with\_spaces.csv

Без пробілів

H2 in text, with interceptions: 4.150657749319352 R2 in text, with interceptions: 0.16219352458497438

*	a	6	В	г	Д	е	ж
а	0.0003861792391526336	0.001583706206717291	0.005820538051651473	0.0017192402666122055	0.003267670485137669	0.0015317205399082824	0.001533577170865747
б	0.0013683370156513989	3.527598819182711e-05	5.941219063886671e-05	2.2279571489575018e-05	2.042294053211043e-05	0.002359777946937487	9.283154787322924e-06
В	0.006986502292939232	0.0002487885483002544	0.0004641577393661462	0.0004121720725571378	0.0007073763947940067	0.005126158073559719	2.9706095319433356e-05
г	0.0012402294795863427	4.270251202168545e-05	0.00016524015521434803	1.2996416702252093e-05	0.0010991255268190342	0.0006980932400066839	-
Д	0.0047993910250459515	9.468817883069382e-05	0.0011752473960750822	9.283154787322924e-05	5.5698928723937545e-05	0.005111305025900002	2.2279571489575018e-05
е	0.0001949462505337814	0.0021963944226806036	0.0034421937951393403	0.004545032583873304	0.0037244017006739572	0.002047863946083437	0.0013367742893745011
ж	0.0015911327305471491	6.683871446872505e-05	2.9706095319433356e-05	3.341935723436253e-05	0.000755648799688086	0.004418781678765712	1.4853047659716678e-05

Повний файл lab1\_bigram\_with\_interceptions\_in\_text\_without\_spaces.csv

#### Біграми, що не перетинаються:

Без пробілів

H2 in text, without interceptions: 4.1497924670210695 R2 in text, without interceptions: 0.1623681810265184

*	a	6	В	г	Д	е	ж
а	0.0004530179536213587	0.0015892760995896845	0.005710996825161063	0.001667254599803197	0.003197118508754015	0.001552143480440393	0.0015447169566105345
б	0.0013961864800133678	2.9706095319433356e-05	7.055197638365422e-05	1.8566309574645847e-05	7.426523829858339e-06	0.00245817938768311	7.426523829858339e-06
В	0.007062624162195281	0.0002487885483002544	0.00047901078702586287	0.0004270251202168545	0.0006721004066021797	0.005020330108984237	2.5992833404504185e-05
г	0.0012847886225654927	4.084588106422086e-05	0.0001559570004270251	1.1139785744787509e-05	0.0011325448840533968	0.0006758136685171089	-
Д	0.004764115036854124	0.00010025807170308757	0.001180817288947476	0.00010397133361801674	7.426523829858339e-05	0.005168860585581404	2.2279571489575018e-05
е	0.0001819498338315293	0.0022242438870425725	0.0034533335808841275	0.0046341508698316035	0.0035610181764170737	0.0020311542674662556	0.0012736488368207052
ж	0.0017155270046972763	6.683871446872505e-05	2.2279571489575018e-05	2.2279571489575018e-05	0.0007760717402201965	0.00438164905961642	1.8566309574645847e-05

Повний файл  $lab1\_bigram\_without\_interceptions\_in\_text\_without\_spaces.csv$ 

3 пробілами

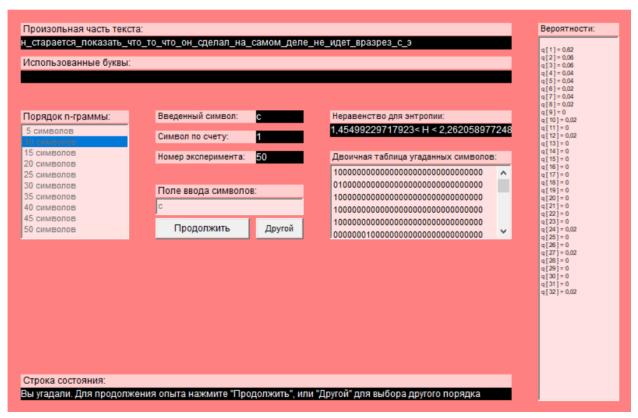
H2 in text, without interceptions: 3.9784288918899122 R2 in text, without interceptions: 0.19695776214018546

*	a	6	В	г	Д	е	ж
а	1.869368527311474e-05	0.0007321693398636607	0.0034988347602846427	0.0008848344362607644	0.002078114679527922	0.0008661407509876497	0.001218205156964644
б	0.0011808177864184146	2.1809299485300533e-05	2.8040527909672112e-05	6.231228424371581e-06	1.869368527311474e-05	0.0019784150247379767	6.231228424371581e-06
В	0.005611221196146608	3.1156142121857905e-06	2.4924913697486324e-05	2.1809299485300533e-05	0.0002804052790967211	0.004181154272753331	-
г	0.0010343839184456824	-	6.542789845590159e-05	-	0.000831868994653606	0.0006324696850737154	-
Д	0.0040908014605999425	3.11561421218579e-05	0.0009222218068069939	6.231228424371581e-06	1.557807106092895e-05	0.0043026632270285765	1.869368527311474e-05
е	3.1156142121857905e-06	0.001025037075809125	0.0012649393701474308	0.0033087822933413093	0.0023554043444124575	0.001473685522363879	0.0009471467205044802
ж	0.0012431300706621303	4.673421318278686e-05	-	3.427175633404369e-05	0.0006418165277102729	0.003623459328772074	1.2462456848743162e-05

Повний файл lab1\_bigram\_without\_interceptions\_in\_text\_with\_spaces.csv

2. За допомогою програми CoolPinkProgram оцінити значення H10, H20, H30.

H10:



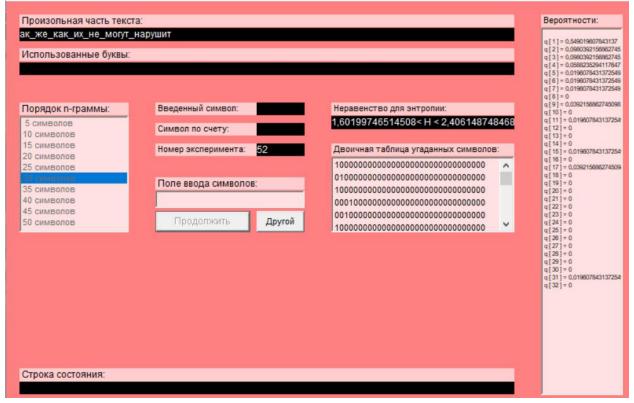
1,45499229717923 < H(10) < 2,262058977248

H20:



2,01197233763589 < H(20) < 2,747919939945

H30:



1,60199746514508 < H(30) < 2,406148748468

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

$$R = 1 - \frac{H_{\infty}}{H_0},$$

$$H(0) = \log 2(32) = 5$$

H(10):

0,709001541 < R < 0,547588205 71% < R < 54,7% H(20):

0,597605532 < R < 0,45041601260% < R < 45%

H(30):

 $\begin{array}{l} 0,679600507 < R < \text{0,51877025} \\ 68\% < R < 52\% \end{array}$ 

#### Висновки:

Ми засвоїли теоретичні поняття ентропії на символ джерела та його надлишковості. Навчилися на практиці визначати частоти монограм та біграм відкритого тексту, а також розраховувати відповідні значення ентропій (H1, H2). Попрацювавши з CoolPinkProgram, ми отримали значення ентропій для H10, H20, H30. У підсумку, ми набули практичних навичок оцінки ентропії.