## Міністерство освіти і науки 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

> Роботу виконали: Касаб О.Р. Косигін О.С. Групи ФБ-06

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

### Хід Роботи

Для виконання даної роботи потрібно було обрати текст російською мовою. Ми обрали текст ранобе "No game no life", коротко – "NGNL"

Найбільшими труднощами особисто для нас стали розбір та зрозуміння завдань лабораторної роботи, підрахунок необхідних п-грам та вивід результатів

## частоти нграм з пробілами і без

1 <b>a</b> (	пастоти играм з прооглами гоез						
N-gramma with spaces					N-gramma without spaces		
ъ	0.00023589313677991353	ы	0.01733996011591426	ъ	0.0002774369097796511	Ы	0.02039374699887958
ф	0.0025313148139075336	у	0.02164228801826176	ф	0.0029771114549431788	у	0.02545376940724537
щ	0.0029468496471583045	П	0.022589489690562643	щ	0.0034658272421704103	П	0.026567785306514432
ц	0.002955922460111378	Д	0.023979444634973517	ц	0.003476497892546551	Д	0.028202528944139146
Э	0.0049319811212908075	M	0.026993433097984566	Э	0.005800565544469936	M	0.03174731899909299
ю	0.005387436331535102	К	0.0274416300578664	ю	0.006336232193352185	К	0.03227444912767433
ш	0.006958847535007449	В	0.03684832252761311	ш	0.008184388838499706	В	0.04333777943765672
х	0.007205628047331051	р	0.039539318849494735	х	0.008474630528730726	p	0.046502694339219976
ж	0.008564735427701475	Л	0.04139743094228421	ж	0.010073093955076561	Л	0.048688043536253535
й	0.00902019063794577	С	0.04785545920228199	й	0.010608760603958812	С	0.05628341247399029
ч	0.012536812938557096	Н	0.053074141212889926	ч	0.01474470468975084	Н	0.06242117057034626
3	0.013975761072914568	Т	0.0553024240741648	3	0.016437069839406713	Т	0.06504188230272635
Γ	0.014603599729267261	И	0.059911413054326186	Γ	0.01717547884543563	И	0.07046257269380568
б	0.014603599729267261	e	0.0678392370127219	б	0.01717547884543563	e	0.07978658699247719
Я	0.01640546038174768	a	0.07059555758786566	Я	0.019294670010137117	а	0.08302833057674865
ь	0.01673389621064894	0	0.09831118659691489	Ь	0.0196809475537534	0	0.11562503334578242
			0.14974133410270787				

# частоти біграм

на скріншоті вказані дані з кінця таблиці, для перегляду повних таблиць значень зверніть увагу на **csv** файли прикріплені до лаби:

	Block with spaces		Block without spaces		Cross with spaces		Cross without spaces
Ы	0.0058682847696779885	за	0.005702383391381547	В	0.005797527477014028	ИС	0.005730139251987409
те	0.0059009468370416875	ис	0.005762138905961893	Ы	0.005801156602195258	за	0.005789894894093795
M	0.005959012734577153	ил	0.005838967424708051	М	0.005810229415148332	ил	0.005915808568532252
го	0.006027965987900519	со	0.0059371371986614765	го	0.006098744867056072	со	0.005969161820412954
во	0.00611143571560775	ит	0.006090794236153794	ка	0.006211247747674184	ИТ	0.006058795283572534
ка	0.00613321042718355	ет	0.00611213549136106	ол	0.0062747574383457	од	0.006220989169289868
ТЬ	0.006202163680506916	ом	0.006197500512190125	ть	0.006312863252748609	об	0.0063511711038787815
ол	0.006394506966093145	ло	0.006308475039267909	во	0.0063527836297421324	ло	0.00635330523395401
ОТ	0.006419910796264911	ва	0.006381035306972615	та	0.006361856442695206	ом	0.006355439364029237
та	0.006419910796264911	об	0.006385303558014068	от	0.006396333131916886	ет	0.006385317185082431
д	0.006459831100820544	од	0.006415181315304241	й	0.006470730198132089	ва	0.006413060876060396
M	0.006477976693800376	ер	0.006419449566345694	д	0.006501577762172539	pe	0.006434402176812677
й	0.006514267879760042	pe	0.006517619340299119	м	0.006565087452844055	ер	0.006498426079069519
ос	0.0068481467905889696	те	0.007059687222563682	ос	0.006672146645690323	те	0.006878301232460119
ОВ	0.0069570203484679675	ат	0.007209076009014546	ов	0.0069080397824702365	ан	0.007115189670810436
ли	0.007015086246003433	го	0.0072346855152632655	пр	0.007062277602672488	го	0.007204823133970016
ен	0.0070259736017913325	ан	0.0072944410298436116	ен	0.00707860866598802	ат	0.007283785946753454
пр	0.007080410380730832	ТЬ	0.007328587038175237	ли	0.007134860106297077	ка	0.007324334418182788
ал	0.00715662187124613	ка	0.007473707573584648	ал	0.007154820294793839	ТЬ	0.0074246385317185085
ор	0.0075739705097822895	ac	0.007495048828791914	ор	0.007470554185560799	та	0.007552686336232193
т.	0.00773365172800482	та	0.007648705866284231	т.	0.007795360889280834	ac	0.007595368937736755
ла	0.00790422030201525	во	0.007721266133988936	ни	0.008131054968544557	ОЛ	0.007717014352024756
ни	0.008205437145480478	ак	0.007853581916273987	ла	0.008145571469269475	ак	0.007898415408419143
ко	0.008299794228975609	ОЛ	0.007951751690227412	к	0.008183677283672385	во	0.007928293229472336
к	0.008390522193874773	ec	0.008203578501673154	ко	0.008452232547083363	ec	0.008263351651283146
не	0.0087788378836432	пр	0.008378576794372738	не	0.00865909268241344	пр	0.008310302512938163
ро	0.008891340560118165	ли	0.008583452844362493	ро	0.009114547892657735	ли	0.00867097049565171
0	0.009203444759371291	ОТ	0.008660281363108653	0	0.0091472100192888	ал	0.008720055487381955
по	0.009217961233755158	ал	0.008664549614150106	ь	0.009237938148819536	ОТ	0.008875846982873607
Ь	0.009486516009856685	ОН	0.009121252475585604	по	0.009288745901356749	ОН	0.00901669956783866
Я	0.009838540513665447	qo	0.009304787270368094	Я	0.009756903049735347	ор	0.009334684949047644
и	0.00986757346243318	ла	0.009787099638052311	на	0.009891180681440835	ла	0.009614255988902524
на	0.009874831699625112	ен	0.009855391654715565	и	0.009956504934702966	ни	0.00975084031371712
ра	0.010281292982373371	ни	0.009868196407839924	ра	0.010232318448476403	ен	0.00991943658966014
но	0.010803886060192561	не	0.009953561428668989		0.010464582460075087	не	0.010226751320492984
СТ	0.012132143466316336	ко	0.01015416922761729	СТ	0.012094059666447105	ко	0.010265165661847089
Н	0.013148296673186984	ос	0.01061087208905279	Н	0.013239048661124993	ос	0.010585285173131303
В	0.01392129893412787	ОВ	0.010751724373420746	В	0.013661841744738222	ОВ	0.01064290668516246
то	0.014759625329796152	по	0.010760260875503653	то	0.014746950173925825	ро	0.010747479058848637
П	0.01489027359925095	-	0.010764529126545106	П	0.014854009366772093	по	0.010924611855092569
и	0.015877393857353865	на	0.011801714129618248	и	0.015915528482281702	на	0.011645947820519661
e	0.017060486519638974	pa	0.012074882196271257	С	0.017096808728771885	pa	0.012042896014512085
С	0.01716210184032604	но	0.012535853308748207	e	0.017118583479859263	но	0.012399295737075175
a	0.017423398379235636	СТ	0.014204739465956429	а	0.017447019308760525	СТ	0.014475804300272102
0	0.021625917713364955	то	0.01770897357098955	0	0.021591480265724546	то	0.017640719201835353
<u> </u>	3.321023317,13304333		1.31,,003,03,030333		5.522552 100205724540		1.31, 0.0, 13201033333

---H1 entropy---

H1:4.401484367688899H1 without spaces:4.460162425027035H1 with redundancy:0.1274503412019915H1 without spaces with redundancy:0.11581801114416501

---H2 block entropy---

H2 block bigrams:
4.003559323053849
H2 block bigrams without spaces:
4.133327222431862
H2 block bigrams with redundancy:
0.20633494760258309
H2 block bigrams without spaces with redundancy:
0.17333455551362764

---H2 cross entropy---

 H2 cross bigrams:
 4.0045129023504415

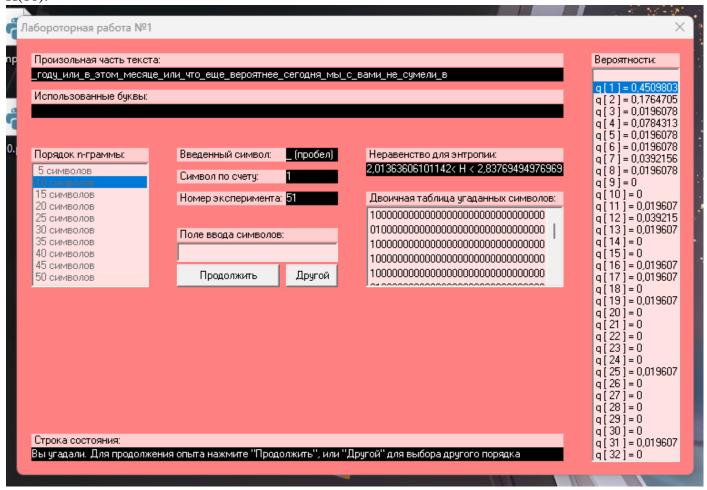
 H2 cross bigrams without spaces:
 4.133213304561436

 H2 cross bigrams with redundancy:
 0.2061459101733042

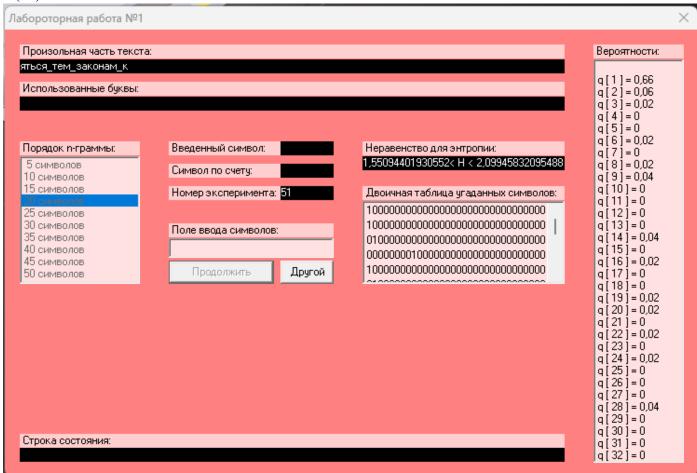
 H2 cross bigrams without spaces with redundancy:
 0.17335733908771278

#### CoolPinkProgram

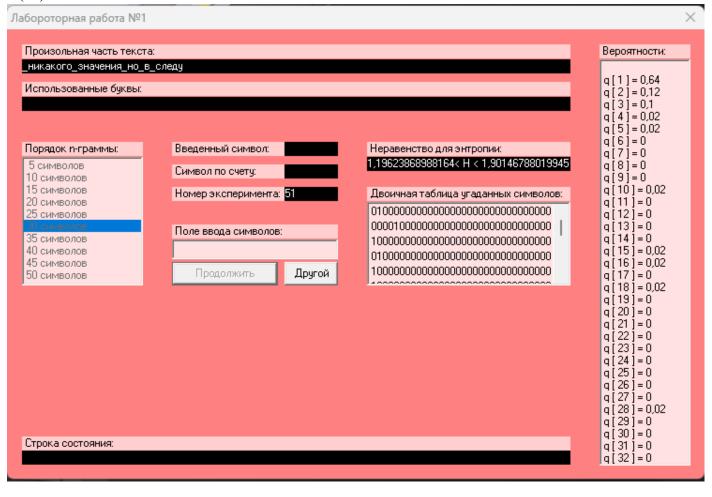
H(10):



#### H(20):



#### H(30):



## Результати:

2.0136 < H(10) < 2.8376

1.5509 < H(20) < 2.0994

1.1962 < H(30) < 1.9014

## Оцінка надлишковості R російської мови у різних моделях відкритого тексту:

	H8	R
H(10)	2,4256	0,51488
H(20)	1,82515	0,63497
H(30)	1,5488	0,69024

#### Висновки:

Протягом даної лабораторної роботи нам необхідно було навчитися працювати з великими масивами інформації, використовуючи математичні функції та формули, а саме ентропія, надлишковість тексту, порівняння різних текстів відносно ентропії та зрозуміння явища ентропії