# МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ "КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ ІМЕНІ ІГОРЯ СІКОРСЬКОГО" ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ

# КРИПТОГРАФІЯ КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

Виконали: Студентки групи ФБ-33 Яремко А.В та Журавльова М.В

### Варіант №3

### 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

1. Написати програми для підрахунку частот букв і частот біграм в тексті, а також підрахунку  ${\rm H_1}$  та  ${\rm H_2}$  за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення  ${\rm H_1}$  та  ${\rm H_2}$  на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення  ${\rm H_1}$  та  ${\rm H_2}$  на тому ж тексті, в якому вилучено всі пробіли.

Обраний текст: "Злочин і кара" Ф. М. Достоєвський (~ 2.1 МБ)

Перед початком обчислень було проведено попереднє очищення тексту за методичними вказівками, а результат збережено у файли clean\_with\_spaces.txt та clean\_no\_spaces.txt . Ймовірність появи кожного символу обчислювалась як відношення кількості входжень символу до загальної кількості символів у тексті:  $p_i = \frac{n_i}{N}$ , де  $n_i$  - кількість появ символу, N - довжина тексту.

Для біграм застосовувалися два підходи: перекривні та неперекривні біграми.

Для обчислення ентропії літер використовувалась формула:  $H_1 = -\sum_{i=1}^n p_i log(p_i)$ , а для

ентропії біграм: 
$$H_2 = -\frac{1}{2} \sum_{i,j} p_{ij} log(p_{ij})$$

Для оцінки надлишковості тексту була застосована формула:  $R=1-\frac{H}{\log |A|}$ , де |A| - розмір алфавіту тексту.

### Результат роботи коду:

```
--- ТЕКСТ БЕЗ ПРОБІЛІВ ---
Аналіз файлу: prestuplenie-i-nakazanie-fedor-dostoevskij.txt
Довжина сирого тексту: 1,248,151 символів
Після фільтрації: 969,927 символів
3бережено очищений текст -> results\clean_no_spaces.txt
 - EHTPONIÏ:
H1 = 4.448816
H2 (перекривні) = 4.131265
H2 (неперекривні)= 4.130387
 - НАДЛИШКОВОСТІ:
|A| = 31, log2(|A|) = 4.954196
R1 (H1) = 0.102011
R2 (H2, перекривні) = 0.166108
R2 (H2, неперекривні) = 0.166285
======== ПОРІВНЯННЯ РЕЗУЛЬТАТІВ =========
             | 3 пробілами | Без пробілів
H1
                                4.355119
                                                   4.448816
Н2 (перекривні)
                                  3.949231
                                                    4.131265
Н2 (неперекривні)
                                  3.948646
                                                    4.130387
R1 (H1)
                                  0.128976
                                                    0.102011
R2 (H2, перекривні)
                                  0.210154
                                                    0.166108
R2 (H2, неперекривні)
                                   0.210271
                                                     0.166285
```

## Монограми з пробілами:

0,167721823 0 0,09501814 e 0,072889887 a 0,065776377 и 0,055329212 H 0,053166842 c 0,045092278 в 0,038981009 л 0,037626095 р 0,035563263 к 0,027648303 д 0,026359461 м 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 b 0,018311498 я 0,01744655 ч 0,014608868 г 0,014414942 б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002503029 ф 0,001180723	Символ	Ймовірність
е       0,072889887         а       0,065776377         и       0,055329212         н       0,054290931         т       0,053166842         с       0,045092278         в       0,038981009         л       0,037626095         р       0,035563263         к       0,027648303         д       0,026359461         м       0,02635947         п       0,022923696         ь       0,018311498         я       0,014608868         г       0,014608868         г       0,01441942         б       0,013722468         з       0,012867817         ж       0,009202085         й       0,008731856         х       0,007284269         ш       0,006584073         ю       0,00246262         щ       0,002662633         ц       0,002503029		
а 0,065776377 и 0,055329212 н 0,054290931 т 0,053166842 с 0,045092278 в 0,038981009 л 0,037626095 р 0,035563263 к 0,027648303 д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,01441942 б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,002662633 ц 0,002503029	0	0,09501814
и         0,055329212           н         0,054290931           т         0,053166842           с         0,045092278           в         0,038981009           л         0,037626095           р         0,035563263           к         0,027648303           д         0,026359461           м         0,02618527           у         0,024091547           п         0,022923696           b         0,018311498           я         0,014608868           г         0,01440942           б         0,014215866           ы         0,013722468           з         0,012867817           ж         0,009202085           й         0,008731856           х         0,007284269           ш         0,006584073           ю         0,0024752923           э         0,002662633           ц         0,002503029	e	0,072889887
H         0,054290931           T         0,053166842           C         0,045092278           B         0,038981009           Л         0,037626095           P         0,035563263           К         0,027648303           Д         0,026359461           М         0,02618527           У         0,024091547           П         0,022923696           Б         0,018311498           Я         0,014608868           Г         0,01440942           Б         0,014215866           Б         0,013722468           З         0,012867817           Ж         0,009202085           Й         0,008731856           х         0,007284269           Ш         0,006584073           Ю         0,004752923           Э         0,002662633           Ц         0,002503029	a	0,065776377
т 0,053166842 с 0,045092278 в 0,038981009 л 0,037626095 р 0,035563263 к 0,027648303 д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,002662633 ц 0,002503029	И	0,055329212
С 0,045092278 В 0,038981009 Л 0,037626095 р 0,035563263 К 0,027648303 Д 0,026359461 М 0,02618527 У 0,024091547 П 0,022923696 Ь 0,018311498 Я 0,01744655 Ч 0,014608868 Г 0,014215866 Ы 0,013722468 З 0,012867817 Ж 0,009202085 Й 0,008731856 X 0,007284269 Ш 0,006584073 Ю 0,004752923 Э 0,002662633 Ц 0,002503029	н	0,054290931
В 0,038981009 л 0,037626095 р 0,035563263 к 0,027648303 д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,002846262 щ 0,002503029	т	0,053166842
л 0,037626095 р 0,035563263 к 0,027648303 д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,002462633 ц 0,002503029	С	0,045092278
р 0,035563263 к 0,027648303 Д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,002462623 ц 0,002503029	В	0,038981009
к 0,027648303 д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,0024752923 э 0,002662633 ц 0,002503029	Л	0,037626095
д 0,026359461 м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,0024752923 э 0,002662633 ц 0,002503029	p	0,035563263
м 0,02618527 у 0,024091547 п 0,022923696 ь 0,018311498 я 0,01744655 ч 0,014414942 б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,00246262 щ 0,002650332 ц 0,002503029	К	0,027648303
у 0,024091547 п 0,022923696 b 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,0026503029	Д	0,026359461
п 0,022923696 b 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,0024752923 э 0,002662633 ц 0,002503029	M	0,02618527
b 0,018311498 я 0,01744655 ч 0,014608868 г 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002650332 ц 0,002503029	у	0,024091547
я 0,01744655 ч 0,014608868 г 0,014414942 б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002650332 ц 0,002503029	п	0,022923696
ч 0,014608868 г 0,014414942 б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	ь	0,018311498
г 0,014414942 6 0,014215866 ы 0,013722468 3 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	я	0,01744655
б 0,014215866 ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	ч	0,014608868
ы 0,013722468 з 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	Γ	0,014414942
3 0,012867817 ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	б	0,014215866
ж 0,009202085 й 0,008731856 х 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	ы	0,013722468
й     0,008731856       х     0,007284269       ш     0,006584073       ю     0,004752923       э     0,002846262       щ     0,002662633       ц     0,002503029	3	0,012867817
x 0,007284269 ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	ж	0,009202085
ш 0,006584073 ю 0,004752923 э 0,002846262 щ 0,002662633 ц 0,002503029	й	0,008731856
ю 0,004752923 ∋ 0,002846262 щ 0,002662633 ц 0,002503029	X	0,007284269
э 0,002846262 щ 0,002662633 ц 0,002503029	Ш	0,006584073
щ 0,002662633 ц 0,002503029	ю	
ц 0,002503029	Э	0,002846262
	Щ	
ф 0,001180723	ц	
	ф	0,001180723

### Монограми без пробілів:

	5
Символ	Ймовірність
0	0,114166324
e	0,087578756
a	0,079031721
И	0,06647923
Н	0,065231713
Т	0,063881096
С	0,054179335
В	0,046836515
Л	0,045208557
p	0,04273002
К	0,033220026
Д	0,031671456
M	0,031462162
У	0,028946508
п	0,02754331
ь	0,022001656
я	0,020962402
ч	0,017552867
Γ	0,01731986
6	0,017080667
ы	0,016487839
3	0,015460957
ж	0,011056502
й	0,010491511
x	0,008752205
Ш	0,007910905
ю	0,005710739
∋	0,003419845
щ	0,00319921
ц	0,003007443
ф	0,001418663

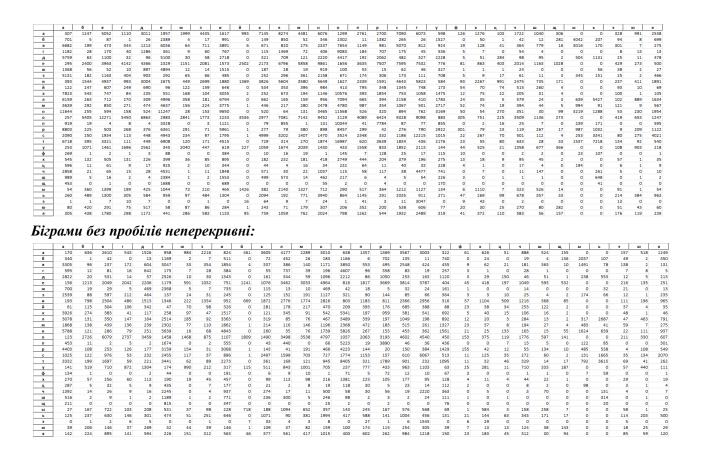
# Біграми з пробілами перекривні:

		a	6	8	r	А	e	ж	3	и	й	K	л	м	н	0	n	р	С	т	у	ф	x	ц	ч	w	щ	ы	ь	9	ю	я
	0	2988	7651	19237	3798	9315	4591	2567	4990	13545	13	9307	2853	6720	18183	12970	18775	5981	18874	11407	5091	619	1852	362	6880	651	71	0	0	3274	76	2820
9	19623	3	597	3123	740	2086	1417	1809	3944	126	993	6163	7965	3753	4231	12	926	2158	5212	4910	87	71	1098	71	1044	974	294	0	0	5	988	2232
6	514	695	1	63	0	24	2372	4	2	954	0	136	850	46	319	2208	0	1379	232	4	1481	0	48	1	20	12	281	4042	207	0	8	664
8	7135	6594	4	29	26	887	5856	0	563	3597	0	141	710	115	1725	7229	418	671	4275	266	726	0	71	21	196	737	12	3016	170	0	0	238
r	1311	1163	0	8	12	1218	339	0	1	697	0	59	1459	39	290	8953	0	644	51	6	500	0	0	0	35	2	0	0	0	0	12	0
Д	1111	5738	31	998	20	47	5071	23	16	2640	0	238	682	87	2067	4359	99	2023	368	286	2208	1	47	282	71	91	2	504	1131	1	11	466
c	21964	41	1455	1647	3702	3283	1930	925	1338	126	2502	1341	6411	4914	8110	334	1281	6866	5353	6266	195	28	515	372	1420	1088	1022	0	0	0	263	253
ж	739	1361	25	0	5	873	4882	11	0	1447	0	107	18	4	828	26	0	5	12	0	316	0	0	1	6	0	0	0	56	0	2	0
3	1477	5114	154	1027	366	815	261	52	12	443	0	141	274	271	2023	596	0	274	35	8	665	0	0	3	22	7	0	345	151	0	2	458
и	20305	137	722	2805	628	1940	2250	253	2129	749	1369	2812	5344	2963	3596	273	275	1005	2458	4724	10	26	2049	948	1859	674	167	0	0	1	399	1610
й	7682	0	0	0	0	280	4	0	4	0	0	137	239	50	396	8	1	2	556	376	0	0	0	46	181	212	2	0	0	0	0	0
K	5594	7768	0	278	0	2	466	42	19	2643	0	11	581	3	599	10240	0	1688	247	623	1348	0	35	18	0	16	0	0	0	0	0	0
л	7593	6003	27	14	139	21	4652	330	5	6239	0	321	75	16	286	6110	62	0	1507	105	1657	2	0	0	219	2	1	639	5427	0	889	1508
M	8972	3516	17	13	79	0	4485	0	0	3097	0	84	105	37	1686	4114	62	68	116	30	2443	11	0	10	52	6	1	994	91	1	3	423
н	4782	11587	11	27	75	345	11157	3	25	9352	0	302	0	0	3059	11275	0	95	359	525	3008	38	5	263	222	8	97	3219	1288	0	230	1913
0	27240	14	4039	9230	4976	5634	2260	2321	1089	732	3534	1852	6674	6139	6997	261	1572	5732	6826	7368	101	224	526	180	2482	1065	266	0	0	22	646	731
n	79	917	13	0	1	0	3326	0	0	1110	0	77	851	0	129	10341	30	7783	80	74	854	0	0	16	22	7	0	239	171	0	0	595
p	895	8782	193	436	256	346	6349	287	54	5873	0	238	74	362	819	8396	120	13	186	751	2907	299	72	22	95	265	17	987	1032	0	208	1111
c	4366	2010	99	1579	31	219	4875	25	10	1603	0	4763	3114	1226	993	3271	1921	238	822	11894	837	6	229	63	407	94	1	253	3341	0	271	3988
T	6697	6521	14	2699	10	114	6717	0	25	4013	0	407	236	47	1279	16566	76	3460	1229	74	1997	2	16	66	353	10	32	1537	7218	10	91	444
y	7765	67	887	718	1578	2184	226	1917	266	8	227	707	1579	1697	571	4	888	638	1137	1626	1	16	450	11	916	656	364	0	0	1	901	70
Ф.	49	268	0	0	0	0	88	0	0	379	0	3	15	6	1	145	0	128	22	15	115	10	0	0	1	1	0	23	107	0	0	0
X	3623	498	0	198	0	0	326	0	8	587	0	1	145	68	156	2539	0	96	30	7	183	0	0	0	1	21	0	0	0	2	0	0
ц	440	586	0	24	0	0	911	0	0	312	0	25	0	3	0	198	0	0	0	5	218	0	0	0	0	1	0	194	0	0	0	0
4	534	2844	0	2	0	0	4524	0	0	1810	0	533	21	4	986	85	0	92	0	4451	731	0	0	0	0	147	0	0	261	0	0	0
	71	988	0	9	0	0	2383	0	0	1549	0	494	571	12	456	213	0	0	0	53	225	0	0	0	0	0	0	0	648	0	1	0
щ	6	452	0	1	0	0	1685	0	0	688	0	0	0	0	55	2	0	3	0	0	170	0	0	0	0	0	0	0	41	0	0	0
ы	4990	0	156	910	121	171	909	24	58	8	1426	172	2095	1091	211	0	115	232	733	774	2	0	1070	1	191	517	12	0	0	0	0	3
ь	13100	0	101	32	146	21	521	0	176	131	0	1376	0	387	2491	5	1	0	939	189	0	13	128	69	63	333	29	0	0	0	381	708
	28	0	0	5	9	6	0	0	1	0	16	63	9	7	24	0	38	3	10	3046	0	9	43	0	0	0	0	0	0	0	0	0
10	3507	0	347	4	13	349	0	18	8	3	1	17	16	47	35	0	0	71	198	392	0	0	6	13	111	65	276	0	0	0	42	0
	13269	0	23	312	68	539	112	113	253	19	95	193	883	402	669	0	55	97	683	1705	0	0	229	78	156	11	156	0	0	0	115	97

# Біграми з пробілами неперекривні:

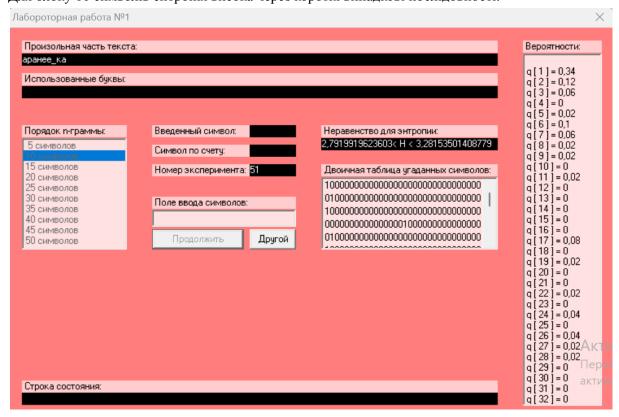
		a	6	0	г	А	e	ж	3	и	й	K	л	м	н	0	п	р	c	т	у	ф	×	ц	ч	w	щ	ы	b	9	ю	я
	0	1496	3857	9643	1881	4685	2236	1297	2504	6733	6	4702	1425	3347	9101	6524	9341	2955	9309	5669	2602	314	945	170	3435	330	32	0	0	1640	35	1373
a	9810	2	285	1551	360	1052	681	860	1918	62	484	3065	3920	1905	2165	5	473	1052	2603	2454	48	41	547	35	528	488	151	0	0	3	505	1105
- 6	241	338	0	31	0	10	1190	2	2	511	0	73	410	24	172	1104	0	667	112	2	749	0	28	1	8	4	133	2006	109	0	3	329
В	3523	3294	2	14	12	447	2890	0	283	1817	0	74	339	56	833	3636	207	334	2171	119	374	0	33	10	90	389	3	1517	89	0	0	139
r	655	568	0	5	5	597	171	0	1	337	0	35	730	21	149	4562	0	331	22	2	244	0	0	0	13	2	0	0	0	0	7	0
Д	546	2866	18	517	10	24	2562	8	7	1337	0	121	320	49	1045	2175	42	1004	183	145	1107	0	25	148	32	51	1	262	542	0	5	234
e	11016	25	713	836	1859	1652	963	456	662	64	1260	691	3132	2511	4003	165	654	3440	2685	3198	97	14	260	192	713	550	521	0	0	0	137	136
340	373	680	12	0	3	447	2458	7	0	718	0	51	7	3	426	10	0	1	6	0	159	0	0	1	4	0	0	0	27	0	1	0
3	754	2577	75	496	175	405	130	29	6	215	0	74	140	127	1021	296	0	123	14	5	343	0	0	8	17	5	0	190	82	0	1	245
и	10179	62	375	1387	316	951	1118	128	1076	366	669	1399	2640	1476	1805	146	136	473	1237	2366	6	12	995	455	901	333	84	0	0	0	197	837
й	3868	0	0	0	0	133	1	0	2	0	0	74	115	28	194	5	0	2	308	188	0	0	0	21	96	102	2	0	0	0	0	0
к	2779	3842	0	144	0	2	241	20	8	1330	0	6	295	1	309	5092	0	809	129	323	671	0	15	9	0	11	0	0	0	0	0	0
л	3828	3085	16	5	77	16	2373	157	1	3143	0	169	36	7	145	3058	29	0	779	43	821	0	0	0	101	1	1	310	2717	0	434	751
M	4500	1757	8	6	28	0	2277	0	0	1500	0	45	40	19	844	2033	35	30	59	16	1195	3	0	5	23	4	1	512	49	1	0	208
н	2456	5796	7	12	38	160	5509	3	14	4706	0	145	0	0	1547	5663	0	51	190	249	1467	17	1	139	118	3	44	1583	641	0	117	965
0	13546	6	2047	4616	2461	2801	1133	1156	558	373	1755	912	3298	3073	3433	117	795	2886	3397	3748	57	105	245	93	1244	532	128	0	0	8	312	381
n	37	458	7	0	1	0	1666	0	0	550	0	44	457	0	68	5196	19	3892	40	31	379	0	0	6	9	3	0	125	83	0	0	305
р	448	4450	87	228	135	157	3187	139	25	2954	0	126	36	168	428	4261	48	6	90	377	1474	128	36	10	57	130	9	486	522	0	108	541
c	2230	992	43	774	14	112	2398	12	6	814	0	2369	1555	626	482	1669	963	123	437	5886	409	3	130	32	210	38	1	130	1693	0	129	2028
т	3315	3291	8	1355	5	55	3316	0	10	2060	0	219	115	21	635	8261	32	1785	592	33	995	2	6	33	167	5	17	775	3653	4	50	208
у	3957	32	438	359	779	1068	102	975	127	6	115	339	836	830	270	0	451	294	584	807	1	8	206	6	464	319	189	0	0	0	478	35
у ф х	27	134	0	0	0	0	52	0	0	214	0	2	9	3	1	70	0	66	12	5	44	4	0	0	0	0	0	12	54	0	0	0
×	1857	269	0	94	0	0	155	0	0	300	0	1	75	28	89	1248	0	48	17	4	90	0	0	0	1	13	0	0	0	2	0	0
ц	224	307	0	12	0	0	477	0	0	152	0	13	0	1	0	87	0	0	0	1	98	0	0	0	0	1	0	103	0	0	0	0
4	268	1463	0	0	0	0	2279	0	0	899	0	291	10	3	479	32	0	35	0	2219	357	0	0	0	0	67	0	0	138	0	0	0
ш	34	485	0	5	0	0	1194	0	0	740	0	255	278	11	240	98	0	0	0	24	120	0	0	0	0	0	0	0	325	0	1	0
щ	3	222	0	0	0	0	828	0	0	368	0	0	0	0	31	1	0	1	0	0	92	0	0	0	0	0	0	0	19	0	0	0
ы	2490	0	73	451	64	84	446	10	31	4	698	80	1059	568	104	0	63	105	366	388	2	0	514	1	98	274	6	0	0	0	0	2
b	6535	0	51	21	78	12	257	0	88	69	0	682	0	188	1245	3	1	0	435	91	0	12	75	32	29	170	12	0	0	0	181	330
	15	0	0	2	2	2	0	0	1	0	8	32	5	3	14	0	20	1	5	1525	0	4	20	0	0	0	0	0	0	0	0	0
ю	1752	0	177	3	7	176	0	11	3	3	1	7	9	16	22	0	0	36	95	196	0	0	2	6	59	31	124	0	0	0	19	0
я	6608	0	9	166	32	285	50	60	115	10	41	89	455	205	329	0	30	44	365	823	0	0	115	33	68	7	79	0	0	0	64	49

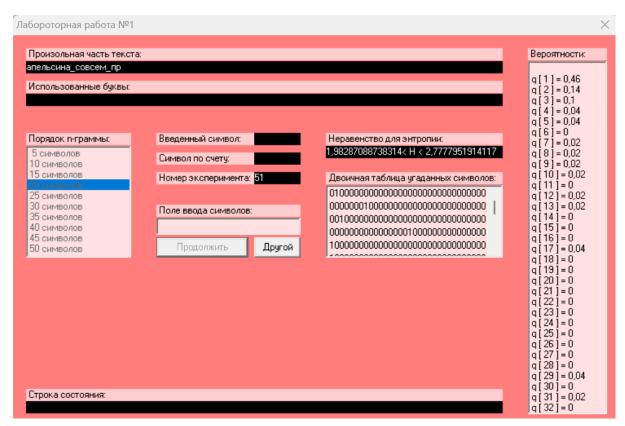
Біграми без пробілів перекривні:



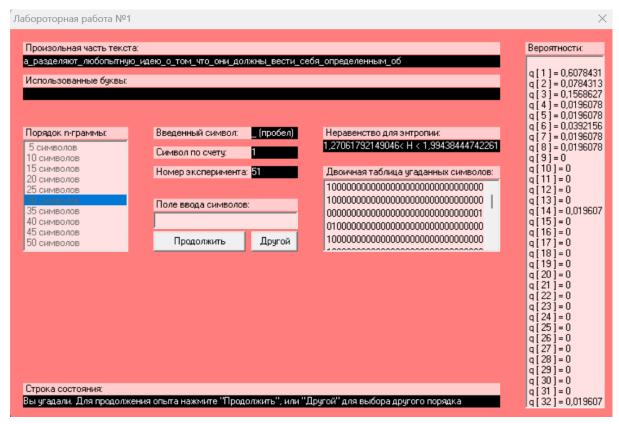
2. За допомогою програми CoolPinkProgram оцінити значення Н <sup>(10)</sup>, Н <sup>(20)</sup>, Н <sup>(30)</sup>

Для блоку 10 символів ентропія висока через короткі випадкові послідовності.





Для 20 символів ентропія зменшується через повтори символів.



Для 30 символів ентропія ще зменшується, оскільки більший блок охоплює більше символів які легше передбачити

Розрахунок ентропії для блоків довжини 10, 20 і 30 символів дозволяє оцінити ступінь випадковості тексту на різних масштабах. Зі збільшенням довжини блоку спостерігається **зменшення ентропії**, що вказує на прояв закономірностей у тексті. Виконання не менше 50 експериментів гарантує статистично надійні результати. Порівняння з CoolPinkProgram дозволяє перевірити коректність власних обчислень. Таку саму закономірність бачимо і в Н1 та H2. H1 > H2 блоки із двох символів менш випадкові, бо з'являються повторювані комбінації.

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

Аналогічно до формули з першого пункту розрахуємо надлишковість: |A| = 32, log 32 = 5

```
$$2.7919919623603 < H < 3.28153501408779$$
0,500802720470174 > R(10) > 0,372661087473526
```

1.98287088738314 < H < 2.7777951914117 0.583395773390644 > R(20) > 0.474633931031182

1.27061792149046 < H < 1.9943844742261 0,378946628238864 > R(30) > 0,3020768941536886

### Висновки.

Аналізуючи результати з першого пункту можна помітити, що H1 та H2 для тексту без пробілів  $\epsilon$  вищими, що можна пояснити тим, що при вилученні пробіла, який  $\epsilon$  найчастішим символом, довжина текста зменшується, а ймовірності стають більш рівномірно розподіленими, що в свою чергу збільшу $\epsilon$  ентропію. Якщо подивитись на надлишковість, то R1 та R2 навпаки для тексту з пробілами вища, що свідчить про те, що наявність пробілу підвищу $\epsilon$  загальну надлишковість.

Аналізуючи результати роботи з CoolPinkProgram можна побачити, що ентропії вийшли меншого розміру, ніж ентропії з пробілами з першого пункту, а надлишковості більше. Це може свідчити про те, що надлишковість зростає зі збільшенням порядку моделі.

Загалом, виконання комп'ютерного практикуму допомогло нам узагальнити поняття ентропії та надлишковості. Ми навчились на практиці визначати частоти п-грам, розраховувати значення ентропії, надлишковості.