

Експериментальна оцінка ентропії на символ джерела відкритого тексту

1 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.

2 Хід роботи

Для виконання поставленого завдання, після короткого аналізу, я вирішив розбити його на три різних частини:

- робота з файлами та вхідним текстом
- робота з літерами
- робота з біграмами

2.1 Робота з файлами та вхідним текстом

Не буду вдаватись в деталі реалізації, так як це не є настільки важливим в даній роботі, лише хочу зазначити, що вхідний текст має назву `boloto.txt`, цей же текст, але вже опрацьований препроцесором має назву `boloto_processed.txt`, а вже оброблений файл та ще й без пробілів має назву `boloto_without_spaces.txt`.

2.2 Робота з літерами

В даній секції ми розв'язуємо декілька задач:

- підрахунок загальної к-ті літер у тексті

```
1 fn letters_count(letter_frequencies: &HashMap<char, i64>) -> i64
2 {
3     let mut count = 0;
4     for (_key, _value) in letter_frequencies {
5         count += _value;
6     }
7     count
8 }
```

- підрахунок к-ті кожної літери

```

1 fn get_letter_frequency(text: &str) -> HashMap<char, i64> {
2     let mut frequencies: HashMap<char, i64> = HashMap::new();
3
4     for c in text.chars() {
5         *frequencies.entry(c).or_insert(0) += 1;
6     }
7
8     frequencies
9 }

```

- підрахунок ймовірності зустріти кожну окрему літеру

```

1 fn count_letters_probabilities(letter_frequencies: &HashMap<char,
2     i64>) -> HashMap<char, f64> {
3     let mut probabilities: HashMap<char, f64> = HashMap::new();
4     let number_of_characters = letters_count(letter_frequencies)
5     as f64;
6
7     for (_key, _value) in letter_frequencies {
8         probabilities.insert(*_key, (*_value as f64) /
9         number_of_characters);
10    }
11    probabilities
12 }

```

- вивід знайдених значень

```

1 fn print_letters_probabilities(probabilities: &HashMap<char, f64
2     >) {
3     let mut sorted_probabilities: Vec<(&char, &f64)> =
4     probabilities.iter().collect();
5     sorted_probabilities.sort_by(|a, b| b.1.partial_cmp(a.1).
6     unwrap());
7     for (&letter, &probability) in sorted_probabilities {
8         println!("{}", letter, probability);
9     }
10    println!();
11 }
12
13 fn print_letter_frequencies(letter_frequencies: &HashMap<char,
14     i64>) {
15     let mut sorted_frequencies: Vec<(&char, &i64)> =
16     letter_frequencies.iter().collect();
17     sorted_frequencies.sort_by_key(|&(_, frequency)| *frequency);
18     for (&letter, &frequency) in sorted_frequencies.iter().rev()
19     {
20         println!("{}", letter, frequency);
21     }
22     println!();
23 }

```

2.3 Робота з біграмами

Аналогічні задачі доводиться розв'язувати і у випадку з біграмами, правда реалізація буде досить сильно відрізнятись.

- підрахунок загальної к-ті біграм у тексті

```

1 fn bigram_count(bigram_frequencies: &HashMap<String, i64>) -> i64
2 {
3     let mut count = 0;
4     for (_key, _value) in bigram_frequencies {
5         count += _value;
6     }
7     count
8 }

```

- підрахунок к-ті кожної з біграм

```

1 fn get_bigram_frequency(text: &str) -> HashMap<String, i64> {
2     let mut frequencies: HashMap<String, i64> = HashMap::new();
3
4     let mut chars = text.chars().peekable();
5     while let (Some(curr), Some(&next)) = (chars.next(), chars.
6     peek()) {
7         if curr.is_alphabetic() && next.is_alphabetic() {
8             let bigram = format!("{}", curr.to_lowercase(),
9             next.to_lowercase());
10            *frequencies.entry(bigram).or_insert(0) += 1;
11        } else if curr.is_alphabetic() && next.is_whitespace() {
12            let bigram = format!("{}", curr.to_lowercase());
13            *frequencies.entry(bigram).or_insert(0) += 1;
14        } else if curr.is_whitespace() && next.is_alphabetic() {
15            let bigram = format!("{}", next.to_lowercase());
16            *frequencies.entry(bigram).or_insert(0) += 1;
17        }
18    }
19    frequencies
20 }

```

- підрахунок ймовірності зустріти кожну окрему біграму

```

1 fn count_bigram_probabilities(bigram_frequencies: &HashMap<String
2 , i64>) -> HashMap<String, f64> {
3     let mut probabilities: HashMap<String, f64> = HashMap::new();
4     let number_of_bigrams = bigram_count(bigram_frequencies) as
5     f64;
6
7     for (_key, _value) in bigram_frequencies {
8         probabilities.insert(_key.clone(), (*_value as f64) /
9         number_of_bigrams);
10    }
11
12    probabilities
13 }

```

- вивід знайдених значень

3 Результати роботи

3.1 Текст з пробілами

3.1.1 Робота з літерами

Літера	Кількість	Ймовірність
' '	248583	0.0149
а	115586	0.06929
б	23656	0.01418
в	65348	0.03918
г	23164	0.01389
д	43783	0.02625
е	113939	0.06831
ё	2136	0.00128
ж	12555	0.00753
з	21398	0.01283
и	107358	0.06436
й	17005	0.01019
к	47901	0.02872
л	56102	0.03363
м	46423	0.02783
н	92312	0.05534
о	158172	0.09482
п	41612	0.02495
р	67768	0.04063
с	77815	0.04665
т	96940	0.05811
у	37574	0.02253
ф	4406	0.00264
х	11742	0.00704
ц	6460	0.00387
ч	20695	0.01241
ш	9662	0.00579
щ	4685	0.00281
ы	25661	0.01538
ь	24874	0.01491
э	6029	0.00361
ю	10427	0.00625
я	25800	0.01547
ъ	513	0.00031

3.1.2 Робота з біграмами

	а	б	в	г	д	е	ж	з	и	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	ш	щ	э	ю	я									
а	0	4885	7320	24519	5747	10578	4509	2332	6963	16014	35115169	4873	10835	22680	14372	28353	377	22587	18313	61232	1968	1720	721	7510	14781	260	8	1	0	5096	912	2347	1	0			
б	23332	0	226	2556	9333	1102	3124	2353	4133	392	1590	7214	10039	5747	9654	194	1720	4654	4803	9438	2851	351	1587	1281	1583	1638	500	0	0	541	2346	2656	54	0			
в	0	758	1028	1455	0	15	2459	21	31	1670	0	457	1879	1481	479	5363	5	1369	393	4	61232	5	68	2	0	80	39	774	335	3088	50	8	177	60	13		
г	13507	15331	45	145	115	327	6189	8	272	4731	0	766	2165	197	2204	946	229	1284	4033	522	937	13	35	115	242	472	21	8	4834	424	10	1	431	75			
д	0	246	268	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92	26	92			
е	д	1800	7324	47	1064	48	518	967	224	100	5128	0	360	1185	327	2077	5516	525	1338	885	152	2495	0	38	69	44	96	0	0	0	828	482	15	35	551	173	
ж	д	27004	402	1172	2434	4112	592	1150	838	1942	18	504	2788	7483	6394	12567	714	1365	890	8034	9624	91	220	974	407	1284	1033	1043	0	0	30	210	416	2	0		
з	ц	237	1331	97	0	5	1558	4395	18	0	1826	0	342	2	16	1842	38	2	14	10	3	736	0	5	13	2	0	0	0	1	46	4	0	17			
и	ш	246	288	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950	68	950		
к	ш	27699	946	596	3050	768	2922	4523	457	4171	3401	2485	4590	5487	4793	8143	1310	541	3132	5282	9360	42	255	3078	1601	2420	704	245	0	0	36	1195	410	0	9		
л	и	11848	29	12	43	42	269	17	3	48	35	0	362	59	1202	632	185	68	19	1087	1326	18	15	0	4	123	643	151	10	0	7	2	1	13	0		
м	и	11848	29	12	43	42	269	17	3	48	35	0	362	59	1202	632	185	68	19	1087	1326	18	15	0	4	123	643	151	10	0	7	2	1	13	0		
н	и	11848	29	12	43	42	269	17	3	48	35	0	362	59	1202	632	185	68	19	1087	1326	18	15	0	4	123	643	151	10	0	7	2	1	13	0		
о	и	4024	6764	47	18	151	78	1807	668	12	11211	0	501	841	43	675	6992	51	42	569	101	1962	21	3	3	101	1	0	1	0	51	106	8223	10	2223	1558	1211
п	и	12343	4922	149	147	132	21	7267	66	38	5418	0	204	593	385	1221	6928	507	57	367	92	2180	79	5	25	45	8	19	0	2485	222	214	7	398	31	0	
р	и	1390	1730	18	126	21	12302	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	
с	и	1390	1730	18	126	21	12302	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	14268	0	131	
т	и	165	2465	8	12	2	1	25	1	25403	0	108	1233	9	209	14298	338	129477	77	86	2617	1	0	1	180	68	13	0	0	0	466	50	10	3	197	4	
у	и	1692	15052	240	625	781	339	10956	837	13	7628	0	483	141	777	1436	1138	194	316	733	1219	430	128	155	59	44	288	36	0	0	2347	771	27	170	1755	66	
ф	и	1305	10649	134	322	61	282	9845	11	84	8768	0	1155	355	199	2168	21426	371	5779	2731	142	2349	29	24	42	156	14	15	8	2286	9895	24	333	752	79		
х	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
ц	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
ч	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
ш	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
щ	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
э	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
ю	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		
я	и	9044	212	1414	578	1317	3377	554	1737	385	62	226	972	1115	1683	71	32	1399	1625	399	3756	20	46	280	103	1891	574	489	0	0	0	20	1801	67	3		

Рис. 1: к-ть біграм

	а	б	в	г	д	е	з	и	к	л	м	н	о	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я					
0	0.000	0.003	0.004	0.015	0.003	0.006	0.003	0.001	0.004	0.010	0.000	0.008	0.002	0.007	0.014	0.009	0.017	0.000	0.014	0.006	0.004	0.001	0.001	0.005	0.001	0.000	0.000	0.000	0.000	0.003	0.001	0.001	0.000	
1	0.014	0.000	0.000	0.000	0.000	0.002	0.002	0.001	0.002	0.001	0.004	0.004	0.000	0.001	0.003	0.004	0.006	0.000	0.000	0.000	0.000	0.001	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.001	0.002	0.000	0.000	
2	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
3	0.008	0.007	0.000	0.000	0.000	0.000	0.000	0.004	0.000	0.000	0.003	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.001	0.002	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.003	0.000	0.000	0.000
4	0.001	0.001	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000	0.001	0.000	0.000	0.007	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
6	0.016	0.000	0.001	0.001	0.002	0.004	0.001	0.001	0.001	0.000	0.003	0.002	0.004	0.004	0.000	0.001	0.005	0.005	0.006	0.000	0.000	0.001	0.000	0.000	0.001	0.000	0.001	0.001	0.001	0.001	0.000	0.000	0.000	0.000
7	0.000	0.001	0.000	0.000	0.000	0.001	0.003	0.000	0.001	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
8	0.001	0.005	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.000	0.001	0.000	0.000	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
9	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
10	0.007	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
11	0.004	0.007	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.003	0.000	0.000	0.001	0.000	0.000	0.000	0.002	0.001	0.001														
12	0.001	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
13	0.007	0.000	0.000	0.000	0.000	0.000	0.004	0.000	0.000	0.003	0.000	0.000	0.000	0.000	0.001	0.004	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000
14	0.002	0.011	0.000	0.000	0.000	0.001	0.007	0.000	0.009	0.000	0.001	0.000	0.000	0.002	0.011	0.000	0.000	0.000	0.001	0.002	0.002													
15	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
16	0.001	0.001	0.000	0.000	0.000	0.000	0.002	0.000	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
17	0.001	0.009	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.004	0.000	0.000	0.000	0.000	0.001	0.009	0.000	0.000	0.000	0.001	0.001	0.003												
18	0.003	0.002	0.000	0.001	0.000	0.000	0.004	0.000	0.000	0.003	0.000	0.004	0.002	0.001	0.003	0.001	0.000	0.000	0.002	0.011	0.001													
19	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
20	0.005	0.000	0.001	0.000	0.001	0.002	0.000	0.001	0.000	0.000	0.000	0.001	0.001	0.001	0.000	0.001	0.001	0.001	0.001	0.002	0.000													
21	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
22	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
23	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
24	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
25	0.003	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
26	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
27	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
28	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
29	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
30	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
31	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
32	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
33	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
34	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
35	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.00

Рис. 2: ймовірність появи біграм

3.2 Текст без пробілів

3.2.1 Робота з літерами

Літера	Кількість	Ймовірність
а	115586	0.08143
б	23656	0.01667
в	65348	0.04604
г	23164	0.01632
д	43783	0.03084
е	113939	0.08027
ё	2136	0.00150
ж	12555	0.00884
з	21398	0.01507
и	107358	0.07563
й	17005	0.012
к	47901	0.03374
л	56102	0.03952
м	46423	0.0327
н	92312	0.06503
о	158172	0.11143
п	41612	0.02931
р	67768	0.04774
с	77815	0.05482
т	96940	0.06829
у	37574	0.02647
ф	4406	0.0031
х	11742	0.00827
ц	6460	0.00455
ч	20695	0.01458
ш	9662	0.00681
щ	4685	0.0033
ы	25661	0.01808
ь	24874	0.01752
э	6029	0.00425
ю	10427	0.00735
я	25800	0.01818
ъ	513	0.00036

3.2.2 Робота з біграмами

	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	я	е		
а	599	3152	3886	1726	4145	3671	2529	4765	1880	1595	8508	10375	6942	11916	11380	4501	4045	8933	10541	903	563	1768	1328	2073	1779	504	0	0	0	625	2393	2935	54	
б	1249	53	185	1	35	2487	22	43	1734	0	7	184	1854	77	555	5451	0	7	1384	432	18	2053	10	74	3	86	39	774	335	3088	56	154	178	63
в	11826	485	86	627	818	6408	85	551	5538	5	1788	2303	1038	3149	10044	116	2157	5460	1287	1238	187	117	215	508	645	23	9	4834	424	385	21	496	75	
г	788	67	18	3	383	185	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
д	7383	94	1261	83	583	708	232	153	5232	1	458	1220	395	10773	5611	735	1545	1018	203	2537	23	50	72	131	104	1	0	50	828	482	54	36	585	173
е	820	2153	5122	4753	734	3191	10671	2862	1779	4504	4056	7916	7643	1440	2380	4727	9968	10586	16681	815	381	1229	487	2661	1184	1046	0	0	50	238	649	208	10	
ж	1328	195	21	11	1578	4401	133	3	1838	0	355	8	21	1808	52	24	19	1	240	740	31	1	5	20	9	1	0	1	40	3	9	1	0	
з	788	67	18	3	383	185	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
и	1473	1450	6016	1373	4143	5101	673	5089	5261	2468	6108	5899	5924	10799	3001	3821	4075	774	10521	788	461	3300	1677	3552	247	12	0	0	558	1229	1366	9	0	
й	3673	550	986	381	431	515	162	335	836	3	1073	212	602	1398	783	1588	766	2380	1590	255	148	57	384	875	258	12	0	1	234	18	115	0	0	
к	1788	67	18	3	383	185	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
л	6855	170	491	238	120	7653	694	133	11463	0	666	892	200	1007	7627	4525	167	932	23	20651	14	24	11	389	30	1	1	516	823	291	8231	2614	121	
м	4868	1485	1253	486	629	6468	441	351	6284	7	994	743	834	3339	7582	1872	525	1433	449	3113	175	74	69	769	63	21	0	2485	222	454	17	507	31	
н	1156	50	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
о	482	7751	19541	7972	9449	3308	3014	2734	1023	5489	4022	10534	1034	10651	2762	6265	10077	14919	12785	551	751	855	581	4043	1223	370	1	1	942	280	1180	106		
п	2471	11	35	41	9	3004	21	59	2417	0	127	1236	16	22	14924	351	12956	88	90	2620	21	1	180	75	14	0	0	0	466	50	14	3	199	4
р	15086	269	788	834	702	6076	856	156	7222	2	583	101	857	1806	15288	4277	360	8362	1259	4362	144	107	62	78	304	36	1	0	2347	771	50	175	192	142
с	10927	588	5352	279	841	10133	114	444	9582	2	1813	582	698	9041	22322	3171	6240	8362	7275	6281	117	105	85	640	56	15	0	0	2286	7995	312	347	862	79
т	452	1595	1516	1490	3709	673	1822	580	777	226	1430	1249	2217	1813	498	2334	1867	2857	4076	191	113	321	121	2367	633	489	0	0	159	1806	190	3	0	
у	159	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
ф	950	1584	843	182	377	173	90	211	786	0	511	273	355	603	3055	838	501	740	303	340	60	36	155	113	74	0	0	0	106	8	40	3	0	
х	494	71	61	8	18	1376	2	13	3403	0	108	15	121	33	165	74	21	126	32	136	41	4	2	4	3	0	0	264	21	0	7	0	0	
ц	3370	4	49	13	10	3852	4	11	2662	0	438	119	10	1078	68	30	81	49	4817	54	1	3	0	37	262	0	1	2	28	3	1	2	125	
ч	1005	9	76	10	0	2768	0	4	2164	0	390	551	47	618	453	44	13	27	634	429	51	1	1	1	340	10	4	1	1	340	10	4	1	
ш	637	0	1	0	1	2125	0	0	1261	0	1	0	1	0	96	1	1	19	4	1	1	83	0	0	2	1	0	0	0	119	1	0	1	377
щ	0	0	2	0	0	140	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
ъ	17	97	262	34	51	4032	24	373	77	2515	515	1438	2482	2411	1158	1454	47	1844	1298	247	51	246	428	1	0	0	0	0	2157	143	216	1	71	
ы	276	507	1516	452	498	783	106	736	1231	4	1974	267	998	5841	739	1440	146	2629	923	288	87	139	87	517	1082	0	0	0	330	671	679	68	0	
ь	0	1	16	11	11	1	1	2	20	359	189	27	27	3	93	247	25	4573	4	319	47	0	0	39	0	1	0	0	1	1	0	0	0	
э	121	769	181	1879	1	1	1	2	20	359	189	27	27	3	93	247	25	4573	4	319	47	0	0	39	0	1	0	0	1	1	0	0	0	
я	647	646	2470	365	1023	735	225	870	987	105	903	137	1083	2565	888	1965	75	1876	1035	466	111	482	376	765	297	355	0	0	282	373	213	0	0	
е	6	18	102	22	63	18	25	50	78	3	57	119	231	200	89	109	154	91	454	24	0	22	1	42	16	0	0	0	28	8	0	0	0	

Рис. 3: к-ть біграм

[illegible]

Рис. 4: ймовірність появи біграм

4 Обчислення ентропій H_1 та H_2

- Обчислення H_1 :

```
1 fn compute_h1(letter_frequencies: &HashMap<char, i64>) -> f64 {
2     let mut h1 = 0.0;
3     let probabilities = count_letters_probabilities(&
4         letter_frequencies);
5
6     for (_key, _value) in probabilities {
7         h1 += _value * f64::log2(_value);
8     }
9
10    h1 = -h1;
11    h1
12 }
```

- Обчислення H_2 :

```
1 fn compute_h2(bigram_frequencies: &HashMap<String, i64>) -> f64 {
2     let mut h2 = 0.0;
3     let probabilities = count_bigram_probabilities(&
4         bigram_frequencies);
5     for (_key, _value) in probabilities {
6         h2 += _value * f64::log2(_value);
7     }
8     h2 = -h2/2.0;
9
10    h2
11 }
```

Якщо у тексті наявні пробіли, то $H_1 = 4.404$, якщо ж їх немає, то $H_1 = 4.461$.
Так само й для H_2 , з пробілами $H_2 = 4.021$, без $H_2 = 4.152$.

5 Оціночні значення величин $H^{(10)}$, $H^{(20)}$ та $H^{(30)}$

Використовуючи програму CoolPinkProgram.exe визначаємо приблизні значення:

Лабораторная работа №1

Произвольная часть текста:
чество_в_

Использованные буквы:

Порядок n-граммы:
5 символов
10 символов
15 символов
20 символов
25 символов
30 символов
35 символов
40 символов
45 символов
50 символов

Введенный символ:
Символ по счету:
Номер эксперимента: 51

Неравенство для энтропии:
 $2.94983312141789 < H < 3.45521058262596$

Двоичная таблица угаданных символов:

00010000000000000000000000000000
10000000000000000000000000000000
00000000000000000001000000000000
00000000010000000000000000000000
10000000000000000000000000000000

Вероятности:

$q[1] = 0.3$
$q[2] = 0.1$
$q[3] = 0.08$
$q[4] = 0.1$
$q[5] = 0$
$q[6] = 0.04$
$q[7] = 0$
$q[8] = 0$
$q[9] = 0.02$
$q[10] = 0.04$
$q[11] = 0$
$q[12] = 0.04$
$q[13] = 0.02$
$q[14] = 0.02$
$q[15] = 0.02$
$q[16] = 0.1$
$q[17] = 0$
$q[18] = 0$
$q[19] = 0.02$
$q[20] = 0.02$
$q[21] = 0$
$q[22] = 0$
$q[23] = 0$
$q[24] = 0$
$q[25] = 0.04$
$q[26] = 0$
$q[27] = 0.02$
$q[28] = 0$
$q[29] = 0$
$q[30] = 0$
$q[31] = 0.02$
$q[32] = 0$

Строка состояния:
Вы угадали. Для продолжения опыта нажмите "Продолжить", или "Другой" для выбора другого порядка

Рис. 5: значения $H^{(10)}$

Лабораторная работа №1

Произвольная часть текста:
потому_что_в_разные_века_различные_цивилизации_придерживались_совершенно_не

Использованные буквы:

Порядок n-граммы:
5 символов
10 символов
15 символов
20 символов
25 символов
30 символов
35 символов
40 символов
45 символов
50 символов

Введенный символ: _ (пробел)
Символ по счету: 1
Номер эксперимента: 50

Неравенство для энтропии:
 $2.40463996775542 < H < 2.97176867487454$

Двоичная таблица угаданных символов:

01000000000000000000000000000000
10000000000000000000000000000000
01000000000000000000000000000000
10000000000000000000000000000000
00010000000000000000000000000000

Вероятности:

$q[1] = 0.38$
$q[2] = 0.18$
$q[3] = 0.1$
$q[4] = 0.04$
$q[5] = 0$
$q[6] = 0.04$
$q[7] = 0$
$q[8] = 0.02$
$q[9] = 0.02$
$q[10] = 0$
$q[11] = 0.02$
$q[12] = 0$
$q[13] = 0$
$q[14] = 0$
$q[15] = 0$
$q[16] = 0$
$q[17] = 0$
$q[18] = 0$
$q[19] = 0.04$
$q[20] = 0$
$q[21] = 0$
$q[22] = 0.02$
$q[23] = 0$
$q[24] = 0.02$
$q[25] = 0$
$q[26] = 0$
$q[27] = 0.06$
$q[28] = 0.04$
$q[29] = 0$
$q[30] = 0.02$
$q[31] = 0$
$q[32] = 0$

Строка состояния:
Вы угадали. Для продолжения опыта нажмите "Продолжить", или "Другой" для выбора другого порядка

Рис. 6: значения $H^{(20)}$

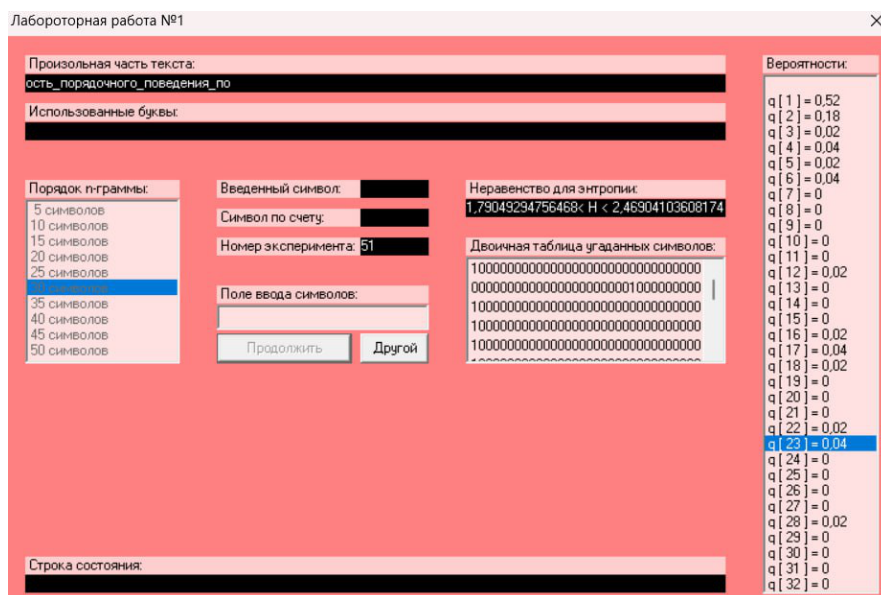


Рис. 7: значения $H^{(30)}$

Маємо, що:

$$\begin{aligned} a &\leq H^{(10)} \leq c \\ a &\leq H^{(20)} \leq c \\ a &\leq H^{(30)} \leq c \end{aligned}$$

Тоді можемо обчислити надлишковість r російської мови R використовуючи формулу:

$$R = 1 - \frac{H_{\inf}}{H_0}, \quad \text{де } H_0 = \log 32 \approx 5$$

$$\begin{aligned} H^{(30)} : 0.642 &\geq R \geq 0.506 \\ R &\sim 0.57 \end{aligned}$$

6 Висновок

З результатів обрахунків бачимо, що r російської мова є досить надлишковою, тому ми можемо ущільнювати тексти цієї мови майже вдвічі без втрати змісту.