

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:
ФБ-31 Аль-Фітурі Асія
ФБ-31 Гриб Вероніка

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 а) частот букв і б) частот біграм в тексті, а також підрахунку с) 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.

А) частота букв

1. Частота літер

Частота окремої літери обчислюється як відношення кількості появ цієї літери до загальної кількості літер у тексті

$$P(l) = \frac{f(l)}{n}$$

Дані збереглись а таблиці

Створено файл 'crypto_5tables.xlsx' з 7 аркушами.

Частота літери з пробілами

	A	B	C
1	Літера	Кількість	Частота
2	пробіл	173144	0.16767
3	о	95774	0.09275
4	а	74368	0.07202
5	е	73361	0.07104
6	е	73361	0.07104
7	и	53873	0.05217
8	н	52690	0.05102
9	т	51049	0.04944
10	с	44734	0.04332
11	л	43373	0.042
12	р	43076	0.04171
13	в	41001	0.0397
14	к	31695	0.03069
15	м	27873	0.02699
16	д	25559	0.02475
17	у	24698	0.02392
18	п	21472	0.02079
19	г	18252	0.01768
20	ь	17094	0.01655
21	ь	17094	0.01655
22	я	16841	0.01631
23	ы	14150	0.0137
24	з	13902	0.01346
25	б	13350	0.01293
26	ч	11986	0.01161
27	й	10369	0.01004
28	ж	8371	0.00811
29	ш	8121	0.00786
30	х	7379	0.00715
31	ю	5223	0.00506
32	ц	3753	0.00363
33	э	2827	0.00274
34	щ	2233	0.00216
35	ф	1050	0.00102

2020

Частота літерів без пробілів

	A	B	C
1	Літера	Кількість	Частота
2	о	95774	0.11143
3	а	74368	0.08653
4	е	73361	0.08535
5	и	53873	0.06268
6	н	52690	0.0613
7	т	51049	0.05939
8	с	44734	0.05205
9	л	43373	0.05046
10	р	43076	0.05012
11	в	41001	0.0477
12	к	31695	0.03688
13	м	27873	0.03243
14	д	25559	0.02974
15	у	24698	0.02874
16	п	21472	0.02498
17	г	18252	0.02124
18	ь	17094	0.01989
19	я	16841	0.01959
20	ы	14150	0.01646
21	з	13902	0.01617
22	б	13350	0.01553
23	ч	11986	0.01395
24	й	10369	0.01206
25	ж	8371	0.00974
26	ш	8121	0.00945
27	х	7379	0.00859
28	ю	5223	0.00608
29	ц	3753	0.00437
30	э	2827	0.00329
31	щ	2233	0.0026
32	ф	1050	0.00122
33			

Без пробілу без перетину

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH
1		а	б	в	г	д	е	є	ж	з	и	й	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	ш	щ	ь	ы	ь	э	ю	я
2		0.000428	0.001364	0.005997	0.001983	0.00333	0.002131	0.002131	0.001459	0.004814	0.001156	0.00088	0.006611	0.009019	0.006392	0.006529	0.001401	0.002571	0.006346	0.007397	0.006797	0.000645	0.000279	0.001112	0.000072	0.001478	0.001964	0.000358	0	0	0.000435	0.001119	0.002464	
3		0.001077	0.000002	0.000135	0	0.00003	0.002276	0.002276	0.000007	0.000002	0.00107	0	0.000133	0.00097	0.000079	0.000328	0.002376	0.000014	0.001594	0.000156	0.000007	0.001426	0.000002	0.000061	0.000005	0	0.000014	0.000086	0.000128	0.00302	0.000128	0.000079	0.000007	0.00057
4		0.0081	0.00172	0.000389	0.000344	0.000068	0.00733	0.00733	0.000035	0.000814	0.001449	0	0.000763	0.001056	0.000365	0.001273	0.008447	0.000914	0.001317	0.003525	0.000842	0.001008	0.000009	0.000137	0.000028	0.000289	0.001233	0.000007	0.000289	0.003828	0.000289	0.000305	0.000019	0.000321
5		0.002669	0.000049	0.000126	0.000084	0.001287	0.002716	0.002716	0.000016	0.000037	0.001035	0	0.000144	0.001126	0.000044	0.00037	0.00914	0.000109	0.000845	0.000168	0.000009	0.001047	0.000005	0.000002	0.000005	0.000033	0.000002	0.000002	0	0	0	0	0.00012	
6		0.005682	0.000074	0.001878	0.000074	0.000074	0.00554	0.00554	0.000016	0.000079	0.002646	0	0.000219	0.000717	0.000128	0.002355	0.004042	0.000207	0.001357	0.000449	0.000024	0.001757	0.000005	0.000004	0.000249	0.000091	0.000056	0.000012	0.000831	0.000484	0.000831	0.000019	0.000026	0.000323
7		0.000335	0.002024	0.004682	0.004251	0.003916	0.001852	0.001852	0.001152	0.001768	0.001017	0.002183	0.002855	0.007283	0.00622	0.010532	0.00155	0.003044	0.00884	0.00833	0.006916	0.0008	0.000081	0.000787	0.00054	0.001617	0.000896	0.000589	0	0	0	0.000261	0.000165	0.000437
8		0.000335	0.002024	0.004682	0.004251	0.003916	0.001852	0.001852	0.001152	0.001768	0.001017	0.002183	0.002855	0.007283	0.00622	0.010532	0.00155	0.003044	0.00884	0.00833	0.006916	0.0008	0.000081	0.000787	0.00054	0.001617	0.000896	0.000589	0	0	0	0.000261	0.000165	0.000437
9		0.001501	0.000081	0.000033	0.000033	0.000033	0.000016	0.001145	0.003933	0.003933	0.000019	0.000007	0.000007	0.001443	0	0.000102	0.000033	0.000019	0.000807	0.000058	0.000023	0.000007	0.000021	0.000014	0.000396	0	0	0.000051	0	0	0.000037	0.000005	0	0.000014
10		0.00619	0.000207	0.001022	0.000538	0.001019	0.000298	0.000298	0.000077	0.000105	0.000365	0	0.000285	0.001247	0.000396	0.002169	0.000789	0.00023	0.000412	0.000226	0.000095	0.000496	0.000009	0	0.000007	0.000004	0.000019	0	0.000209	0.000363	0.000209	0.000337	0	0.000361
11		0.000296	0.001359	0.004284	0.001136	0.002388	0.002841	0.002841	0.00057	0.002286	0.001415	0.001841	0.00333	0.006478	0.003567	0.005499	0.001503	0.001971	0.001419	0.004086	0.006192	0.000442	0.000095	0.000293	0.001031	0.00199	0.00071	0.000179	0	0	0	0.00024	0.000286	0.001485
12		0.000158	0.000037	0.000856	0.000291	0.000798	0.000168	0.000168	0.000182	0.00023	0.000477	0	0.000877	0.00023	0.000517	0.001056	0.000672	0.000977	0.000309	0.001515	0.001073	0.000219	0.000077	0.000074	0.000163	0.000484	0.000163	0.000007	0	0	0	0.000091	0.000014	0.0001
13		0.009603	0.000291	0.000686	0.000182	0.000258	0.000691	0.000691	0.000133	0.000109	0.003106	0	0.000454	0.001087	0.000237	0.001282	0.01297	0.000375	0.002152	0.000586	0.000775	0.001559	0.000028	0.000056	0.000012	0.000337	0.000056	0.000002	0.000002	0	0.000002	0.000114	0.000049	0.000109
14		0.009126	0.000333	0.00081	0.0008	0.000363	0.005117	0.005117	0.0004	0.000226	0.0078	0	0.001485	0.000461	0.000333	0.01042	0.006132	0.000264	0.002064	0.00034	0.001978	0.000049	0.000063	0.000012	0.000377	0.000067	0.000012	0.005373	0.000752	0.005373	0.00004	0.001519	0.002234	
15		0.03849	0.000335	0.000912	0.000389	0.000517	0.004468	0.004468	0.000121	0.000138	0.003216	0	0.000828	0.000351	0.000523	0.002685	0.006132	0.001083	0.001298	0.001022	0.000377	0.00393	0.000054	0.000065	0.000333	0.000361	0.000037	0.000061	0.000054	0.000949	0.000054	0.000216	0.000002	0.000742
16		0.012472	0.000235	0.000424	0.000147	0.000353	0.010152	0.010152	0.00033	0.000184	0.007602	0	0.000421	0.000088	0.0002	0.003507	0.009715	0.000465	0.001475	0.001661	0.000777	0.003451	0.000026	0.00003	0.000593	0.000358	0.00003	0.000198	0.000999	0.00018	0.000998	0.000042	0.000182	0.002145
17		0	0.00023	0.00471	0.009641	0.006455	0.006504	0.003207	0.003207	0.00045	0.001983	0.004905	0.004198	0.009582	0.006485	0.008545	0.001882	0.003383	0.008472	0.008752	0.008989	0.000668	0.000116	0.000898	0.000163	0.002853	0.001664	0.000216	0	0	0	0.000496	0.000458	0.001368
18		0.001587	0	0.00007	0.000005	0.000005	0.002003	0.002003	0.000002	0.000002	0.000919	0	0.000037	0.000712	0.000005	0.00112	0.009278	0.00004	0.008275	0.00003	0.000154	0.000803	0	0	0	0.000007	0	0	0.000154	0.00033	0.000154	0.000002	0	0.000249
19		0.008133	0.000128	0.000738	0.00114	0.000356	0.006543	0.006543	0.000354	0.000079	0.008624	0	0.000424	0.00101	0.000289	0.010129	0.010071	0.000184	0.00084	0.001101	0.000824	0.00299	0.000084	0.000142	0.001096	0.000121	0.000221	0.000035	0.000852	0.001378	0.000852	0.00007	0.000128	0.00128
20		0.001864	0.0002	0.002292	0.000219	0.000338	0.003809	0.003809	0.000144	0.000079	0.002243	0	0.006176	0.00357	0.001096	0.001382	0.003509	0.002783	0.000261	0.000856	0.011476	0.00101	0.000037	0.000195	0.000063	0.000389	0.000093	0	0.000686	0.000386	0.000686	0.000091	0.00017	0.000656
21		0.006602	0.000298	0.004856	0.000188	0.004444	0.007593	0.007593	0.000074	0.000128	0.004712	0	0.000912	0.00377	0.000349	0.001924	0.01428	0.000554	0.003016	0.00175	0.000493	0.001801	0.000023	0.000049	0.000126	0.000316	0.000337	0.000023	0.006427	0.001475	0.006427	0.00142	0.000093	0.000507
22		0.000577	0.000926	0.002087	0.00181	0.002464	0.00043	0.00043	0.00154	0.00077	0.001066	0.00024	0.001631	0.002269	0.001589	0.00093	0.00043	0.001347	0.000538	0.001901	0.001773	0.000163	0.000051	0.000551	0.000009	0.00134	0.000805	0.00023	0	0	0.000116	0.001098	0.000235	
23		0.000147	0.000002	0.000007	0.000002	0.000037	0.00007	0.00007	0.000002	0	0.000156	0	0.000016	0.000116	0.000007	0.000005	0.000135	0.000009	0.000375	0.000009	0.000007	0.00004	0.000012	0	0	0	0.000007	0.000012	0.000007	0.000002	0.000033	0.000005		
24		0.000789	0.000151	0.000577	0.000116	0.000312	0.000119	0.000119	0.00004	0.000088	0.000379	0	0.000263	0.000175	0.000202	0.000061	0.002736	0.000433	0.000361	0.000442	0.000135	0.000265	0.000035	0.000009	0.000007	0.000095	0.000019	0.000007	0.000002	0	0.000002	0.000042	0.000002	0.000061
25		0.000873	0.000014	0.000116	0.000012	0.000016	0.000908	0.000908	0.000002	0.000012	0.000254	0	0.000007	0.000005	0.000021	0.000033	0.001287	0.000063	0.000009	0.000054	0.000023	0.000326	0.000009	0.000005	0.000002	0.000012	0	0	0.000235	0	0.000009	0	0.000002	
26		0.002099	0.000005	0.000019	0	0.000009	0.004223	0.004223	0	0	0.001487	0	0.000226	0.000042	0	0.000798	0.000081	0.000016	0.000019	0.000012	0.003725	0.000721	0	0	0	0.000012	0.000151	0	0.00034	0.000012	0	0	0	
27		0.001124	0.000028	0.000037	0.000012	0.000026	0.00379	0.00379	0.000002	0.000007	0.001429	0	0.000377	0.000561	0.000019	0.000377	0.000421	0.000282	0.000009	0.000016	0.000012	0.000356	0.000002	0	0	0.000005	0.000002	0	0.000507	0.000007	0.000002	0	0	
28		0.000414	0.000002	0.000005	0	0.000002	0.00115	0.00115	0	0.000002	0.000821	0	0.000012	0	0.000005	0.000051	0.000028	0.000009	0.000005	0.000007	0.000002	0.000093	0	0	0	0.000002	0	0	0.000051	0	0	0	0	
29		0.0002	0.000461	0.001613	0.000349	0.0006	0.00134	0.00134	0.000114	0.000572	0.000935	0	0.002036	0.000242	0.000775	0.002355	0.000882	0.001163	0.000216	0.002192	0.000663	0.000347	0.000047	0.000102	0.000086	0.000524	0.000551	0.000088	0	0	0.000179	0.000603	0.000593	
30		0.000079	0.000463	0.001412	0.000226	0.000351	0.000924	0.000924	0.000098	0.000305	0.000356	0.001878	0.000356	0.001657	0.001478	0.000933	0.000461	0.000719	0.000347	0.001287	0.000994	0.000233	0.000028	0.001075	0.000021	0.000268	0.000582	0.000005	0	0	0.000054	0	0.000109	
31		0.0002	0.000461	0.001613	0.000349																													

Надлишковість джерела відкритого тексту (мови) дорівнює $R = 1 - \frac{H_{\infty}}{H_0}$ і

характеризує величину можливого ущільнення тексту деякою схемою кодування символів без втрати його змісту.

H_{∞} — це ентропія «нескінченного порядку» (умовна ентропія, яку оцінюють через експеримент з n-грамами, наприклад $H(10), H(20), H(30)$)

$H_0 = \log_2 m$, де m — кількість символів в алфавіті (для російської мови без пробілу ≈ 32 , з пробілом ≈ 33).

$H(10)$

Ентропія для H_{10} : $2.028 < H < 2.894$

Надлишковість (R) $0.426 < R < 0.598$

Лабораторная работа №1

Произвольная часть текста:
онцов_я_сам_такой_же_то_есть_мне_самому_не_удается_как_следует_соблюдать_ес

Использованные буквы:
к, т, о, г, л, ы, и, _, ч, н,

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

Введенный символ: а

Символ по счету: 11

Номер эксперимента: 4

Неравенство для энтропии:
 $2.68365034650669 < H < 2$

Двоичная таблица угаданных символов:
01000000000000000000000000000000
00000100000000000000000000000000
10000000000000000000000000000000
00000000001000000000000000000000

Вероятности:
q[1] = 0,25
q[2] = 0,25
q[3] = 0
q[4] = 0
q[5] = 0
q[6] = 0,25
q[7] = 0
q[8] = 0
q[9] = 0
q[10] = 0
q[11] = 0,25
q[12] = 0
q[13] = 0
q[14] = 0
q[15] = 0
q[16] = 0
q[17] = 0
q[18] = 0
q[19] = 0
q[20] = 0
q[21] = 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
q[30] = 0
q[31] = 0
q[32] = 0

Поле ввода символов:
а

Продолжить Другой

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

тут ми багато разів вгадували що може бути після с...

Произвольная часть текста:
амое_вам_то_мое_место_я_его_первый_занял_оставьте_его_в_покое_он_не_делает_

Использованные буквы:
н, д, п, в, и, с, б, о, у, ж, р,

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: т

Символ по счету: 12

Номер эксперимента: 19

Неравенство для энтропии:
 $1,74349081867847 < H < 2,39428114771339$

Двоичная таблица угаданных символов:

01000000000000000000000000000000	^
00000100000000000000000000000000	
10000000000000000000000000000000	
00000000001000000000000000000000	
10000000000000000000000000000000	v

Поле ввода символов:
т

Продолжить Другой

Найважче було вгадати перше слово, бо там було «амое вам_» і треба було вгадати, що буде після пробілу

Лабораторная работа №1

×

Произвольная часть текста:
ите_как_этот_же_человек_сам_возвращается_к_отвергнутым_им_принципам_он_может_

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: т

Символ по счету: 1

Номер эксперимента: 52

Неравенство для энтропии:
 $2,02854294639459 < H < 2,89435841102112$

Двоичная таблица угаданных символов:

01000000000000000000000000000000	^
00000100000000000000000000000000	
10000000000000000000000000000000	
00000000001000000000000000000000	
10000000000000000000000000000000	v

Поле ввода символов:
т

Продолжить Другой

Вероятности:

q[1] = 0,4423076
q[2] = 0,1153846
q[3] = 0,0961538
q[4] = 0,0576923
q[5] = 0,0384615
q[6] = 0,0384615
q[7] = 0,0192307
q[8] = 0
q[9] = 0,0192307
q[10] = 0,038461
q[11] = 0,038461
q[12] = 0,038461
q[13] = 0
q[14] = 0
q[15] = 0
q[16] = 0,019230
q[17] = 0,019230
q[18] = 0
q[19] = 0
q[20] = 0
q[21] = 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
q[30] = 0
q[31] = 0
q[32] = 0,019230

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

Надлишковість (R) $0.538 < R < 0.712$

Лабораторная работа №1

Произвольная часть текста:
ое_утверждение_заявив_что_договор_который_они_собираются_нарушить_несправед

Использованные буквы:
л,

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: и Символ по счету: 2 Номер эксперимента: 51 Поле ввода символов: <input type="text"/> и <div>Продолжить Другой</div>	Неравенство для энтропии: $1,45262844902627 < H < 2,33073909065524$ Двоичная таблица угаданных символов: <div> 10000000000000000000000000000000 ^ 10000000000000000000000000000000 10000000000000000000000000000000 10000000000000000000000000000000 10000000000000000000000000000000 v </div>
--	---	---

Вероятности:

```
q[1] = 0,5882352
q[2] = 0,0980392
q[3] = 0,0980392
q[4] = 0,0196078
q[5] = 0,0196078
q[6] = 0,0196078
q[7] = 0
q[8] = 0,0196078
q[9] = 0,0196078
q[10] = 0,019607
q[11] = 0,019607
q[12] = 0,019607
q[13] = 0,019607
q[14] = 0
q[15] = 0,019607
q[16] = 0
q[17] = 0
q[18] = 0
q[19] = 0
q[20] = 0
q[21] = 0,019607
q[22] = 0
q[23] = 0
q[24] = 0
q[25] = 0
q[26] = 0
q[27] = 0
q[28] = 0
q[29] = 0
q[30] = 0
q[31] = 0
q[32] = 0
```

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

$H(30)$

Ентропія $1.484 < H < 2.159$

Надлишковість (R) $0.572 < R < 0.706$

Порядок n-грам	Ентропія (H)	Надлишковість (R)
10 символів	$2.028 < H < 2.894$	$0.426 < R < 0.598$
20 символів	$1.453 < H < 2.331$	$0.538 < R < 0.712$
30 символів	$1.484 < H < 2.159$	$0.572 < R < 0.706$

Висновок:

У роботі на основі книги російською мовою підраховано частоти літер і біграм для двох варіантів (з пробілом/без пробілу) та побудовано матриці біграм з перетином і без, що загалом дало 7 таблиць результатів.

За програмою CoolPinkProgram було досліджено ентропію та надлишковість відкритого тексту за допомогою n-грамного аналізу. Метою було визначити міру невпорядкованості тексту (ентропію H) та частку повторюваної інформації (надлишковість R).

Отримані результати:

$$H10: 2.028 < H < 2.894$$

$$R10: 0.426 < R < 0.598$$

$$H20: 1.452 < H < 2.330$$

$$R20: 0.538 < R < 0.712$$

$$H30: 1.483 < H < 2.159$$

$$R30: 0.572 < R < 0.706$$

Аналіз показав, що зі збільшенням порядку n-грам ентропія зменшується, а надлишковість зростає.