

**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

Виконали:
Кравченко Владислав
Жмур Назар
Група: ФБ-04

Київ - 2022

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 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 видалено всі пробіли. 2. За допомогою програми CoolPinkProgram оцінити значення $H(10)$, $H(20)$, $H(30)$. 3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Хід роботи

Під час виконання роботи для читання файлу використовували функцію `open()`

Всі результати виконання роботи були виписані в окремі файли. Далі наведені таблиці з даними підписані та знаходяться нижче.

Частота монограм з та без пробілів

texts > ≡ monogram_space.txt	texts > ≡ monogram_no_space.txt
1 : 0.16464678482139297	1 о : 0.10734058657145787
2 о : 0.08966730411162493	2 е : 0.08725345925984002
3 е : 0.07288745772816296	3 а : 0.08523007118546619
4 а : 0.07119721399468072	4 и : 0.06683563414570405
5 и : 0.05583136187211497	5 н : 0.06597370902211956
6 н : 0.05511134994888545	6 т : 0.06446446054041449
7 т : 0.053850594377189684	7 с : 0.05295838836587038
8 с : 0.044238959992107155	8 л : 0.052550293205316095
9 л : 0.04389805638763931	9 в : 0.04124625829006118
10 в : 0.03445519447668988	10 р : 0.04094471014186835
11 р : 0.03420329526156585	11 к : 0.03480066662243961
12 к : 0.02907084875341376	12 м : 0.03134241194291493
13 м : 0.026181984587966355	13 д : 0.031145903066342604
14 д : 0.026017830266110528	14 у : 0.028407845880751778
15 у : 0.02373058539278435	15 п : 0.026100499966829702
16 п : 0.02180313656506032	16 ь : 0.02035198710178054
17 ь : 0.017001097860745915	17 я : 0.019851417175780456
18 я : 0.016582945163640027	18 ы : 0.01916740545962974
19 ы : 0.01601155377733369	19 б : 0.01752698353346079
20 б : 0.014641222047058972	20 г : 0.017433503607521012
21 г : 0.014563133290370523	21 з : 0.015979036372737634
22 з : 0.01334813940942229	22 ч : 0.01467232773056874
23 ч : 0.012256576143884831	23 ж : 0.011706601693092336
24 ж : 0.009779147363140008	24 й : 0.010854728174447614
25 й : 0.009067532080414626	25 х : 0.009357541618670253
26 х : 0.007816852477323822	26 ш : 0.008913763260579816
27 ш : 0.007446140799066293	27 ю : 0.005729917395910605
28 ю : 0.004786504919381755	28 э : 0.004586044753765834
29 э : 0.003830967230011272	29 щ : 0.003143136864663181
30 щ : 0.002625629485642794	30 ц : 0.0030571956424282265
31 ц : 0.002553838209332446	31 ф : 0.0010735114075664462
32 ф : 0.000896761205841542	32
33	

Частота неперехресних біграм з та без пробілів

texts > ≡ bigram_space.txt	texts > ≡ bigram_no_space.txt
1 а : 0.02057681585823778	1 то : 0.015251300175098959
2 и : 0.02044246955376875	2 на : 0.013123375409351611
3 о : 0.0196825732691158	3 не : 0.012645924174712976
4 е : 0.01926777905406767	4 ст : 0.011793548075821267
5 с : 0.01629704639649624	5 он : 0.011453803828857353
6 н : 0.01626849780679658	6 ла : 0.011396509680700717
7 п : 0.01517189609656812	7 по : 0.010862769458399422
8 в : 0.01499808556516131	8 но : 0.010462715581796945
9 то : 0.01199292666706970	9 ен : 0.01013402810026677
10 о : 0.01130020353465127	10 ка : 0.010118950692857129
11 ь : 0.01104074723414545	11 ал : 0.009715881334772723
12 на : 0.01094418582780834	12 ли : 0.009683716198965488
13 не : 0.01070824013058460	13 ни : 0.009305775853230485
14 и : 0.01002559297100135	14 ов : 0.009259538470507585
15 ст : 0.00981315787705969	15 ко : 0.009231393976676256
16 я : 0.00972835177236362	16 ос : 0.008985129655652117
17 ла : 0.00943111057372589	17 ет : 0.008931856149471387
18 по : 0.00911035877180608	18 от : 0.00827548634690501
19 но : 0.00873167012608400	19 ра : 0.008142805161700168
20 к : 0.00851587637453062	20 го : 0.00799404140859171
21 ка : 0.00831519658223001	21 ть : 0.007736720322133835
22 д : 0.00790963867561413	22 ас : 0.0077236532357121464
23 ал : 0.00788780740113791	23 ер : 0.007604039136928993
24 ли : 0.00779460465241252	24 ло : 0.0074693476307362
25 ни : 0.00765438069712297	25 во : 0.007416074124555467
26 он : 0.00752423271466860	26 ро : 0.0073818986677602816
27 ко : 0.00737393278654387	27 та : 0.0072833929393506264
28 м : 0.00730172164789177	28 ес : 0.007222078149218086
29 т : 0.00728240936662434	29 ол : 0.006978829309675876
30 т : 0.00715897869939342	30 ан : 0.006967772544242139
31 м : 0.00711111782842633	31 ле : 0.00695470545782045
32 й : 0.00701539608649215	32 пр : 0.006930581605965025
33 б : 0.00690288105649933	33 ат : 0.0069205300010252635
34 ра : 0.00679120569090945	34 ом : 0.006894395828181886
35 ет : 0.00669128562696061	35 ак : 0.006804936544218015
36 го : 0.00667785099651371	36 ел : 0.006728544346675833
37 ен : 0.00663334878315834	37 ва : 0.006395836223169754
38 ть : 0.00623702718497470	38 ре : 0.006250087951543223
39 у : 0.00613962611423466	39 од : 0.006219933136723941
40 ро : 0.00613122947020534	40 ор : 0.006208876371290204
41 от : 0.00605817866715031	41 ит : 0.006119417087326333

Частота перехресних біграм з та без пробілів

texts > ≡ cross_bigram_space.txt	texts > ≡ cross_bigram_no_space.txt
1 а : 0.02070780350509508	1 то : 0.015110082719724425
2 и : 0.02042567626571012	2 на : 0.013064077505954322
3 о : 0.01967039813527329	3 не : 0.012698701482763248
4 е : 0.01941094183476748	4 ст : 0.011924727513417643
5 с : 0.01640956142648906	5 ла : 0.011511103762267362
6 н : 0.01640284411126561	6 он : 0.011327661879922462
7 п : 0.01520842149809564	7 по : 0.010893934908789179
8 в : 0.01498549059911734	8 но : 0.010561729198570059
9 то : 0.01212727297153873	9 ен : 0.010111417070950797
10 о : 0.01124730467726659	10 ка : 0.010043066123008587
11 ь : 0.01088121099758848	11 ал : 0.009704326866294992
12 на : 0.01087827217217822	12 ли : 0.009615370117870205
13 не : 0.01054030724999832	13 ни : 0.009379659863569497
14 и : 0.01002265414559109	14 ко : 0.009359556643586495
15 ст : 0.00969308586744050	15 ов : 0.009105753491301086
16 я : 0.00958980714587993	16 ос : 0.009045946411851653
17 ла : 0.00957175436121690	17 ет : 0.008935378701945136
18 по : 0.00909986296676944	18 от : 0.008346354356443158
19 но : 0.00865400116881284	19 ра : 0.008198595689568088
20 к : 0.00840336134453781	20 го : 0.008087525399161997
21 ка : 0.00834542450073554	21 ть : 0.007655306169527439
22 д : 0.00785967864363971	22 ас : 0.007620628115056759
23 ли : 0.00782567223532098	23 ер : 0.007607561022067807
24 ал : 0.00778620800838320	24 во : 0.007466838482186788
25 ни : 0.00759266536350751	25 ро : 0.007416580432229281
26 он : 0.00755865895518879	26 ло : 0.00740451850023948
27 ко : 0.00751373690963195	27 та : 0.00733013658630237
28 т : 0.00732187359356212	28 ес : 0.007209517266404353
29 м : 0.00725218144811881	29 ан : 0.007033614091553079
30 т : 0.00713672759271574	30 ом : 0.006951693470122343
31 й : 0.00703764719316983	31 ол : 0.006950688309123193
32 м : 0.00702967038134198	32 ле : 0.0069114870301563375
33 б : 0.00687727129220993	33 ат : 0.006909979288657613
34 ра : 0.00683654756866775	34 пр : 0.00685519801420393
35 ен : 0.00672277304207054	35 ак : 0.00677780061726937
36 ет : 0.00671605572684709	36 ел : 0.006701408381333959
37 го : 0.00671059790822803	37 ва : 0.006461677483036651
38 ть : 0.00639488409272581	38 ре : 0.0062787381811913266
39 ро : 0.00613626745662293	39 од : 0.0062350136777282955
40 у : 0.00607959010942506	40 ит : 0.006142036285306908
41 от : 0.00602921024524917	41 ор : 0.006137010480311158

Ентропія та надлишковість

Монограми:

Ентропія з пробілами: 4.3673644204588555

Надлишковість з пробілами: 0.12652711590822885

Ентропія без пробілів: 4.455664035773609

Надлишковість без пробілів: 0.10886719284527824

Біграми без перетину:

Ентропія з пробілами: 3.9694337769074917

Надлишковість з пробілами: 0.20611324461850167

Ентропія без пробілів: 4.14573870702101

Надлишковість без пробілів: 0.1708522585957979

Біграми з перетином:

Ентропія з пробілами: 3.969475796555487

Надлишковість з пробілами: 0.20610484068889023

Ентропія без пробілів: 4.145451584049106

Надлишковість без пробілів: 0.17090968319017874

coolpinkprogram

Лабораторная работа №1

Произвольная часть текста:
и_нигде_не_считался_похвальным_качеством_разного_мнения_держались_люди_и_по

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: e

Символ по счету: 1

Номер эксперимента: 50

Поле ввода символов:
e

Продолжить Другой

Неравенство для энтропии:
2,61702845507456 < H < 3,24888737212111

Двоичная таблица угаданных символов:
10000000000000000000000000000000
00000000000010000000000000000000
00000000000010000000000000000000
00000001000000000000000000000000
00000000000000000000100000000000
.....

Вероятности:
q[1] = 0,42
q[2] = 0,12
q[3] = 0,02
q[4] = 0,04
q[5] = 0
q[6] = 0
q[7] = 0,04
q[8] = 0,02
q[9] = 0
q[10] = 0
q[11] = 0,02
q[12] = 0,02
q[13] = 0,04
q[14] = 0
q[15] = 0,02
q[16] = 0
q[17] = 0
q[18] = 0,04
q[19] = 0,02
q[20] = 0
q[21] = 0,04
q[22] = 0
q[23] = 0,02
q[24] = 0
q[25] = 0,02
q[26] = 0
q[27] = 0,02
q[28] = 0,04
q[29] = 0,02
q[30] = 0
q[31] = 0
q[32] = 0,02

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

$R = 0.41341$

Произвольная часть текста:

ую и Богу известно то, что я не пытаюсь показаться лучше других, я просто ста

Использованные буквы:

в, л, й, ц, у, м,

Порядок n-граммы:

5 символов

10 символов

15 символов

25 символов

30 символов

35 символов

40 символов

45 символов

50 символов

Введенный символ:

т

Символ по счету:

7

Номер эксперимента:

50

Поле ввода символов:

т

Продолжить

Другой

Неравенство для энтропии:

2,00847251710399 < H < 2,71243118972888

Двоичная таблица угаданных символов:

01000000000000000000000000000000

10000000000000000000000000000000

10000000000000000000000000000000

00010000000000000000000000000000

00100000000000000000000000000000

00000000000000000000000000000000

Вероятности:

q[1] = 0,46

q[2] = 0,18

q[3] = 0,04

q[4] = 0,08

q[5] = 0,02

q[6] = 0,02

q[7] = 0,02

q[8] = 0

q[9] = 0,02

q[10] = 0

q[11] = 0

q[12] = 0

q[13] = 0

q[14] = 0

q[15] = 0

q[16] = 0,02

q[17] = 0

q[18] = 0

q[19] = 0

q[20] = 0,04

q[21] = 0

q[22] = 0

q[23] = 0

q[24] = 0

q[25] = 0

q[26] = 0,04

q[27] = 0,02

q[28] = 0,02

q[29] = 0,02

q[30] = 0

q[31] = 0

q[32] = 0

Строка состояния:

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

$$R = 0.52792$$

Произвольная часть текста:
ие_мы_не_поясняем_внешними_причинами_мы_ставим_его_исключительно_в_заслугу_

Использованные буквы:

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

Введенный символ: и

Символ по счету: 1

Номер эксперимента: 50

Поле ввода символов:
и

Продолжить

Другой

Неравенство для энтропии:
1,93704438485278 < H < 2,53624384458448

Двоичная таблица угаданных символов:
01000000000000000000000000000000
10000000000000000000000000000000
01000000000000000000000000000000
00000010000000000000000000000000
10000000000000000000000000000000
.....

Вероятности:
q[1] = 0,5
q[2] = 0,18
q[3] = 0,06
q[4] = 0
q[5] = 0
q[6] = 0
q[7] = 0,04
q[8] = 0,02
q[9] = 0
q[10] = 0
q[11] = 0
q[12] = 0
q[13] = 0
q[14] = 0,02
q[15] = 0,04
q[16] = 0
q[17] = 0
q[18] = 0,02
q[19] = 0
q[20] = 0
q[21] = 0
q[22] = 0
q[23] = 0,04
q[24] = 0,02
q[25] = 0
q[26] = 0
q[27] = 0,02
q[28] = 0,02
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

$$R = 0.55268$$

Висновок

Під час роботи ми навчилися рахувати ентропію, та надлишковість російської мови на прикладі вибраного тексту. З отриманих даних можемо зробити висновок що найчастіше у нашому тексті зустрічаються букви «о», «е» та «а»