КРИПТОГРАФІЯ

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

Виконали:

Винник Михайло та Кузнєцов Олексій ФБ-12

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:

- 0. Уважно прочитати методичні вказівки до виконання комп'ютерного практикуму.
- 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 H1() 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.
- 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:

- 1. Спочатку в коді програми ми відфільтрували текст та очистили його від зайвих символів, знаків переносу слова, подвійних пробілів, пробілів взагалі та чисел та інших символів, що не належать російському алфавіту. Таким чином отримуємо другу версію нашого текстового документу. Відповідний крок продемонстровано в коді методом *filter()*.
- 2. Далі починається виконання основного завдання:

В наступній частині програми ми виконуємо задачу, яка відповідає підрахунку частот букв, а потім і частот біграм в тексті. Щоб це здійснити, ми спочатку рахуємо кількість, яку трапляється кожна літера з алфавіту в нашому тексті. Також рахуємо кількість пробілів, тому що будемо порівнювати два випадки: коли в нашому алфавіті є пробіл, та коли його нема. Для цієї задачі використовуємо словники. Вони дуже зручні для того, щоб зберігати інформацію про елемент, та його кількість.

Після того, як ми підрахували, скільки в тексті трапляється кожна літера, ми можемо порахувати імовірність, або ж частоту: ділимо кількість кожного окремого символу алфавіту на загальну кількість символів в тексті, що відповідають символам алфавіту.

3. Тепер можна легко обчислити H(1). Для цього скористаємось формулою:

$$H_1 = -\sum_{i=1}^{n} p(i) \log_2 p(i)$$

Де n – кількість літер алфавіту, p(i) – імовірність (частота) появи літери в тексті.

Цю формулу ми теж реалізували в своєму коді.

Далі обчислюємо H(2). Для цього треба порахувати частоту біграм. Робимо це так само, як і з частотою символів, але з урахуванням того, що частота біграм — відношення кількості появ деякої біграми до загальної кількості біграм у тексті. Щоб порахувати біграми, ми розробили фрагмент коду, що їх формує з тексту та рахує в словнику.

4. Н(2) обчислюємо за такою формулою:

$$H_2 = -\sum_{i,j} p(i,j) \log_2 p(i,j)/2$$

де p(i,j) — частота появи деякої біграми в тексті.

Таку формулу ми отримуємо з загальної формули для розрахунку ентропії п-грам:

$$H_n = \frac{1}{n} * H(x_1, x_2, x_3 ..., x_n)$$

Де $H(x_1, x_2, x_3 ..., x_n)$ в свою чергу, — ентропія n-грами відкритого тексту $(x_1, x_2, x_3 ..., x_n)$.

Результати виконання програми:

(на наступній сторінці)

Таблиця частот літер

3 пробілом:

Буква	Кількість	Частота
a	47917	0.06147
б	10629	0.01363
В	29047	0.03726
Γ	11416	0.01464
Д	19428	0.02492
e	56840	0.07291
ë	1	0.0
ж	6235	0.008
3	12203	0.01565
И	49448	0.06343
й	7911	0.01015
К	20658	0.0265
Л	31596	0.04053
M	23876	0.03063
Н	45985	0.05899
О	73763	0.09462
П	18248	0.02341
p	29569	0.03793
c	36491	0.04681
Т	37317	0.04787
y	17100	0.02194
ф	1448	0.00186
X	7790	0.00999
Ц	1981	0.00254
Ч	9675	0.01241
Ш	5103	0.00655
Щ	3072	0.00394
ъ	150	0.00019
ы	14333	0.01839
Ь	11523	0.01478
Э	2296	0.00295
Ю	3876	0.00497
Я	14067	0.01804
_	118574	0.1521

Таблиця частот літер

Без пробіла:

Буква	Кількість	Частота
a	47917	0.07249
б	10629	0.01608
В	29047	0.04394
Γ	11416	0.01727
Д	19428	0.02939
e	56840	0.08599
ë	1	0.0
ж	6235	0.00943
3	12203	0.01846
И	49448	0.07481
й	7911	0.01197
К	20658	0.03125
Л	31596	0.0478
M	23876	0.03612
Н	45985	0.06957
0	73763	0.11159
П	18248	0.02761
p	29569	0.04473
c	36491	0.05521
Т	37317	0.05646
у	17100	0.02587
ф	1448	0.00219
X	7790	0.01179
Ц	1981	0.003
Ч	9675	0.01464
Ш	5103	0.00772
Щ	3072	0.00465
ъ	150	0.00023
Ы	14333	0.02168
Ь	11523	0.01743
Э	2296	0.00347
Ю	3876	0.00586
R	14067	0.02128

Таблиця частот біграм з повторами 3 пробілом:

Біграма	Кількість	Частота
И_	15616	0.02003
0_	14985	0.01922
e_	13857	0.01778
_п	12465	0.01599
_c	11913	0.01528
B	11166	0.01432
_H	11063	0.01419
a	9498	0.01218
ст	9300	0.01193
И	8901	0.01142
Я	8746	0.01122
НО	8703	0.01116
то	8703	0.01116
o	8700	0.01116
не	7405	0.0095
на	7334	0.00941
по	7285	0.00934
ни	6967	0.00894
Ь_	6875	0.00882
 M	6600	0.00847

Таблиця частот біграм з піторами Без пробілу:

		ı
Біграма	Кількість	Частота
ст	9436	0.01428
то	8912	0.01348
НО	8871	0.01342
ен	7677	0.01161
не	7429	0.01124
на	7349	0.01112
ПО	7287	0.01102
НИ	7127	0.01078
oc	6824	0.01032
ОВ	6672	0.01009
ли	6631	0.01003
pa	6469	0.00979
ко	6323	0.00957
ал	6183	0.00935
po	5582	0.00844
ec	5487	0.0083
ГО	5480	0.00829
OM	5474	0.00828
ep	5451	0.00825
ОН	5316	0.00804

Таблиця частот біграм без повторів 3 пробілом:

	ı	ı
Біграма	Кількість	Частота
и_	7736	0.01985
0_	7522	0.0193
e_	6956	0.01785
_п	6283	0.01612
С	5992	0.01537
Н	5575	0.0143
В	5573	0.0143
a	4767	0.01223
СТ	4548	0.01167
и	4478	0.01149
_	4444	0.01145
я_		
ТО	4380	0.01124
_0	4344	0.01114
НО	4279	0.01098
не	3704	0.0095
на	3683	0.00945
по	3592	0.00922
b_	3404	0.00873
ни	3399	0.00872
M_	3313	0.0085
ен	3281	0.00842
pa	3202	0.00821

Таблиця частот біграм без повторів Без пробіла:

Біграма	Кількість	Частота
СТ	4670	0.01413
но	4531	0.01371
то	4443	0.01344
ен	3788	0.01146
не	3742	0.01132
по	3663	0.01108
ни	3654	0.01106
на	3617	0.01094
ос	3389	0.01025
ОВ	3289	0.00995
ли	3266	0.00988
ра	3151	0.00953
ко	3134	0.00948
ал	3097	0.00937
ро	2845	0.00861
ес	2794	0.00845
ОМ	2790	0.00844
ер	2761	0.00835
го	2706	0.00819
от	2674	0.00809
ОН	2662	0.00805
pe	2583	0.00782

Надлишковість відкритого тексту обчислюється за формулою:

$$R=1-\frac{H_{\infty}}{H_0}$$

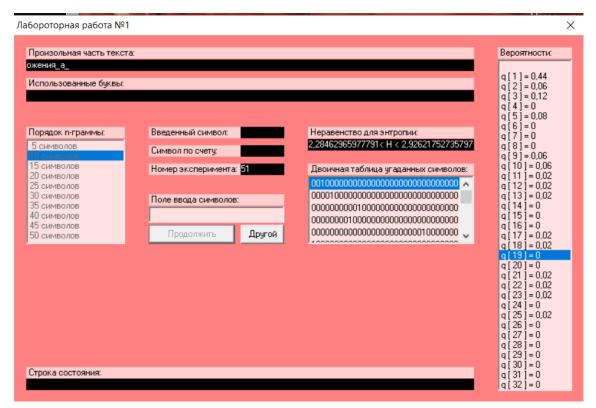
 $H_0 = \log_2 34 = 5.08746$ для тексту з пробілами, $\log_2 33 = 5.04439$ для тексту без пробілів

	Ентропія	Надлишковість
Н1 (з пробілом)	4,4032187	0,134417397
Н1 (без пробіла)	4,4676885	0,11425684
Н2 (з пробілом, з повторами)	4,0015652	0,213374248
Н2 (без пробіла, з повторами)	4,1504491	0,177151249
Н2 (з пробілом, без повторів)	4,0017929	0,213329487
Н2 (без пробіла, без повторів)	4,1496428	0,177311102

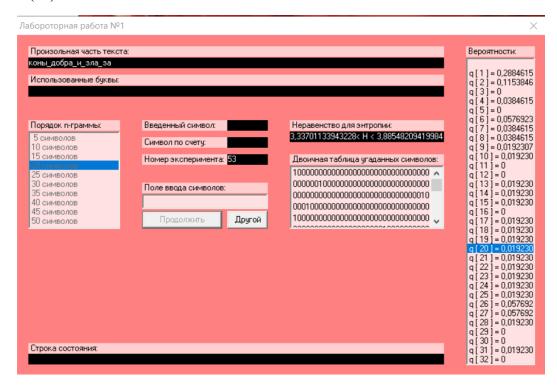
Експеременти в програмі CoolPinkProgram

	Ентропія	Надлишковість
H(10)	2.2846297 <h<2.9262175< th=""><th>0.4147565<r<0,54307406< th=""></r<0,54307406<></th></h<2.9262175<>	0.4147565 <r<0,54307406< th=""></r<0,54307406<>
H(20)	3.3370113 <h<3.8854821< th=""><th>0.2229035<r<0.3325977< th=""></r<0.3325977<></th></h<3.8854821<>	0.2229035 <r<0.3325977< th=""></r<0.3325977<>
H(30)	2.7737683 <h<3.4553927< th=""><th>0.3089214<r<0.4452463< th=""></r<0.4452463<></th></h<3.4553927<>	0.3089214 <r<0.4452463< th=""></r<0.4452463<>

H(10)



H(20)



H(30)



Висновок:

При виконанні цього комп'ютерного практикуму ми навчились якісно очищувати текст для подальшої роботи над ним. Також ми навчились визначати експериментальним шляхом значення імовірностей(частот) літер та біграм(а за аналогіює й інших п-грам). Користуючись цими значиннями, ми тепер можемо вираховувати ентропію та надлишковість у різних моделях відкритого тексту.

Також ми навчилися працювати з програмою CoolPinkProgram, за допомогою якої змогли наближено обчислити значення H(10), H(20) та H(30).