КРИПТОГРАФІЯ

КОМП'ЮТЕРНИЙ ПРАКТИКУМ №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

Виконав

ФБ-12 Сущенко Олександр

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 1 Н та 2 Н за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення 1 Н та 2 Н на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення 1 Н та 2 Н 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

Створюємо програму, яка спочатку буде фільтрувати за критеріями текст, створюючи два нових з пробілами та без відповідно. Далі обчислюємо значення частот символів та їх ентропію, використовуючи формулу:

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

Далі рахуємо частоти біграм з перетинами та без; з пробілами та без. Наступний етап – обчислення ентропії, використовуємо наступну формулу:

$$H_2 = -\sum_i i, j p(i, j) \log_2 p(i, j)/2$$

Потім рахуємо надлишковість за допомгою іншої формули:

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

Результати виконання:

Частота символів без пробілів:

- o 0.11223137474889013
- e 0.08425977281232766
- a 0.0795172677347457
- и 0.07499155384689636
- н 0.06438302730698661
- т 0.05865942341455618
- c 0.05533713799233952
- л 0.04760544636435636
- p 0.04738865517462372
- в 0.04736808374786077
- м 0.032534502634329436
- к 0.03093626101659244
- д 0.03041801930391039
- п 0.02819709642223328
- y 0.025236393385<u>81157</u>
- я 0.019904437810598873

- ы 0.019145668646534626
- г 0.01809336104673749
- ь 0.017603602848034916
- 6 0.016868569945620224
- з 0.016729317210609475
- ч 0.012958416443216136
- й 0.011432966027871118
- x 0.010306284808238699
- ж 0.009619515637844779
- ш 0.007282285074084829
- ю 0.0060923071567202496
- ц 0.005410285238656243
- ш 0.0037985930726511584
- a 0.0035588568299906084
- Ф 0.0018632965394904141
- ъ 0.00026584613047506546
- ë 2.3736261649559415e-06

```
Частота символів з пробілами:
  0.15449145620195515
o 0.09489258623238676
e 0.07124235781130525
a 0.06723252924920409
и 0.06340599949024202
н 0.05443639966364001
т 0.049597043671274335
c 0.04678802296185445
л 0.04025081163238288
p 0.040067512829243784
в 0.040050119512157596
м 0.027508199945545538
к 0.026156873002695296
д 0.025718695206870094
п 0.023840885935295524
v 0.021337586222352153
```

я 0.0168293722283582

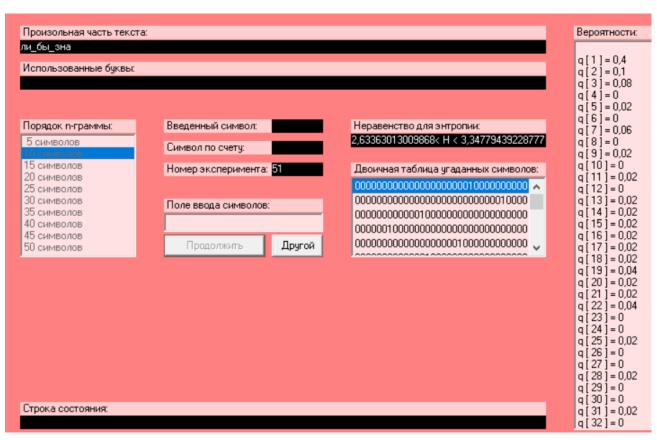
```
ы 0.016187826417371376
г 0.015298091351039284
ь 0.014883996609641115
6 0.014262520010676821
з 0.014144780633477988
ч 0.010956451816832316
й 0.009666670457517826
x 0.008714051860181813
ж 0.008133382658996659
ш 0.006157234248511701
ю 0.0051510977524489455
ц 0.004574442393668297
ш 0.003211742897338622
э 0.0030090438559110854
Ф 0.0015754331437684763
ъ 0.00022477517465231762
ë 2.0069212022528358e-06
```

```
Частоти біграм без пробілу без перетину:
{'cт': 9032, 'то': 8801, 'на': 7671, 'но': 7281, 'ен': 7012, 'ов': 6964, 'по': 6758, 'ни': 6708, 'ра': 6584, 'ос': 6423, 'ко': 6102, 'не': 6080, 'го': 5892, 'во': 5886, 'ли': 51876, 'го': 8.01492405656197384, 'то': 8.013926866937576747, 'на': 8.012138733811856747, 'но': 8.011521590528679847, 'ен': 8.011095919891635967, 'ов': 8.01101996379426025, 'по': 8.01101906379426025, 'по': 8.0110190637942602, 'по': 8.01101906379426025, 'по': 8.01101906379426025, 'по': 8.0110190637942602, 'по': 8
```

```
Ентропія Н1 без пробілів: 4.457223768425505
Ентропія Н1 з пробілами: 4.389587744688843
Ентропія Н2 без пробілів без перетину: 4.149596664532538
Ентропія Н2 без пробілів з перетином: 4.150057053053292
Ентропія Н2 з пробілами без перетину: 3.9976662350370633
Ентропія Н2 з пробілами з перетином: 3.9983739015880295

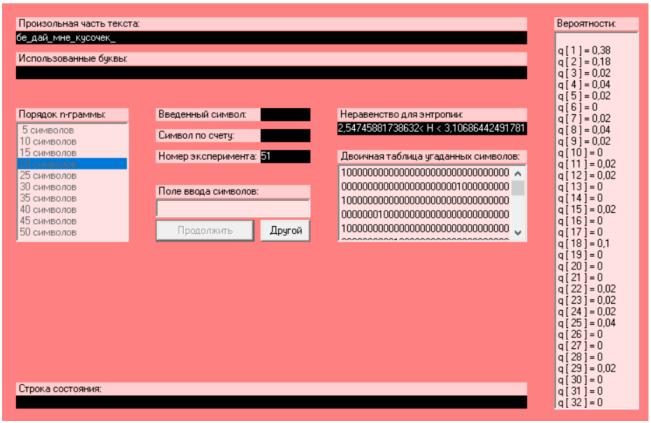
Надлишковість Н1 без пробілів: 0.11640057002675763
Надлишковість Н1 з пробілами: 0.13717546807476644
Надлишковість Н2 без пробілів без перетину: 0.1773845250100553
Надлишковість Н2 без пробілів з перетином: 0.17729325765269577
Надлишковість Н2 з пробілами без перетину: 0.21421219971907213
Надлишковість Н2 з пробілами з перетином: 0.2140730996267376
```

$H^{(10)}$:



0.330441121542446 < R < 0.473273973980264

Лабороторная работа №1



0.378627115016438 < R < 0.490508236522736

 $H^{(30)}$:

Лабороторная работа №1



0,432443812140248 < R < 0.574843682564354

Висновки

Під час виконання цієї лаборатоної роботи, я набув навичок оцінки частот букв та біграм. За допомогою практичної частини я навчився розраховувати ентропію та надлишковість. Також цікавою виявилася можливість застосування цих знань при роботі з CoolPinkProgram. Я використовував ці дані для передбачення наступного символу, вибираючи його на основі ймовірності, хоча в деяких випадках всеодно виникали складності.