

Тарасов Микита

Лабораторна №1

**Мета:** 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.

Код для виконання поставленої задачі міститься у файлі **tarasov\_lab1.py**. Даний код рахує к-сть заданих літер чи біграм у тексті, їх частоту, і на основі частот – ентропію тексту. Для початку підготовлюємо наш текст (залишаємо лише літери а-я та '\_' замість пробілів):

```
5     #prepare text
6     file = "text.txt"
7     with open(file, "r", encoding="utf-8") as file:
8         text = file.read().replace(' ', '_')
9     text = text.lower()
10    text = text.replace('-', '_')
11    text = re.sub(r'\n+', '_', text)
12    text = re.sub(r'_+', '_', text)
13
14    def text_without_spaces(text):
15        return text.replace('_', '')
16    text_without_spaces = text_without_spaces(text)
```

Після цього рахуємо частоту літер, що трапляються в тексті, а далі їх частоти. Так само робимо із біграмами що перетинаються на тексті із пробілами та без, і потім із біграмами, що не перетинаються на тексті із пробілами та без. Результати виконання даного коду наведені у файлі **freqAndEntropy.txt** (нижче наведено лише частину відповідного файлу із результатами)

```
1   _ : Amount: 112535 [+] Frequency: 0.162586
2   о : Amount: 64668 [+] Frequency: 0.09343
3   е : Amount: 48471 [+] Frequency: 0.070029
4   а : Amount: 47298 [+] Frequency: 0.068334
5   н : Amount: 38392 [+] Frequency: 0.055467
6   и : Amount: 38155 [+] Frequency: 0.055125
7   т : Amount: 35364 [+] Frequency: 0.051092
8   л : Amount: 30226 [+] Frequency: 0.043669
9   с : Amount: 30014 [+] Frequency: 0.043363
10  р : Amount: 26852 [+] Frequency: 0.038795
11  в : Amount: 23792 [+] Frequency: 0.034374
12  к : Amount: 20352 [+] Frequency: 0.029404
13  м : Amount: 18890 [+] Frequency: 0.027291
14  д : Amount: 18425 [+] Frequency: 0.02662
15  у : Amount: 16853 [+] Frequency: 0.024348
16  п : Amount: 16196 [+] Frequency: 0.023399
17  я : Amount: 11957 [+] Frequency: 0.017275
18  ы : Amount: 11172 [+] Frequency: 0.016141
19  ь : Amount: 10642 [+] Frequency: 0.015375
20  г : Amount: 10237 [+] Frequency: 0.01479
21  б : Amount: 10081 [+] Frequency: 0.014565
22  з : Amount: 10011 [+] Frequency: 0.014463
23  ч : Amount: 8171 [+] Frequency: 0.011805
24  й : Amount: 6529 [+] Frequency: 0.009433
25  ж : Amount: 6023 [+] Frequency: 0.008702
26  х : Amount: 5107 [+] Frequency: 0.007378
27  ш : Amount: 5077 [+] Frequency: 0.007335
28  ю : Amount: 3193 [+] Frequency: 0.004613
29  щ : Amount: 2167 [+] Frequency: 0.003131
30  э : Amount: 2151 [+] Frequency: 0.003108
31  ц : Amount: 1906 [+] Frequency: 0.002754
32  ф : Amount: 1079 [+] Frequency: 0.001559
33  ь : Amount: 170 [+] Frequency: 0.000246
34  ё : Amount: 2 [+] Frequency: 3e-06
35
36  Entropy: 4.379028773060594
```

Відповідно, частоти наведені у файлі, згаданому вище. Ентропія також обрахована у цьому файлі, але для подальшого обрахування надлишковості запишемо результати нижче:

	Ентропія	Надлишковість
<b>H1 (із пробілом)</b>	4.379028773060594	0.1392504760606287
<b>H1 (без пробілу)</b>	4.464420005196533	0.1149732663024602
<b>H2 (із пробілом, з перетином)</b>	3.9903784772111153	0.2156442552450309
<b>H1 (без пробілів, з перетином)</b>	4.162145373626576	0.1748961968391469
<b>H1 (із пробілом, без перетину)</b>	3.9903267703185494	0.2156544188419076
<b>H1 (без пробілів, без перетину)</b>	4.160968545981837	0.1751294911809283

Надлишковість обчислюється за формулою

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

Де  $H_0$  це 5.08746 для тексту із пробілами ( $\log_2(34)$ ), 5.04439 для текстів без пробілів ( $\log_2(33)$ ).

Переходимо до програми **CoolPinkProgram**

Результати виконання наведені нижче (**H10, H20 та H30 відповідно**):

Лабораторная работа №1

Произвольная часть текста:

\_у\_него\_б

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

Поле ввода символов:

Продолжить

Другой

Неравенство для энтропии:

2.47329333847461 < H < 3.25455575012416

Двоичная таблица угаданных символов:

00000000000000010000000000000000

10000000000000000000000000000000

00000000100000000000000000000000

00000000000001000000000000000000

00000000000000000000000000000001

Вероятности:

q[1] = 0.36

q[2] = 0.14

q[3] = 0.08

q[4] = 0.06

q[5] = 0.04

q[6] = 0.02

q[7] = 0.04

q[8] = 0

q[9] = 0.08

q[10] = 0

q[11] = 0.02

q[12] = 0.02

q[13] = 0

q[14] = 0.02

q[15] = 0.02

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] = 0.02

q[26] = 0.02

q[27] = 0

q[28] = 0.02

q[29] = 0.02

q[30] = 0

q[31] = 0

q[32] = 0.02

Строка состояния:



Отримуємо наступне:

	Ентропія	Надлишковість
<b>H(10)</b>	$2.47329333847461 < H < 3.25455575012416$	$0.349088849975168 < R < 0.505341332305078$
<b>H(20)</b>	$1.6494715378956 < H < 2.39603171251419$	$0.520793657497162 < R < 0.67010569242088$
<b>H(30)</b>	$1.98699074895968 < H < 2.57640234716675$	$0.48471953056665 < R < 0.602601850208064$

**Висновки:**

Під час виконання даного практикуму, я навчився підготовлювати текст для подальшого його аналізу, а також ознайомився із поняттям ентропії та навчився обраховувати частоти літер та біграм, а також, користуючись цими значеннями, обраховувати ентропію та надлишковість у різних текстах.