Лабораторна робота 1 з Симетричної Криптографії

Команда: Бондар, Дигас Група: ФІ-03

Підготовка даних

- 1. В якості середньостатистичного тексту російською мовою ми взяли дописи з телеграм-каналу терориста і військового злочинця ігоря гіркіна
- 2. Повний вхідний текст можна знайти у відповідному файлі
- 3. Обробка тексту і підготовка до аналізу відбувається у розділі "Text reading and preprocessing"

Програмна частина

Text reading and preprocessing

```
In [ ]: filename = "girkin crying.txt"
In [ ]: def get text( filename):
           f = open( filename, "r", encoding='utf-8')
           text = f.read()
           f.close()
            return text
        def transform_symbol(_c):
           if 'a' <= _c and _c <= 'я':
              return _c
           elif _c \leftarrow '9' and _c \rightarrow '4':
               return _c.lower()
           elif _c == 'Ë' or _c == 'ë':
               return 'e'
            else:
               return ' '
        def preprocess_text(_text):
            text = get text(filename)
            text formatted = ""
            # Change symbols according to requirements
            for c in _text:
               text formatted += transform symbol(c)
            # Remove consequtive spaces
            text_formatted = ' '.join(text_formatted.split())
            return text_formatted
In [ ]: text = preprocess_text(get_text(filename))
```

Text processing (singular char count and bigram count)

```
c\_count[c] = c\_count[c] + 1
           return dict(sorted(c_count.items()))
       # Bigrams with intersection (ex: [1, 2], [2, 3], [3, 4])
       def count_bigrams_w_i(_text):
           b_count = {}
           prev_char = _text[0]
           for c in _text[1:]:
               bg = prev_char + c
               prev_char = c
               if bg not in b_count:
                   b_count[bg] = 1
               else:
                   b_count[bg] = b_count[bg] + 1
           return dict(sorted(b_count.items()))
       # Bigrams without intersection (ex: [1, 2], [3, 4])
       def count_bigrams_wo_i(_text):
           b_count = {}
           i = 1
           while i < len(_text):</pre>
               bg = _text[i - 1] + _text[i]
               if bg not in b_count:
                   b count[bg] = 1
               else:
                   b_count[bg] = b_count[bg] + 1
               i = i + 2
           return dict(sorted(b_count.items()))
In [ ]: chars_freq_wspaces = count_chars(text)
       chars_freq_wospaces = chars_freq_wspaces.copy()
       del chars_freq_wospaces[' ']
       bigrams_freq_w_intersect = count_bigrams_w_i(text)
       bigrams_freq_wo_intersect = count_bigrams_wo_i(text)
```

Show symbol frequencies

```
In [ ]: for k, v in chars_freq_wspaces.items():
    print(f"{k} : {v}")
```

```
: 86793
a: 37904
6 : 8174
в: 22994
г: 7676
д: 14226
e: 43426
ж : 4720
з: 7408
и: 37250
й: 6380
к: 16426
л : 18174
м: 15420
н: 36030
o: 56742
п : 15740
p: 24918
c : 26566
т: 31194
y: 12712
ф: 1612
x : 5356
ц: 2422
ч: 6828
ш: 3240
щ: 1846
ъ: 164
ы: 8952
ь: 7650
э: 1236
ю: 3406
я: 9048
```

Show bigram frequencies

```
In []: alph = list(chars_freq_wspaces.keys())
print("    ", end='|')
for l_c in alph:
    print(f" '\[1_c\]' ", end='|')
for l_c in alph:
    print(f\[1^n\]'\[1_c\]' ", end='|')
    for r_c in alph:
        k = l_c + r_c
        if k in bigrams_freq_w_intersect:
            print(f\[1^n\]'\[1_c\]' ", end='|')
        else:
            print(f\[1^n\]'\[1_c\]' ", end='|')
```

```
'д' |
                                    'e' |
                                             | 'з'
                                                   | 'и' |
                                                          'й' | 'к' | 'л' | 'м' |
                                                                                'н' | 'о' | 'п' |
                                                                                                'p' | 'c' |
                                                                                                           'т' | 'у' | 'ф' | 'х' | 'ц' |
                                                                                                                                      '4'
                                                                                                                                            'ш'
                                                                                                                                                       'ъ' |
               3314
                    9052 | 1754 |
                                    1146
                                               2046 | 5760
                                                            54 | 4112 |
                                                                     1376
                                                                          3208
                                                                                9268
                                                                                     5460 | 10444 |
                                                                                                3268
                                                                                                           4094 | 2800 |
                                                                                                                       728
                                                                                                                            749
                                                                                                                                             382
                                                                                                                                                                              72
                               3626
                                           514
                                                                                                     7726
                                                                                                                                  320
                                                                                                                                      2080
                                                                                                                                                    6
                                                                                                                                                               0
'a'| 7586|
           10
                488
                    1636
                          560
                               1036
                                     954
                                           856
                                               1386
                                                     298
                                                           468 | 2572 |
                                                                     2688
                                                                           1642
                                                                                3256
                                                                                            594
                                                                                                1914
                                                                                                     2356
                                                                                                           2938
                                                                                                                  48
                                                                                                                        48
                                                                                                                            672
                                                                                                                                  804
                                                                                                                                       690
                                                                                                                                             664
                                                                                                                                                  122
                                                                                                                                                         0
                                                                                                                                                                              762
                                                                                                                                                                                   798
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          41
'6'|
    152
          500
                      12
                            0
                                 20
                                    1048
                                            36
                                                  6
                                                      566
                                                             0
                                                                 32
                                                                      494
                                                                            74
                                                                                 264 1704
                                                                                             20
                                                                                                 582
                                                                                                      192
                                                                                                              6
                                                                                                                 666
                                                                                                                         0
                                                                                                                            106
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                                  294
                                                                                                                                                       140
                                                                                                                                                            1070
                                                                                                                                                                              16
                                                                                                                                                                                   136
                 2
                                                                                                                                                                    32
'B'
   4476
                                 94
                                                                      798
                                                                            108
                                                                                1220
                                                                                                 598
                                                                                                      1558
                                                                                                                 304
                                                                                                                                                                                   196
         2924
                 8
                      50
                            8
                                    2190
                                             2
                                                216 1430
                                                             0|
                                                                372
                                                                                     4218
                                                                                            168
                                                                                                            196
                                                                                                                              2
                                                                                                                                   28
                                                                                                                                        66
                                                                                                                                             302
                                                                                                                                                   12
                                                                                                                                                            1386
                                                                                                                                                                               0
'r'|
     256
                                                                      472
                                                                                 186
                                                                                     3774
          678
                 0
                      28
                            4
                                358
                                     288
                                            0
                                                  01
                                                     508
                                                             0|
                                                                 42
                                                                             2
                                                                                                 656
                                                                                                       26
                                                                                                             14
                                                                                                                 366
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                         6
                                                                                                                                              10
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          01
                                                                                                                                                                               2
                                                                                                                                                                                     0
     670
         2872
                10
                     278
                           78
                                 70
                                    2566
                                            18
                                                 14 | 1236 |
                                                             0
                                                                160
                                                                      532
                                                                            26
                                                                                 984
                                                                                     2134
                                                                                           122
                                                                                                 506
                                                                                                       412
                                                                                                             38|
                                                                                                                 754
                                                                                                                             14
                                                                                                                                   12
                                                                                                                                        68
                                                                                                                                             46
                                                                                                                                                    0
                                                                                                                                                         6
                                                                                                                                                             296
                                                                                                                                                                   134
                                                                                                                                                                         10
                                                                                                                                                                               8|
                                                                                                                                                                                   152
'д'|
                                                                                                                         01
'e' | 10900 |
          116
                340
                    1060 | 1234 |
                               1878
                                     810
                                           338
                                                972
                                                     132
                                                          1462
                                                                756
                                                                     2488
                                                                          2242
                                                                                5884
                                                                                      246
                                                                                            606
                                                                                                3992
                                                                                                     2790
                                                                                                           2974
                                                                                                                  74
                                                                                                                        52
                                                                                                                            498
                                                                                                                                  264
                                                                                                                                       428
                                                                                                                                             352
                                                                                                                                                  286
                                                                                                                                                         0
                                                                                                                                                                         14
                                                                                                                                                                             164
                                                                                                                                                                                    74
     70
          558
                30
                                                     558
                                                                                 608
                                                                                                       36
                                                                                                                         2
                                                                                                                                                         0
                                                                                                                                                                                     0
'ж'|
                       0
                            6|
                                372
                                    2300
                                                             0
                                                                 54
                                                                        0
                                                                                       10
                                                                                             2
                                                                                                              0
                                                                                                                                    0
                                                                                                                                        12
                                                                                                                                                               0|
                                                                                                                                                                    26
                                                                                                                                                                          0
                                                                                                                                                                               0
's'|
    784
         2248
                          110
                                                     496
                                                                                                 146
                                                                                                              2
                                                                                                                                                         8
                                                                                                                                                             154
                                                                                                                                                                              58
                                                                                                                                                                                  142
                102
                     458
                                472
                                     226
                                            32
                                                  6
                                                             0|
                                                                 86
                                                                       78
                                                                           336
                                                                                 546
                                                                                      496
                                                                                             60
                                                                                                       18
                                                                                                                 268
                                                                                                                         0|
                                                                                                                              2
                                                                                                                                  16
                                                                                                                                         2
                                                                                                                                                                    52
'u' | 10652 |
          278
                252
                    1570
                          278
                                592
                                    1938
                                           162
                                               1438 | 1206 |
                                                          1036 | 1560 |
                                                                     1846
                                                                           1914
                                                                                1914
                                                                                      282
                                                                                            176
                                                                                                 872
                                                                                                     1306
                                                                                                           2506
                                                                                                                  34
                                                                                                                        46 1334
                                                                                                                                  388 | 1140 |
                                                                                                                                             162
                                                                                                                                                   86
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0|
                                                                                                                                                                              488
                                                                                                                                                                                  1794
'й' | 4376|
            0
                 0
                           38
                                 68
                                                             0
                                                                 54
                                                                            42
                                                                                 330
                                                                                      222
                                                                                             10
                                                                                                       798
                                                                                                             82|
                                                                                                                   0
                                                                                                                         2
                                                                                                                              2
                                                                                                                                   44
                                                                                                                                       162
                                                                                                                                             138
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                               0
                                                                                                                                                                                     0
                       0 l
                                                                                                   41
'ĸ'
   2094
                                 2
                                     378
                                                  2 | 2264 |
                                                             0
                                                                            12
                                                                                             12
                                                                                                1302
                                                                                                                 728
                                                                                                                                   30
                                                                                                                                             10
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                                     0
         3156
                 8
                      76
                            2 |
                                           124
                                                                 16
                                                                      280
                                                                                  54
                                                                                      4862
                                                                                                       184
                                                                                                            814
                                                                                                                         0|
                                                                                                                             16
                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                               0
'л'
     984
         1970
                38
                       4
                          102
                                 58
                                    3106
                                           230
                                                  0 3374
                                                             0
                                                                216
                                                                      248
                                                                             0
                                                                                 430
                                                                                     2002
                                                                                             20
                                                                                                       160
                                                                                                             16
                                                                                                                 690
                                                                                                                         0
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                        60
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                             248
                                                                                                                                                                 2662
                                                                                                                                                                          0
                                                                                                                                                                              540
                                                                                                                                                                                  1010
'm'
    4472
         1430
                54
                       6
                            2 |
                                 6
                                    2322
                                                 10
                                                    2022
                                                             01
                                                                 52
                                                                      128
                                                                            72
                                                                                 568
                                                                                     2250
                                                                                            136
                                                                                                  18
                                                                                                       96
                                                                                                              0
                                                                                                                 944
                                                                                                                        20
                                                                                                                              0
                                                                                                                                   10
                                                                                                                                        10
                                                                                                                                              6
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                             590
                                                                                                                                                                   18
                                                                                                                                                                          0
                                                                                                                                                                               0
                                                                                                                                                                                   174
'H'
   1070
         6156
                50
                      14
                           98
                                338
                                    5634
                                             41
                                                 20
                                                    5800
                                                             0
                                                                352
                                                                        0
                                                                            18
                                                                                2452
                                                                                     6742
                                                                                             8|
                                                                                                 130
                                                                                                       608
                                                                                                           1052
                                                                                                                 902
                                                                                                                       114
                                                                                                                              2
                                                                                                                                  202
                                                                                                                                        62
                                                                                                                                             14
                                                                                                                                                   10
                                                                                                                                                         0
                                                                                                                                                            3244
                                                                                                                                                                   336
                                                                                                                                                                          0|
                                                                                                                                                                              38
                                                                                                                                                                                   560
               2788
                    4806 | 2624 |
                                                          2662
                                                                                                3440
                                                                                                                       142
                                                                                                                            202
                                                                                                                                             2861
                                                                                                                                                  128
                                                                                                                                                                         78
'o'|11081|
                               3082
                                    1624
                                         1222
                                                860
                                                     444
                                                                952
                                                                     3306 | 2980 |
                                                                                3056
                                                                                      362
                                                                                          1142
                                                                                                      4006
                                                                                                           3878
                                                                                                                  46
                                                                                                                                  106
                                                                                                                                       936
                                                                                                                                                         0
                                                                                                                                                               01
                                                                                                                                                                    01
                                                                                                                                                                              84
                                                                                                                                                                                   412
            6
'n'|
    214 | 1004 |
                 0
                      36
                           20
                                 2
                                    1618
                                                  2
                                                     460
                                                             0
                                                                 20
                                                                      720
                                                                                 118
                                                                                     6148
                                                                                           126
                                                                                                4460
                                                                                                       22
                                                                                                             40
                                                                                                                 394
                                                                                                                         0
                                                                                                                              0
                                                                                                                                   10
                                                                                                                                         4
                                                                                                                                              2
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                             230
                                                                                                                                                                    2
                                                                                                                                                                          2
                                                                                                                                                                               0
                                                                                                                                                                                    86
                          194
                                122 | 3738 |
                                                                                                       418
                                                                                                                             58
                                                                                                                                                                              32
'p'|
     628
         5210
                56
                     302
                                           214
                                                 38 2414
                                                             0
                                                                108
                                                                       24
                                                                           546
                                                                                 554
                                                                                     5806
                                                                                             94
                                                                                                 130
                                                                                                            516
                                                                                                                 1468
                                                                                                                       336
                                                                                                                                   16
                                                                                                                                        20
                                                                                                                                             132
                                                                                                                                                    8|
                                                                                                                                                         0
                                                                                                                                                             880
                                                                                                                                                                   284
                                                                                                                                                                         10
                                                                                                                                                                                   562
'c'| 1684|
          774
                62
                     758
                           20
                                160
                                    1846
                                             8
                                                  0 | 1212 |
                                                             0 | 2840 |
                                                                     1436
                                                                           270
                                                                                 564 | 1878 |
                                                                                          1008
                                                                                                 328
                                                                                                       766
                                                                                                           6960
                                                                                                                 842
                                                                                                                        18
                                                                                                                            118
                                                                                                                                   26
                                                                                                                                       160
                                                                                                                                             92
                                                                                                                                                    0
                                                                                                                                                         6
                                                                                                                                                             138
                                                                                                                                                                  742
                                                                                                                                                                          0|
                                                                                                                                                                              38
                                                                                                                                                                                  1812
'T'| 3914|
         3734
                18
                    1694
                           16
                                 52
                                    2902
                                             0
                                                  4 | 3088 |
                                                             0
                                                                402
                                                                       68
                                                                            38
                                                                                 836 | 6562 |
                                                                                             20
                                                                                                1550
                                                                                                      1180
                                                                                                             20
                                                                                                                 940
                                                                                                                         2
                                                                                                                             32
                                                                                                                                    0
                                                                                                                                        20
                                                                                                                                              01
                                                                                                                                                         0
                                                                                                                                                             622
                                                                                                                                                                  3074
                                                                                                                                                                          8
                                                                                                                                                                              16
                                                                                                                                                                                   380
                                                                                                                                                    2
'y' | 2916
          136
                296
                     188
                           392
                               1192
                                     176
                                           784
                                                 80
                                                       8
                                                            36
                                                                730
                                                                      318
                                                                            432
                                                                                 182
                                                                                        4
                                                                                            822
                                                                                                 388
                                                                                                       828
                                                                                                            878
                                                                                                                   0
                                                                                                                         0
                                                                                                                             94
                                                                                                                                   22
                                                                                                                                       658
                                                                                                                                             130
                                                                                                                                                  254
                                                                                                                                                         0
                                                                                                                                                               2
                                                                                                                                                                    0
                                                                                                                                                                          2
                                                                                                                                                                              736
                                                                                                                                                                                    28
    356
          122
                 0
                       0
                            0
                                  0
                                     198
                                             0|
                                                  0|
                                                     160
                                                             0|
                                                                  0
                                                                       60
                                                                             2
                                                                                  10
                                                                                      220
                                                                                             0
                                                                                                 392
                                                                                                        10
                                                                                                              8|
                                                                                                                  22
                                                                                                                        46
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                         2
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                               4
                                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                               0
                                                                                                                                                                                     01
'x'| 3228|
          274
                 2
                     124
                            0
                                 2
                                     176
                                             0|
                                                  0|
                                                      80
                                                             0
                                                                  0
                                                                       18
                                                                            92
                                                                                 116 | 1020 |
                                                                                             0
                                                                                                  70
                                                                                                        44
                                                                                                             30|
                                                                                                                  76
                                                                                                                         01
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                              0
                                                                                                                                                    0
                                                                                                                                                         4
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0|
                                                                                                                                                                               0
                                                                                                                                                                                     0
'ц'|
     92
          188
                 0
                      16
                            0
                                 62
                                     558
                                             0
                                                  0 | 1188
                                                             0
                                                                102
                                                                        2
                                                                             0
                                                                                  12
                                                                                       60
                                                                                             12
                                                                                                        2
                                                                                                              2
                                                                                                                  38
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                              0
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                              88
                                                                                                                                                                               0
                                                                                                                                                                                     0
                                                                                                                                                                          01
                                                                                                        0
'4'
     102
         1304
                 0
                      20
                                 0
                                    2062
                                                  0
                                                     976
                                                             0
                                                                102
                                                                        4
                                                                             2
                                                                                 694
                                                                                       12
                                                                                             0
                                                                                                  20
                                                                                                           1242
                                                                                                                 116
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                              54
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                   116
                                                                                                                                                                               0
                                                                                                                                                                                     0
'ш'|
     62
          282
                 0
                                  0 | 1020 |
                                                  01
                                                     920
                                                             0
                                                                 82
                                                                      144
                                                                             4
                                                                                 220
                                                                                      138
                                                                                             2|
                                                                                                        0
                                                                                                            204
                                                                                                                                         01
                                                                                                                                                               0
                                                                                                                                                                                     0
'щ'|
          180
                 0
                       0
                            0
                                  0
                                     958
                                             0|
                                                  0|
                                                     586
                                                             0
                                                                  0
                                                                        0
                                                                             0
                                                                                  48
                                                                                        01
                                                                                                        0
                                                                                                              0
                                                                                                                  36
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    32
                                                                                                                                                                          0
                                                                                                                                                                               0|
                                                                                                                                                                                     0
'ъ'|
            0
                 0
                       0
                            0
                                 0
                                      86
                                             0
                                                  0
                                                       0
                                                             0
                                                                  0
                                                                        0
                                                                             0
                                                                                  0
                                                                                        0
                                                                                             01
                                                                                                   0
                                                                                                        0
                                                                                                              0
                                                                                                                   0
                                                                                                                         0
                                                                                                                              0
                                                                                                                                    0
                                                                                                                                         0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                               0
                                                                                                                                                                                    74
'ы'|
   2588
            0
                50
                     548
                           58
                                 70
                                             2
                                                           638
                                                                            878
                                                                                  76
                                                                                        0
                                                                                                       276
                                                                                                            372
                                                                                                                   4
                                                                                                                         0 | 1254 |
                                                                                                                                    0
                                                                                                                                                                               0
                                                                                                                                                                                    36|
                                    1222
                                                 26
                                                      24
                                                                 32
                                                                      462
                                                                                            120
                                                                                                  48
                                                                                                                                        46
                                                                                                                                             122
                                                                                                                                                         0
                                                                                                                                                               0|
                                                                                                                                                                    0
                                                                                                                                                                          0
'ь'
    4180
            0
                40
                      10
                            8
                                 16
                                     158
                                             01
                                                 84
                                                       40
                                                             0
                                                                 504
                                                                        0
                                                                            56
                                                                                1260
                                                                                       72
                                                                                                   01
                                                                                                       476
                                                                                                             90
                                                                                                                   0
                                                                                                                         4
                                                                                                                              0
                                                                                                                                   56
                                                                                                                                        12
                                                                                                                                             306
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                             218
                                                                                                                                                                                    52
'э'|
     26
            0
                 2
                      22
                            01
                                 2
                                       0
                                             0
                                                  0
                                                       0|
                                                             4
                                                                 54
                                                                       54
                                                                            14
                                                                                  16
                                                                                        0
                                                                                             4
                                                                                                  18
                                                                                                       18
                                                                                                            930
                                                                                                                   2 |
                                                                                                                        50
                                                                                                                              0
                                                                                                                                    4
                                                                                                                                         0
                                                                                                                                             12
                                                                                                                                                    0
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0|
                                                                                                                                                                               0
                                                                                                                                                                                     4
'ю' | 1598 |
                                                             2
                                                                            58
                                                                                        0
                                                                                                  40
                                                                                                       88
                                                                                                                   0
                                                                                                                              0
                                                                                                                                                                              20
                                                                                                                                                                                     0
            4
               124
                       0
                           26
                               194
                                       4
                                            36
                                                 16
                                                       0
                                                                  8
                                                                       4
                                                                                  6
                                                                                             0
                                                                                                            556
                                                                                                                         0
                                                                                                                                   4
                                                                                                                                       120
                                                                                                                                              2
                                                                                                                                                  496
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0
            0
                40
                     226
                           44
                               280
                                    170
                                          130
                                                178
                                                                           348
                                                                                298
                                                                                                            736
                                                                                                                   0
'я' | 5574
                                                       0
                                                            18
                                                                 54
                                                                      126
                                                                                        4
                                                                                             10
                                                                                                  40
                                                                                                      166
                                                                                                                            200
                                                                                                                                   60
                                                                                                                                        64
                                                                                                                                              6
                                                                                                                                                  136
                                                                                                                                                         0
                                                                                                                                                               0
                                                                                                                                                                    0
                                                                                                                                                                          0
                                                                                                                                                                             114
                                                                                                                                                                                    26
```

Calculate H_1 and H_2

H2 = 3.9881215496418245

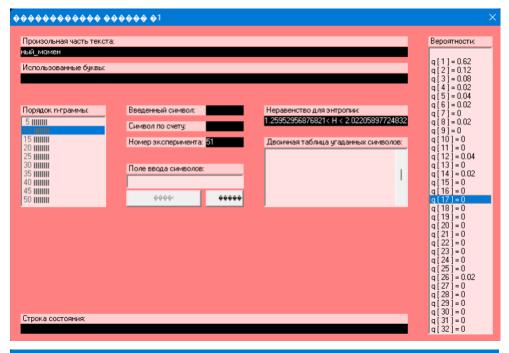
```
In []: import math

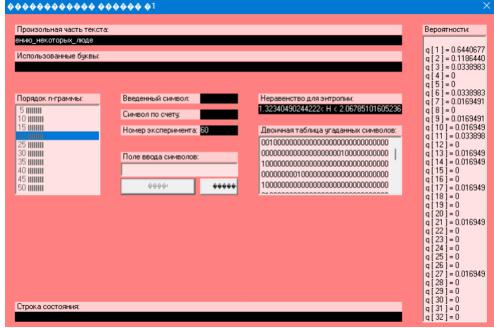
char_amount = sum(chars_freq_wspaces.values())
t1 = [chars_freq_wspaces[k] / char_amount for k in chars_freq_wspaces.keys()]
H1 = -sum(a * math.log2(a) for a in t1)

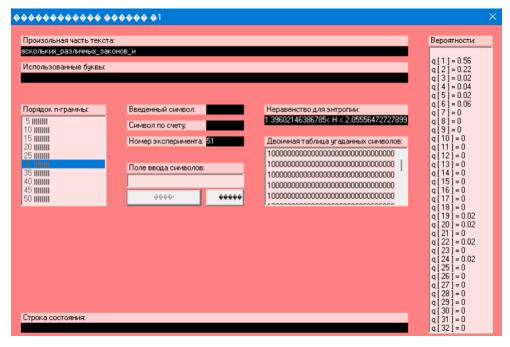
bg_amount = sum(bigrams_freq_w_intersect.values())
t2 = [bigrams_freq_w_intersect[k] / bg_amount for k in bigrams_freq_w_intersect.keys()]
H2 = -sum(a * math.log2(a) for a in t2) / 2
print(f"H1 = {H1}")
print(f"H2 = {H2}")

H1 = 4.385129362944809
```

Обчислення $H^{(10)}$, $H^{(20)}$, $H^{(30)}$ за допомогою Cool Pink Program







Отримані результати:

$$1.2595 \le H^{(10)} \le 2.0221 \tag{1}$$

$$1.3234 \le H^{(20)} \le 2.0679 \tag{2}$$

$$1.3960 \le H^{(30)} \le 2.0556 \tag{3}$$

Обчислимо надлишковість:

Так як $H_0 = \log_2(32) = 5$, то за формулою надлишковості $R = 1 - rac{H_\infty}{H_0}$:

Розглядаємо $H^{(i)}$ як наближення H_{∞} .

$$H^{(10)}: 0.404 \le R \le 0.748$$

 $H^{(20)}: 0.464 \le R \le 0.735$
 $H^{(30)}: 0.588 \le R \le 0.721$ (4)