Міністерсво освіти і науки України Національний технічний університет України Київський політехнічний інститут імені Ігоря Сікорського Навчально-науковий фізико-технічний інститут

Симетрична криптографія Комп'ютерний практикум №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

Виконали:

Медведцький Костянтин ФІ-04 Сковрон Роман ФІ-04

Мета роботи

Засвоєння понять ентропії на символ джерела та його надлишковості, вивчення та порівняння різних моделей джерела відкритого тексту для наближеного визначення ентропії, набуття практичних навичок щодо оцінки ентропії на символ джерела.

Порядок виконання роботи

- 0. Уважно прочитати методичні вказівки до виконання комп'ютерного практикуму.
- 1. Написати програми для підрахунку частот букв і частот біграм в тексті, а також підрахунку Н1 та Н2 за безпосереднім означенням. Підрахувати частоти букв та біграм, а також значення Н1 та Н2 на довільно обраному тексті російською мовою достатньої довжини (щонайменше 1Мб), де імовірності замінити відповідними частотами. Також одержати значення Н1 та Н2 на тому ж тексті, в якому вилучено всі пробіли.
- 2. За допомогою програми CoolPinkProgram оцінити значення H(10), H(20), H(30).
- 3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Хід роботи

Для роботи було обрано твір "Війна і мир"

1.1 Кількість та частота літер, якщо вважати, що пробіл не входить до алфавіту:

```
Monogram frequency without spaces:
Monograms without spaces:
                            o: 11.540391238792298 %
o: 265208
                            e: 8.236945108644807 %
e: 189292
                            a: 8.190993805712148 %
a: 188236
                            и: 6.786302508392858 %
и: 155955
                            н: 6.479264256979181 %
н: 148899
                            т: 5.861445507890265 %
T: 134701
                            c: 5.329828966291499 %
c: 122484
                            л: 5.016045968708729 %
л: 115273
                            B: 4.6268088430149446 %
B: 106328
                            p: 4.4283392476779575 %
p: 101767
                            κ: 3.3624517805042027 %
к: 77272
                            л: 3.0659005215211796 %
д: 70457
                            m: 2.9640766116135824 %
м: 68117
                            y: 2.722136039354506 %
y: 62557
                            n: 2.6334970203452004 %
п: 60520
                            я: 2.2014416351005295 %
я: 50591
                            ь: 1.9919628734359258 %
ь: 45777
                            r: 1.9695964248493854 %
r: 45263
                            ы: 1.9370040707806717 %
ы: 44514
                            6: 1.7395788232376086 %
6: 39977
                            з: 1.7139052733036417 %
з: 39387
                            ч: 1.4454643757737422 %
ч: 33218
                            й: 1.1201935524578073 %
й: 25743
                            ж: 1.0659744961565827 %
ж: 24497
                            ш: 0.92768544244447 %
ш: 21319
                            x: 0.8606731256676754 %
x: 19779
                            ю: 0.6139024448617001 %
ю: 14108
                            ц: 0.36029998890380466 %
ц: 8280
                            э: 0.30677716446519604 %
э: 7050
                            щ: 0.2962901720345418 %
щ: 6809
                            φ: 0.20482271108335856 %
φ: 4707
```

1.2 Кількість та частота літер, якщо вважати, що пробіл входить до алфавіту:

```
Monograms with spaces:
                         Monogram frequency with spaces:
 : 443926
                           : 16.189796466899658 %
o: 265208
                         o: 9.672025385747904 %
e: 189292
                         e: 6.903400460464966 %
a: 188236
                         a: 6.864888579950993 %
и: 155955
                         и: 5.687613944655948 %
н: 148899
                         н: 5.430284561221672 %
T: 134701
                         T: 4.912489410144598 %
c: 122484
                         c: 4.466940504615043 %
л: 115273
                         л: 4.203958335688661 %
B: 106328
                         в: 3.8777379084183106 %
p: 101767
                         p: 3.711400136615061 %
к: 77272
                         κ: 2.8180776809429284 %
д: 70457
                         д: 2.5695374672092854 %
м: 68117
                         m: 2.4841986410703676 %
y: 62557
                         y: 2.281427755030888 %
п: 60520
                         n: 2.2071392127894454 %
я: 50591
                         я: 1.845032715040166 %
ь: 45777
                         ь: 1.6694681385304433 %
г: 45263
                         r: 1.6507227724469375 %
ы: 44514
                         ы: 1.6234070541657202 %
6: 39977
                         6: 1.4579445523741517 %
3: 39387
                         з: 1.436427497920322 %
ч: 33218
                         ч: 1.2114466353344315 %
й: 25743
                         й: 0.9388364962795555 %
ж: 24497
                         ж: 0.8933953948397726 %
ш: 21319
                         ш: 0.7774950574596527 %
x: 19779
                         x: 0.7213318983767752 %
ю: 14108
                         ю: 0.5145128885332699 %
ц: 8280
                         ц: 0.30196815403001664 %
э: 7050
                         э: 0.25711056593135473 %
щ: 6809
                         щ: 0.24832139623072264 %
                         φ: 0.17166233104097686 %
φ: 4707
```

1.3 Кількість та частота біграм без перетинів та без пробілів:

```
BiGram with no crossing and no spaces:
                                        Bigram frequency without spaces and crossing:
                                         то: 1.7593786824154383 %
то: 20216
                                         ст: 1.349907140034916 %
ст: 15511
                                        на: 1.2535660141230693 %
на: 14404
                                         oB: 1.1573989462526173 %
ов: 13299
                                         ал: 1.0764619570041827 %
ал: 12369
                                        oc: 1.0620151395684405 %
oc: 12203
                                         го: 1.055749050078239 %
                                        не: 1.053921440643597 %
не: 12110
                                         но: 1.0478294091947902 %
но: 12040
                                        он: 1.0369507816076349 %
он: 11915
                                        ен: 1.0252888928342045 %
ен: 11781
                                        по: 1.0133659169986824 %
по: 11644
ко: 11570
                                        ко: 1.0069257694670863 %
                                         pa: 0.971504958043309 %
                                        от: 0.9584506049387228 %
от: 11013
ни: 11010
                                        ни: 0.9581895178766311 %
                                        во: 0.9299921151707249 %
ка: 9741
                                        ка: 0.8477496906118314 %
                                        ол: 0.8290384511619244 %
                                        ло: 0.818420910636861 %
ло: 9404
                                        ep: 0.8060627896978527 %
                                         ла: 0.8022335127871739 %
ла: 9218
                                        ли: 0.7984042358764953 %
ли: 9174
                                        op: 0.7895272757653768 %
                                        po: 0.7844795925649368 %
                                        ел: 0.7548026965071772 %
                                         ан: 0.7539324063002049 %
ан: 8663
                                        np: 0.7393985598437655 %
пр: 8496
                                        pe: 0.6884865827358791 %
                                        од: 0.6646406310648348 %
                                         та: 0.6613335282783397 %
та: 7599
```

1.4 Кількість та частота біграм з перетинами, але без пробілів:

```
BiGram with crossing and without spaces:
                                          Bigram frequency with crossing and without spaces:
                                          то: 1.757246471408356 %
ст: 31163
                                          ст: 1.3560426859940717 %
на: 28685
                                           на: 1.2482137293501891 %
ов: 26796
                                          OB: 1.1660148193016444 %
                                          ал: 1.0696301788794491 %
                                          го: 1.0612753928925138 %
не: 24252
                                          не: 1.0553139049747529 %
oc: 24107
                                          oc: 1.049004300974203 %
но: 24062
                                          но: 1.047046148008515 %
ен: 23647
                                          ен: 1.0289876262138373 %
он: 23634
                                           on: 1.0284219375793051 %
по: 23274
                                          no: 1.0127567138538016 %
ко: 23103
                                          ко: 1.0053157325841875 %
pa: 22360
                                          pa: 0.9729844513951622 %
ни: 21891
                                          ни: 0.9525761460416591 %
во: 21300
                                          во: 0.9268590704256242 %
ка: 19484
                                          ка: 0.8478367196325286 %
ол: 19141
                                          ол: 0.8329112425829517 %
                                          ep: 0.8169849317953565 %
ло: 18707
                                          ло: 0.8140259450916503 %
ла: 18483
                                          ла: 0.8042786947735591 %
                                          ли: 0.7918770593242022 %
op: 18154
                                          op: 0.7899624208688629 %
                                          po: 0.7889615871308446 %
ел: 17340
                                          ел: 0.7545416094450855 %
                                          ан: 0.7499725858584804 %
ан: 17235
                                          пр: 0.7319575785741513 %
                                          pe: 0.6972765138263005 %
                                          од: 0.6647276600855321 %
од: 15276
```

1.5 Кількість та частота біграм без перетинів, але з пробілами:

```
BiGram with no crossing and with spaces: Bigram frequency with spaces and no crossing:
0: 29644
                                          o : 2.16220947407194 %
и: 26048
                                          и: 1.8999201315823062 %
e : 23013
                                          e: 1.6785496770617174 %
a: 22764
                                          a : 1.660387817695778 %
c: 21769
                                           c: 1.5878133194262603 %
п: 21220
                                           п: 1.5477697017881045 %
н: 21158
                                           н: 1.5432474717451796 %
 в: 20292
                                           B: 1.4800821295327151 %
то: 19707
                                          то: 1.437412700901893 %
0: 16425
                                           o: 1.198026265403846 %
и: 15365
                                           и: 1.1207107195086816 %
                                          CT: 1.1202001451489965 %
я: 14706
                                          я: 1.0726437905040462 %
на: 14339
                                          на: 1.0458751062177016 %
ь: 13522
                                          ь: 0.9862837845230324 %
                                          к: 0.9788439867104789 %
к: 13420
го: 12147
                                          го: 0.8859923924420408 %
не: 12034
                                          не: 0.8777502634928392 %
ал: 11978
                                         ал: 0.8736656686153588 %
но: 11684
                                          но: 0.8522215455085868 %
по: 11571
                                         по: 0.8439794165593854 %
pa: 11139
                                          pa: 0.8124696846473937 %
ко: 10988
                                         ко: 0.8014558663170447 %
ов: 10450
                                          oB: 0.7622145798155368 %
ни: 10429
                                          ни: 0.7606828567364816 %
                                          д: 0.7501796127658178 %
                                          й: 0.7455115043344116 %
во: 10170
                                          во: 0.7417916054281348 %
```

1.6 Кількість та частота біграм з перетинами та з пробілами:

BiGram with crossing and spaces: Bigram frequency with spaces and crossing: o : 2.1647623458703653 % o: 59358 и : 1.9050623447762773 % и : 52237 e: 1.6786955554501988 % e: 46030 a: 1.6590384426023244 % a: 45491 c: 43324 c: 1.5800088256425031 % п: 1.546493265888892 % п: 42405 н: 1.5463838570975306 % н: 42402 B: 40699 B: 1.4842761332015566 % то: 1.4443419243547617 % то: 39604 o: 1.204189627317187 % o: 33019 и: 1.1239200440552732 % и: 30818 ст: 1.1088945700416848 % CT: 30406 я: 1.0792447875828317 % я: 29593 на: 1.0422281465056655 % на: 28578 к: 26929 κ: 0.9820897808541909 % ь: 26912 ь: 0.9814697977031448 % го: 0.8831112942695322 % го: 24215 не: 24072 не: 0.8778961418813206 % ал: 23882 ал: 0.8709669184284521 % но: 23322 но: 0.8505439440410502 % по: 23267 no: 0.8485381161994303 % pa: 22261 pa: 0.8118497014963475 % ко: 22145 ко: 0.8076192282303857 % ни: 0.7683779417288777 %

1.7 Н1 та Н2:

Entropy MonoGram without spaces: 4.459495560105348

Entropy MonoGram with spaces: 4.376343121574429

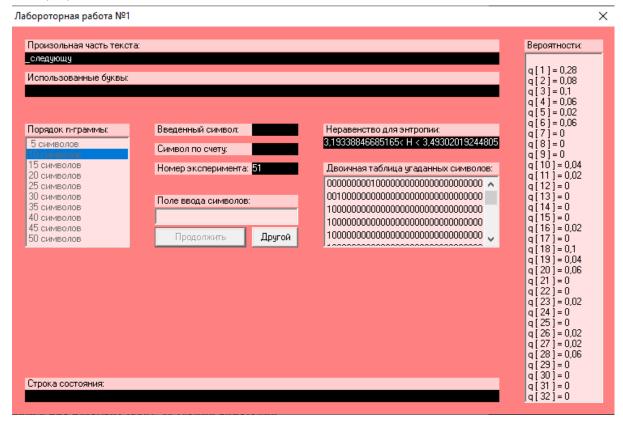
Entropy BiGram without spaces and crossing: 4.138908270465847

Entropy BiGram without spaces and with crossing: 4.138711256669181

Entropy BiGram with spaces and without crossing: 3.966535460585095

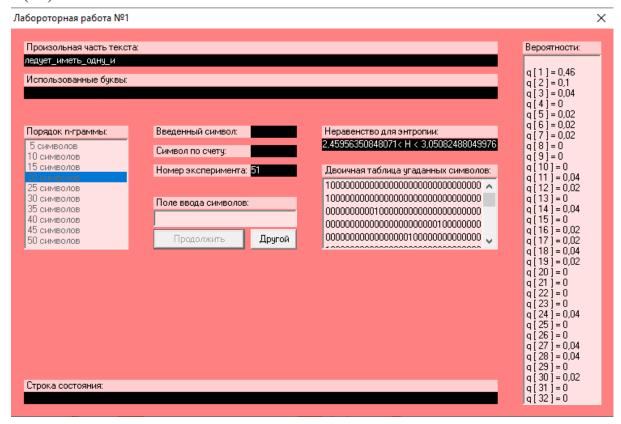
Entropy BiGram with spaces and crossing: 3.966562055656355

2. H(10)



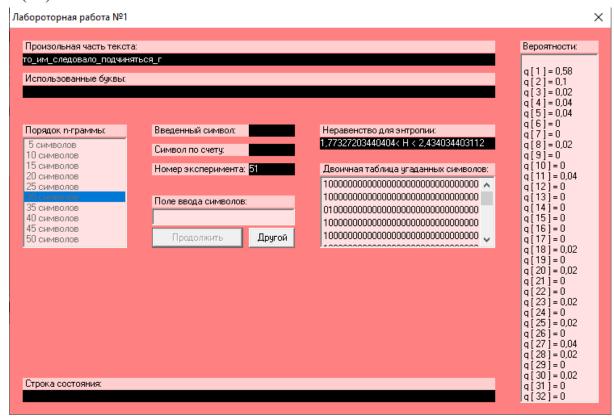
R = 0.33137

H(20)



R = 0.44897

H(30)



R = 0.57928

Висновок

Виконую цю лабораторну роботу, ми дізналися про існування російської мови. Також виявилось, що у творі "Війна і мир" майже половина символів - це розділові знаки, латиниця та цифри. Тим самим з трохи більше 5 мільйонів символів оригінального тексту, після обробки, у нас залишилось приблизно 2.3 мільйони символів кирилиці.