Міністерство освіти і науки України Національний технічний університет України "Київський політехнічний інститут імені Ігоря Сікорського" Фізико-технічний інститут

Криптографія

Комп'ютерний практикум №1

Експериментальна оцінка ентропії на символ джерела відкритого тексту

Виконали:

Студенти III курсу групи ФБ-95 Пашинський М. О. Бурчак Б. Ю. **Перевірила:** Селюх К. І.

Мета роботи: Засвоєння понять ентропії на символ джерела та його надлишковості, вивчення та порівняння різних моделей джерела відкритого тексту для наближеноговизначення ентропії, набуття практичних навичок щодо оцінки ентропії на символ джерела.

Завдання:

- 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 всіпробіли.
- 2. За допомогою програми CoolPinkProgram оцінити значення (10) H, (20) H, (30) H.
- 3. Використовуючи отримані значення ентропії, оцінити надлишковість російської мови в різних моделях джерела.

Хід роботи:

За основу тексту було взято науково-фантастичний роман «Дюна». Програма була розроблена на мові Руthon, вона фільтрує текст від усього лишнього і підраховує частоти монограм та біграм. Інформація записується в окремі файли. За допомогою отриманих даних було підраховано значення ентропій і надлишковості за даними формулами:

Ентропія(**H**):

$$-0.5 \times \sum_{i=0}^{n} (x(i) \times \log_{2}(x(i)))$$

Надлишковість(\mathbf{R}):

$$1-\left(\frac{h}{\log_2(31)}\right)$$

Монограми і їх частоти:

Символи з пробілом	Символи без пробілу
a> 0.06932653754842562	a> 0.08179723672781458
б> 0.014045151669829198	б> 0.016571642499996302
в> 0.036196250027695664	в> 0.04270735762775485
г> 0.015252473693297566	г> 0.01799614181660163
д> 0.02611210118763197	д> 0.030809236950762535
e> 0.07135448738770714	e> 0.08418998125165546
ж> 0.009208338715698657	ж> 0.010864766775230482
3> 0.014073578919689257	з> 0.01660518334263105
и> 0.058503280212000394	и> 0.06902705414231873
й> 0.00925516006840934	й> 0.01092001051604066
к> 0.02652680459735518	к> 0.03129853865508127
л> 0.043232412415852205	л> 0.05100920942754167
м> 0.026674793515744307	м> 0.031473148335856294
н> 0.05660994176176209	н> 0.06679313537330711
o> 0.09319163185373679	o> 0.10995526736737433
п> 0.02546663539669182	п> 0.030047662523879352
p> 0.0376268095721239	p> 0.04439525120858014
c> 0.04626911157737375	c> 0.05459216061723043
т> 0.05319240106167417	т> 0.06276083554184984
y> 0.023863004066350873	y> 0.02815556440113072
ф> 0.0017980235536487002	ф> 0.0021214582966480365
x> 0.007820001998268445	x> 0.009226691210671512
ц> 0.003035863065937424	ц> 0.0035819646943171454
ч> 0.011567800454083512	ч> 0.013648656713914863
ш> 0.006299144130755316	ш> 0.007432256129712303
щ> 0.0034004007406134686	щ> 0.0040120766763392535
ы> 0.016297593173446773	ы> 0.019229261031114556
ь> 0.01554803348228376	ь> 0.018344867930465886
э> 0.0032310913848292966	э> 0.00381231136358816
ю> 0.005216818397112795	ю> 0.006155237871162718
я> 0.017345638988140403	я> 0.020465832979428118
> 0.15245868538183022	

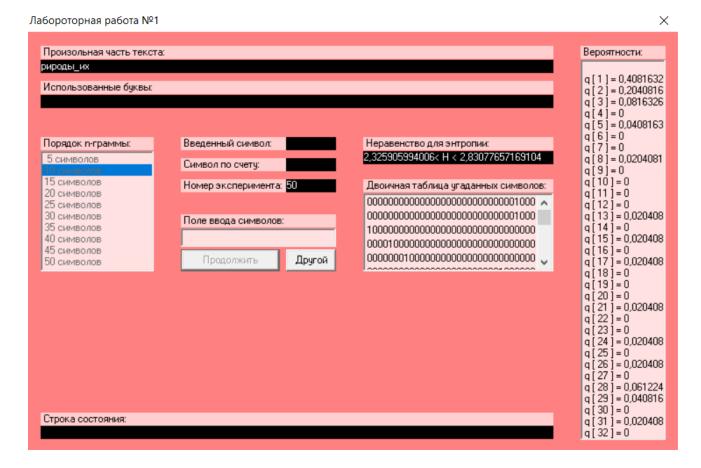
Перші 10 бірам і їх частот (за найбільшою частотою появи):

Частоти з кроком	Частоти з кроком	Частоти з кроком	Частоти з кроком
1 без пробілів	2 без пробілів	1 з пробілами	2 з пробілами
то: 0.015386	то: 0,015386	o_: 0,019159	o_: 0,019159
ст: 0,013524	ст: 0,013524	$\Pi_{-}:0,017552$	$\Pi_{-}: 0.017552$
на: 0,012524	на: 0,012524	a_: 0,016985	a_: 0,016985
по: 0,011953	по : 0,011954	e_: 0,015536	e_: 0,015536
но: 0,011815	но : 0,011815	и_: 0,015524	и_: 0,015524
ен: 0,011453	ен: 0,011453	c_: 0,015309	c_: 0,015309
ни : 0,010609	ни : 0,010609	в_: 0,013683	в_: 0,013683
ов: 0,010382	ов: 0,010382	н_: 0,01337	н_::0,01337
ал: 0,010292	ал: 0,010292	то: 0,012565	то: 0,012565
не: 0,010067	не: 0,010067	o_: 0,01178	o_: 0,01178

	Ентропія:	Надлишковість:
Монограми без пробілу.	2,229394	0,549999
Монограми з пробілом.	2,197483	0,560503
Біграми без перетину та без пробілів.	4,14984	0,162359
Біграми без перетину з пробілами.	4,003918	0,199216
Біграми з перетином та без пробілів.	4,149933	0,16234
Біграми з перетином та пробілами.	4,00399	0,199202

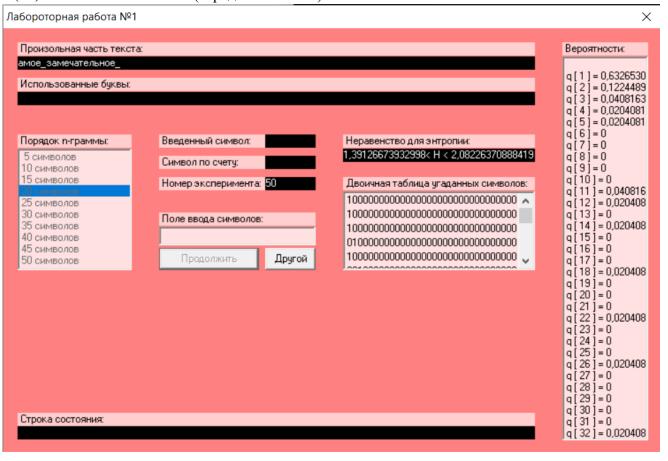
Ентропія для 10 символів:

H(10) = 2.57834128284852 (середнє значення)



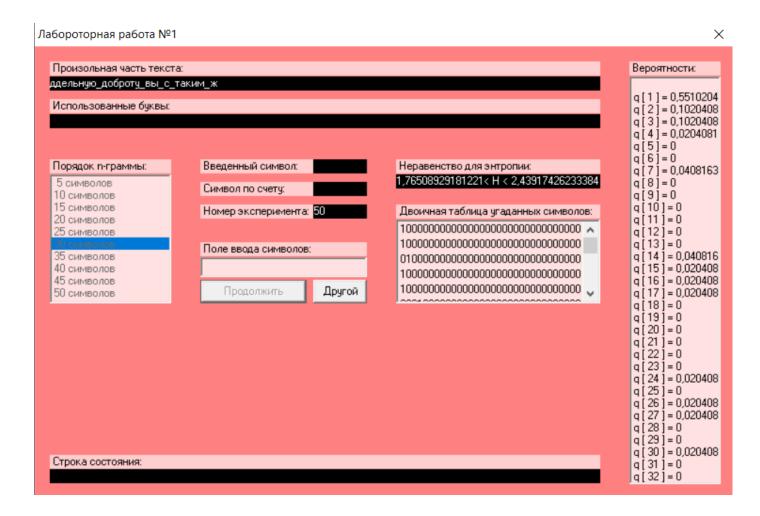
Ентропія для 20 символів:

H(20) = 1.736765224107085 (середнє значення)



Ентропія для 30 символів:

H(30) = 2.102131777073025 (середнє значення)



	Ентропія:	Надлишковість:
H(10)	2.57834128284852	0.484332
H(20)	1.736765224107085	0.652647
H(30)	2.102131777073025	0.579574

Висновок:

В ході цієї лабораторної роботи ми засвоїли поняття ентропії та надлишковості, навчилися знаходити частоти біграм та монограм в тексті. А також вдосконалили свої практичні навички щодо мови програмування Python.