

Київський національний університет імені Тараса Шевченка
Факультет радіофізики, електроніки та комп'ютерних систем

Лабораторна робота №1
з курсу «Комп'ютерні системи»
студента 3 курсу
групи СА-КІ
Колоскова Нікіти

Київ
2019

1. Дослідження кількості інформації в тексті

Обрані тексти:

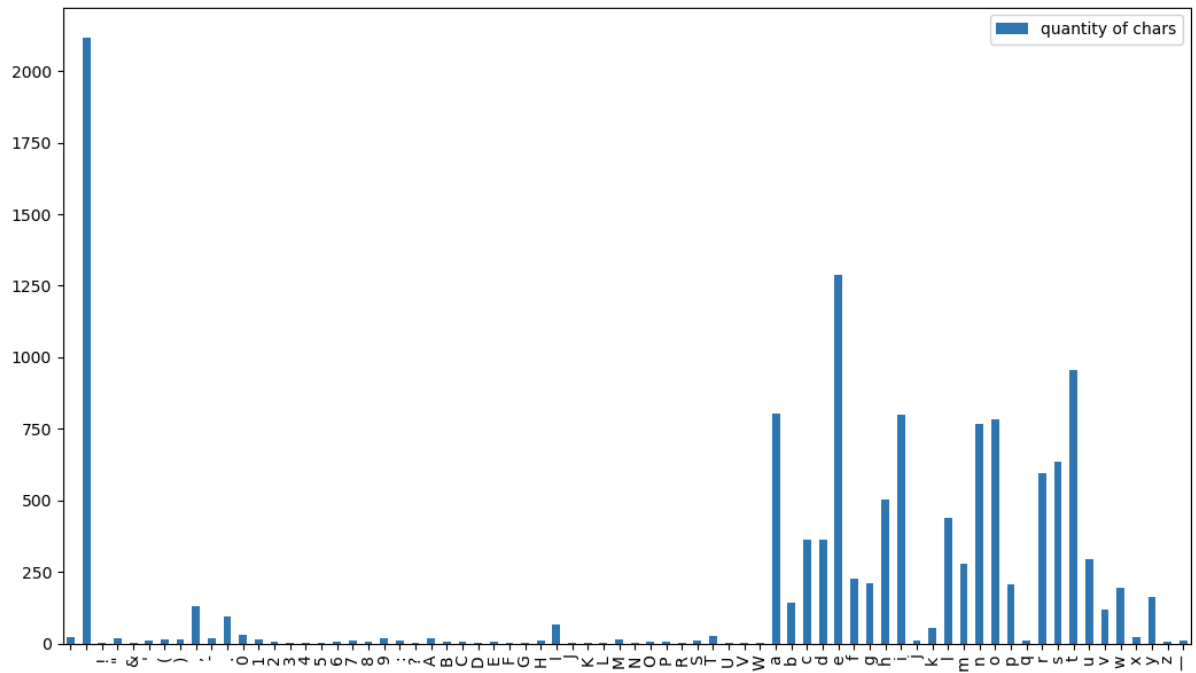
1. Ray Kurzweil «The Singularity is Near» (prologue)
2. Structured Procreation
3. 3D Reconstruction of Human Face Based on Single or Several Images (up to section 2.4)

Посилання на GitHub репозиторій: <https://github.com/Carterochka/cs>

1. Ray Kurzweil «The Singularity is Near» (prologue)

```
Chars amounts in "The Singularity is Near" Prologue
: 24
: 2117
!: 1
": 19
&: 1
': 10
(: 13
): 13
,: 129
-: 20
.: 96
0: 29
1: 16
2: 8
3: 4
4: 3
5: 2
6: 8
7: 9
8: 5
9: 17
[: 11
?: 1
A: 18
B: 6
C: 5
D: 4
E: 5
F: 2
G: 2
H: 10
I: 68
J: 4
K: 2
L: 1
M: 16
N: 1
O: 7
P: 7
R: 4
S: 12
T: 28
U: 1
V: 1
W: 4
a: 802
b: 141
c: 362
d: 364
e: 1288
f: 225
g: 209
h: 502
i: 799
j: 10
k: 54
l: 441
m: 278
n: 768
o: 782
p: 205
q: 12
r: 595
s: 635
t: 956
u: 294
v: 117
w: 193
x: 21
y: 164
z: 8
-: 11
```

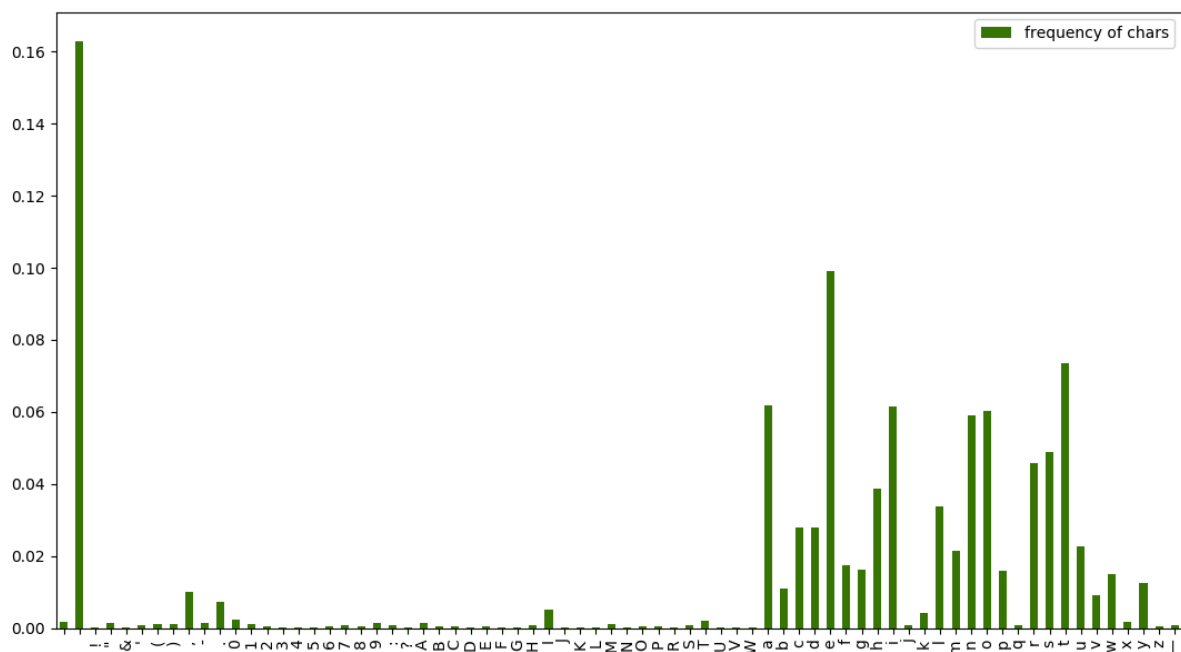
Кількість символів у тексті



Гістограма кількості символів у тексті

Chars frequency in "The Singularity is Near" Prologue	
U:	0.000077
V:	0.000077
W:	0.000308
a:	0.061692
b:	0.010846
c:	0.027846
d:	0.028000
e:	0.099077
f:	0.017308
g:	0.016077
h:	0.038615
i:	0.061462
j:	0.000769
k:	0.004154
l:	0.033923
m:	0.021385
n:	0.059077
o:	0.060154
p:	0.015769
q:	0.000923
r:	0.045769
s:	0.048846
t:	0.073538
u:	0.022615
v:	0.009000
w:	0.014846
x:	0.001615
y:	0.012615
z:	0.000615
-:	0.000846
:	0.001846
:	0.162846
!:	0.000077
":	0.001462
&:	0.000077
':	0.000769
(:	0.001000
):	0.001000
,:	0.009923
-:	0.001538
.:	0.007385
0:	0.002231
1:	0.001231
2:	0.000615
3:	0.000308
4:	0.000231
5:	0.000154
6:	0.000615
7:	0.000692
8:	0.000385
9:	0.001308
::	0.000846
?:	0.000077
A:	0.001385
B:	0.000462
C:	0.000385
D:	0.000308
E:	0.000385
F:	0.000154
G:	0.000154
H:	0.000769
I:	0.005231
J:	0.000308
K:	0.000154
L:	0.000077
M:	0.001231
N:	0.000077
O:	0.000538
P:	0.000538
R:	0.000308
S:	0.000923
T:	0.002154

Частота символів у тексті



Гістограма розподілу символів у тексті

Гістограми кількості та частот виглядають однаково, тому що по суті це ті самі ж значення, але у другому випадку пронормовані, щоб сума по усім символам дорівнювала одиниці.

```
Entropy: 4.412461
Quantity of information: 7170.248740
```

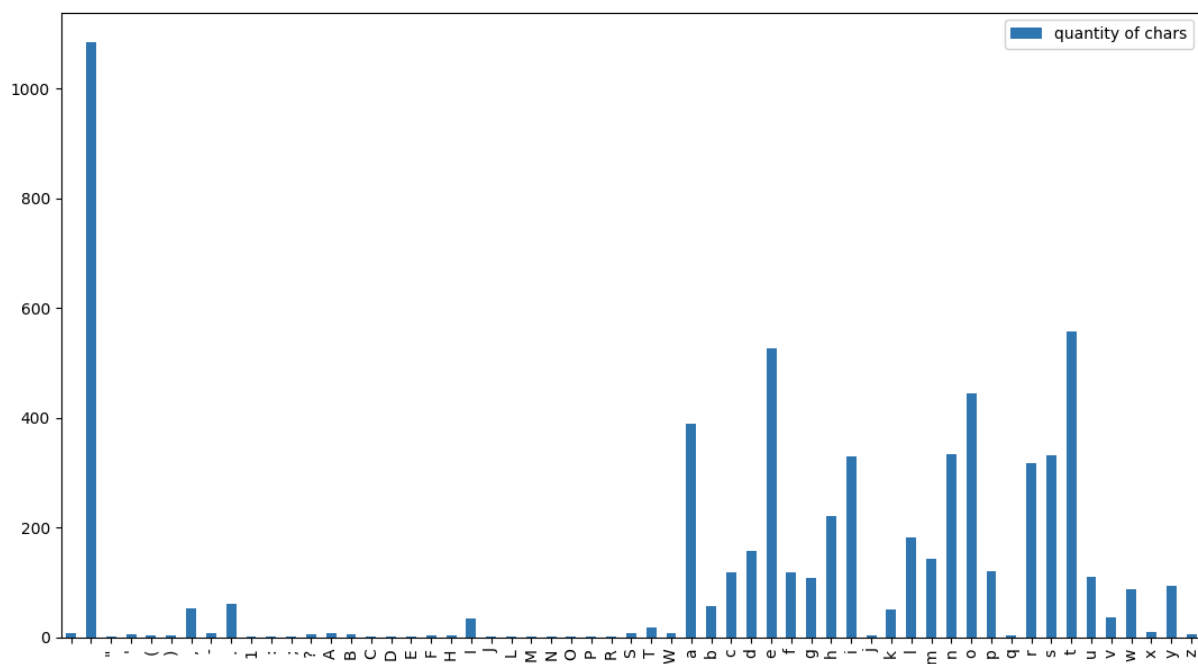
Ентропія та кількість інформації в оригінальному тексті

(цьо ще не все, див. далі)

2. Structured Procrastination

Chars amounts in "Structured Procrastination"		W: 7
:	8	a: 390
:	1085	b: 56
"	2	c: 119
'	6	d: 158
(3	e: 527
)	3	f: 118
,	53	g: 109
-	7	h: 222
.	60	i: 329
1	1	j: 4
:	1	k: 50
;	2	l: 181
?	6	m: 144
A	8	n: 334
B	5	o: 444
C	1	p: 120
D	1	q: 3
E	2	r: 317
F	4	s: 331
H	3	t: 557
I	34	u: 110
J	1	v: 36
L	1	w: 87
M	1	x: 9
N	2	y: 94
O	2	z: 5
P	2	
R	2	
S	7	
T	18	

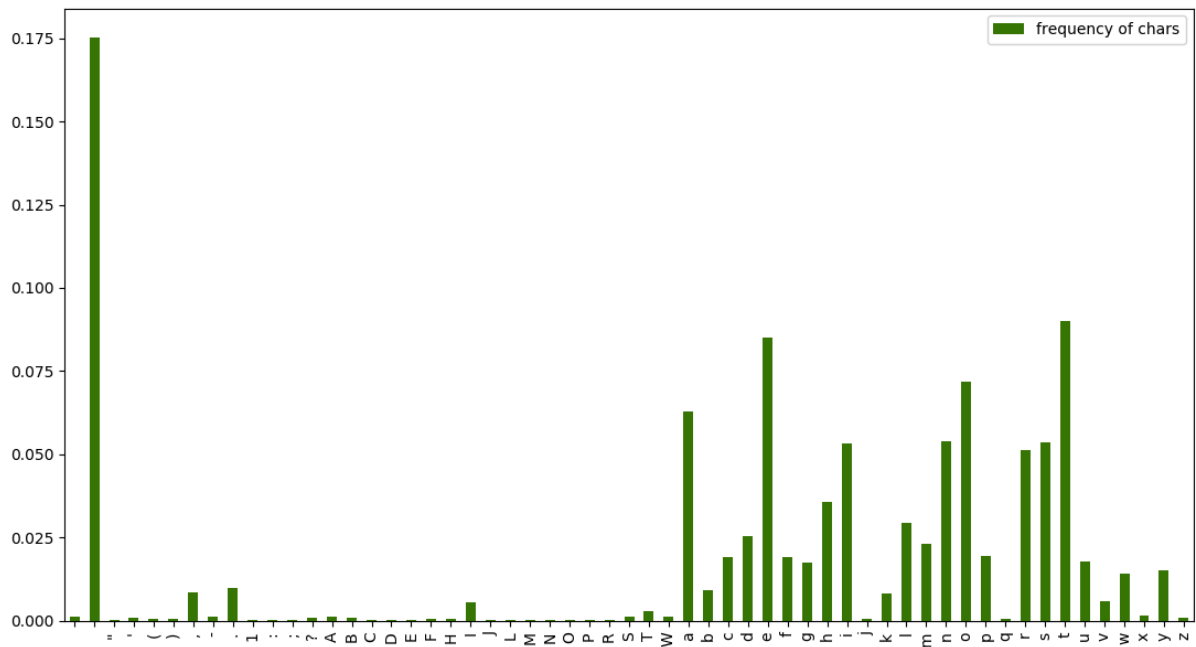
Кількість символів у тексті



Гістограма кількості символів у тексті

Chars frequency in "Structured Procrastination"		W: 0.001130
:	0.001292	a: 0.062984
:	0.175226	b: 0.009044
"	0.000323	c: 0.019218
'	0.000969	d: 0.025517
(0.000484	e: 0.085110
)	0.000484	f: 0.019057
,	0.008559	g: 0.017603
-	0.001130	h: 0.035853
.	0.009690	i: 0.053133
1	0.000161	j: 0.000646
::	0.000161	k: 0.008075
;	0.000323	l: 0.029231
?	0.000969	m: 0.023256
A	0.001292	n: 0.053941
B	0.000807	o: 0.071705
C	0.000161	p: 0.019380
D	0.000161	q: 0.000484
E	0.000323	r: 0.051195
F	0.000646	s: 0.053456
H	0.000484	t: 0.089955
I	0.005491	u: 0.017765
J	0.000161	v: 0.005814
L	0.000161	w: 0.014050
M	0.000161	x: 0.001453
N	0.000323	y: 0.015181
O	0.000323	z: 0.000807
P	0.000323	
R	0.000323	
S	0.001130	
T	0.002907	

Частота символів у тексті



Гістограма розподілу символів у тексті

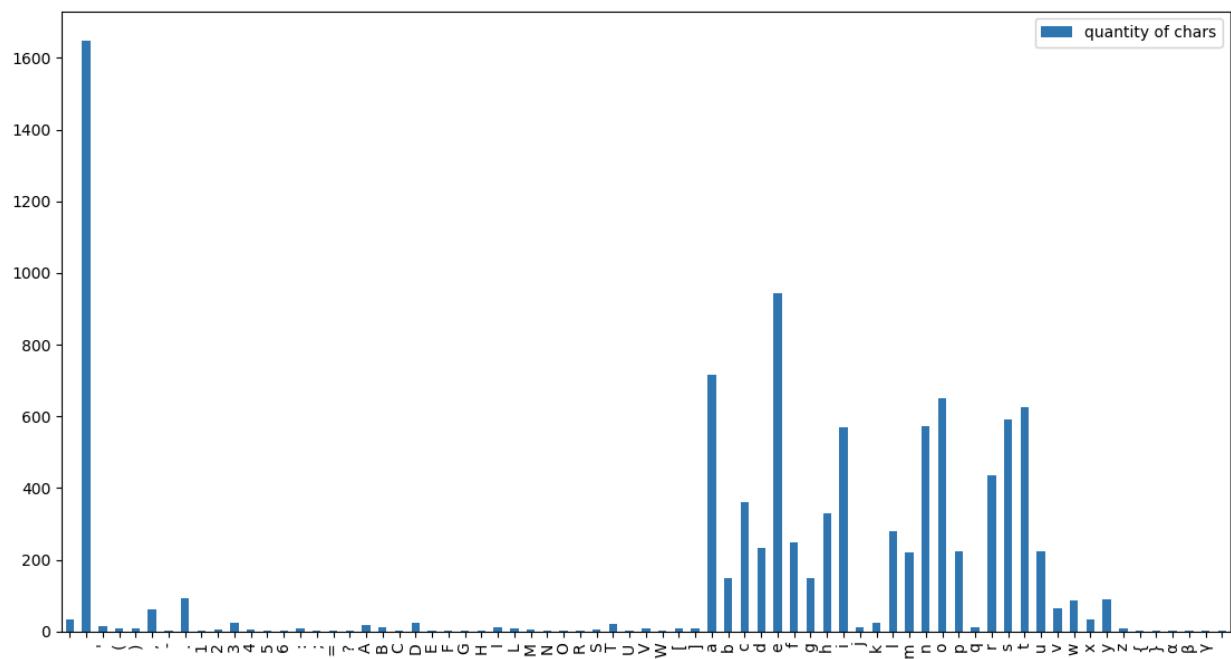
```
Entropy: 4.305188
Quantity of information: 3332.215782
```

Ентропія та кількість інформації в оригінальному тексті

3. 3D Reconstruction of Human Face Based on Single or Several Images (up to section 2.4)

```
Chars amounts in "3D Reconstruction of Human Face..."
: 32
: 1648
': 14
(: 9
): 9
,: 62
-: 3
.: 91
1: 2
2: 4
3: 24
4: 4
5: 1
6: 2
:: 9
;: 2
=: 1
?: 1
A: 19
B: 12
C: 3
D: 24
E: 1
F: 1
G: 1
H: 2
I: 12
L: 9
M: 5
N: 2
O: 1
R: 1
S: 6
T: 22
U: 1
V: 7
W: 3
[: 7
]: 7
a: 715
b: 149
c: 362
d: 232
e: 945
f: 248
g: 148
h: 329
i: 569
j: 10
k: 24
l: 278
m: 219
n: 574
o: 650
p: 223
q: 13
r: 435
s: 590
t: 626
u: 225
v: 65
w: 86
x: 32
y: 90
z: 8
{: 1
}: 1
α: 2
β: 1
γ: 2
: 3
```

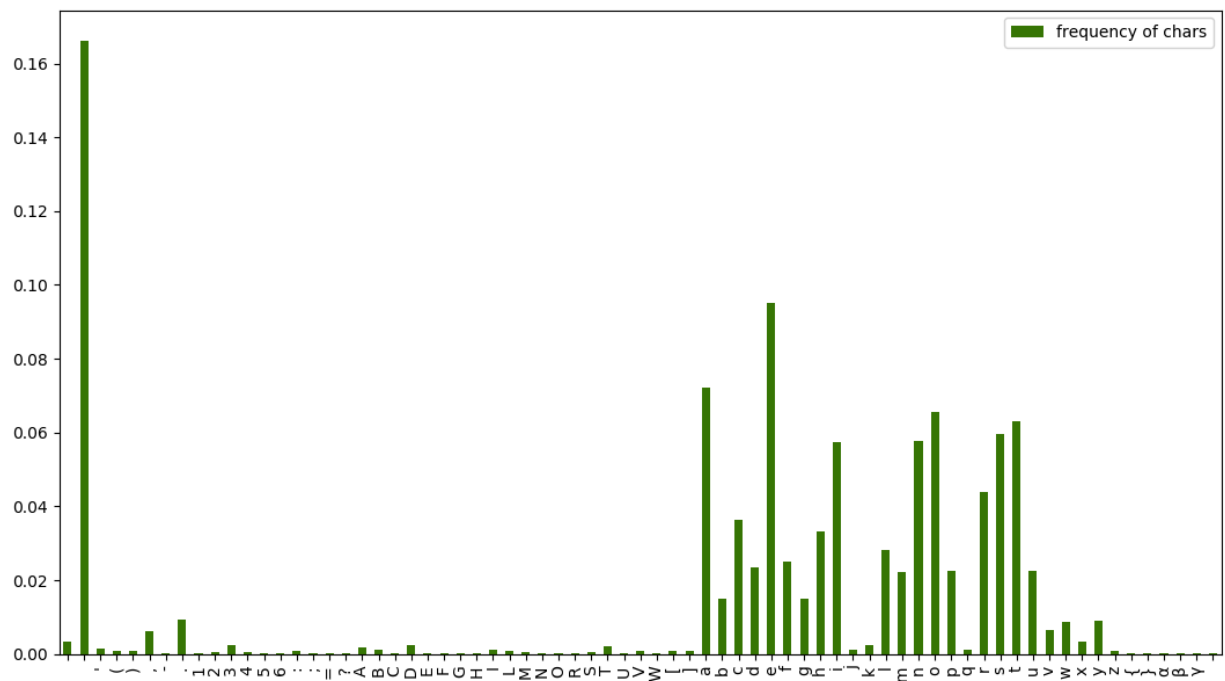
Кількість символів у тексті



Гістограма кількості символів

Chars frequency in "3D Reconstruction of Human Face..."		U: 0.000101
		V: 0.000706
		W: 0.000302
:	0.003226	[: 0.000706
:	0.166146]: 0.000706
':	0.001411	a: 0.072084
(:	0.000907	b: 0.015022
):	0.000907	c: 0.036496
,:	0.006251	d: 0.023389
-:	0.000302	e: 0.095272
.:	0.009174	f: 0.025003
1:	0.000202	g: 0.014921
2:	0.000403	h: 0.033169
3:	0.002420	i: 0.057365
4:	0.000403	j: 0.001008
5:	0.000101	k: 0.002420
6:	0.000202	l: 0.028027
::	0.000907	m: 0.022079
;	0.000202	n: 0.057869
=:	0.000101	o: 0.065531
?:	0.000101	p: 0.022482
A:	0.001916	q: 0.001311
B:	0.001210	r: 0.043855
C:	0.000302	s: 0.059482
D:	0.002420	t: 0.063111
E:	0.000101	u: 0.022684
F:	0.000101	v: 0.006553
G:	0.000101	w: 0.008670
H:	0.000202	x: 0.003226
I:	0.001210	y: 0.009073
L:	0.000907	z: 0.000807
M:	0.000504	{: 0.000101
N:	0.000202	}: 0.000101
O:	0.000101	a: 0.000202
R:	0.000101	β: 0.000101
S:	0.000605	γ: 0.000202
T:	0.002218	: 0.000302

Частота символів у тексті

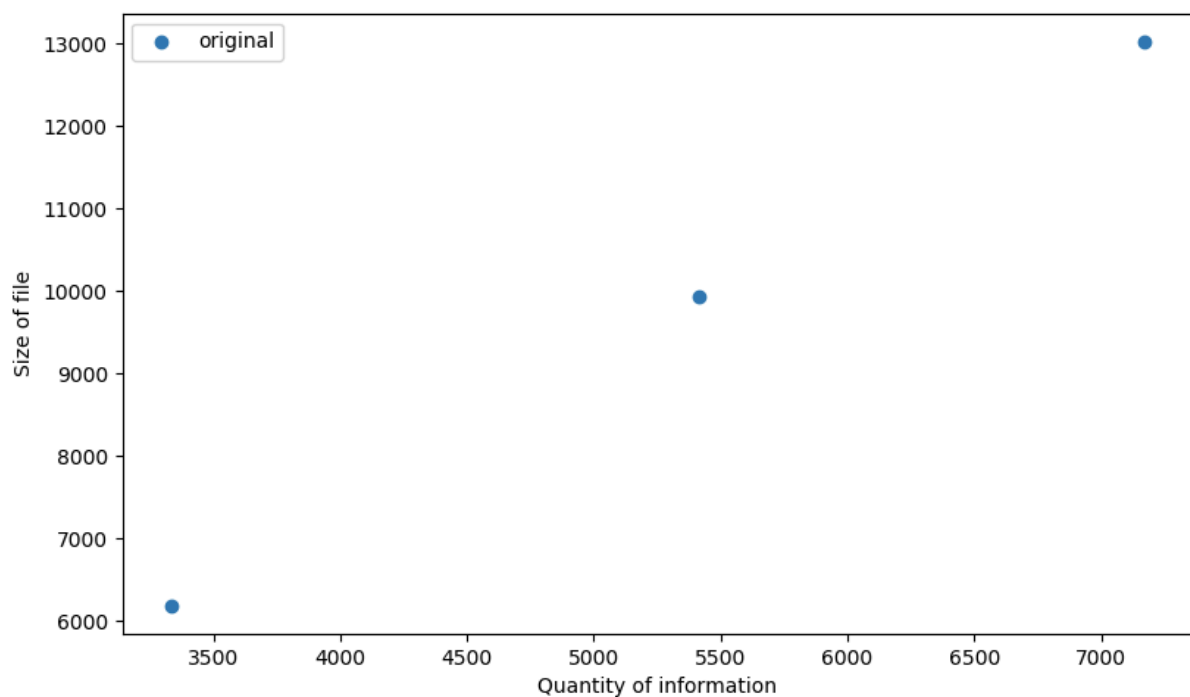


Гістограма розподілу символів у тексті

```
Entropy: 4.366933
Quantity of information: 5414.451134
```

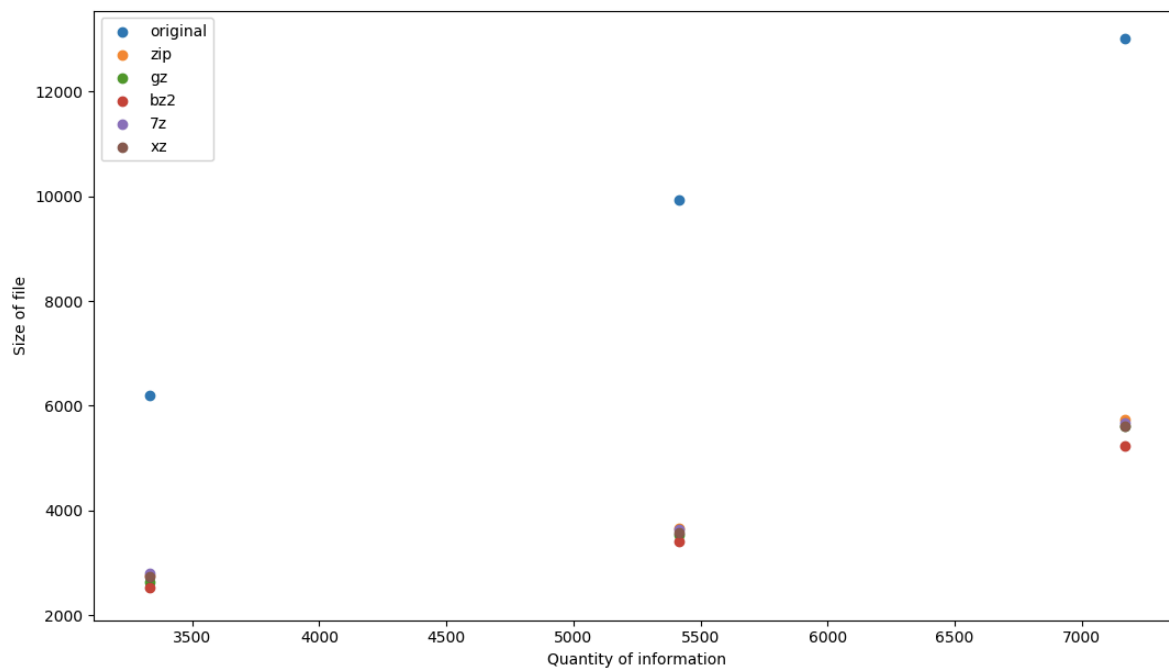
Ентропія та кількість інформації в оригінальному тексті

Для того, щоб можна було наглядно побачити тенденцію зміни розміру файлу в залежності від кількості інформації, тексти були вибрані таким чином, щоб їх довжина була різною.

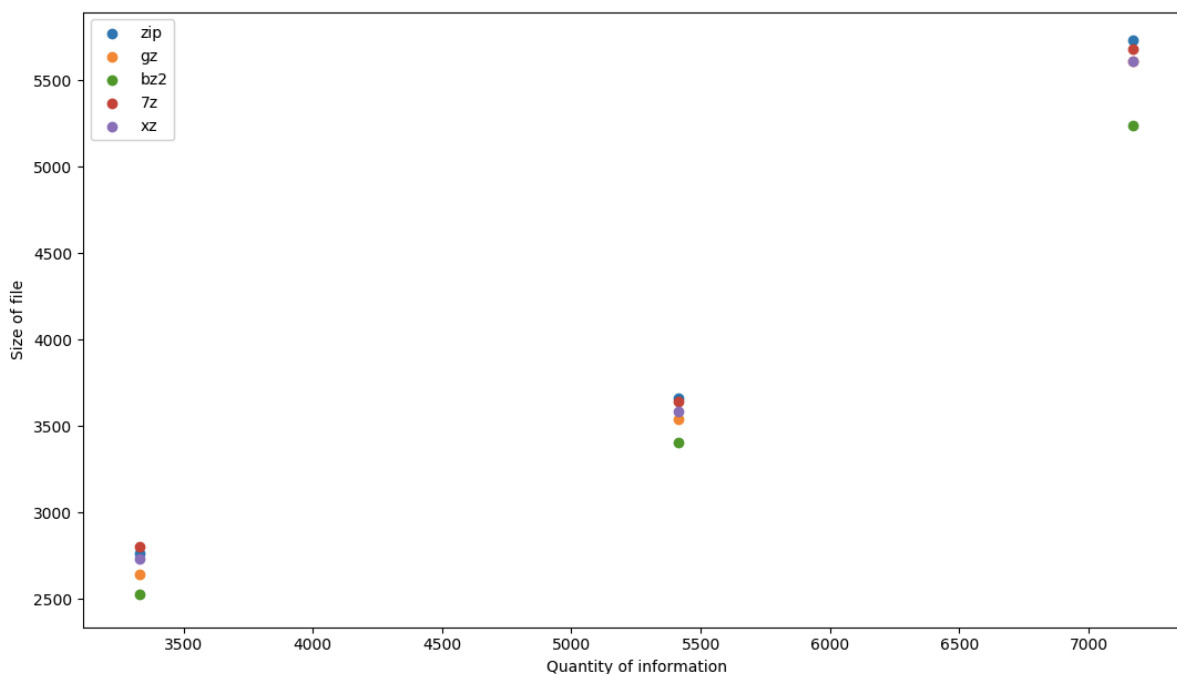


Залежність розміру файла (у байтах) від кількості інформації (у байтах)

Обрані алгоритми стиснення файлів: zip, gzip, bzip2, 7-zip, xz



Залежність розміру файла від кількості інформації в порівнянні з оригінальним файлом



та без порівняння з оригінальним файлом

З побаченого на останньому графіку можна зробити висновок, що:

1. ефективність різних алгоритмів змінюється із зміною розміру оригінального файлу (так, для малих розмірів оригінального файлу zip ефективніше за 7-zip, але зі збільшенням розміру 7-zip реалізує краще стискання. Аналогічна ситуація із gzip та xz - для малих розмірів xz стискає краще, але для найбільшого з наведених файлів їх результативність стала майже однаковою)

2. bzip2 усіх випереджає взагалі без шансів, і ця тенденція збільшується зі збільшенням розміру файлу

bzip2 - НАШ ПЕРЕМОЖЕЦЬ

3. для того, щоб зрозуміти ефективніше кореляцію між усіма алгоритмами, потрібно провести порівняння для ще декількох більших значеннях розміру файлу (більше даних - більш коректні висновки)

Загалом, алгоритм bzip2 і повинен стискати ефективніше за zip та gzip. Це зумовлено тим, що цей алгоритм використовує для стискання інший математичний апарат. Але за краще стискання доводиться платити більшим часом очікування та навантаженням на процесор.

Згідно інформації з вікіпедії, bzip2 повинен поступатися в ефективності алгоритму 7-zip, проте з отриманих результатів цього побачити не можна. Можливо, для більших розмірів файлів це дійсно було б так.

2. Дослідження способів кодування інформації на прикладі Base64

Перевірку роботи програми-кодувальника робив за допомогою онлайн-сервіса для кодування в base64 <https://www.motobit.com/util/base64-decoder-encoder.asp> та сервісу для порівняння двох текстів <http://text.num2word.ru>

Приклад роботи програми для кодування наведений в додатку А.

Результат порівняння:

(скріншоти фрагменту виводу порівняння на наступній сторінці)

Второй текст :

Поменять текст местами

Стереть

Настройка результата отображения сравнения:

- ## Сравнить

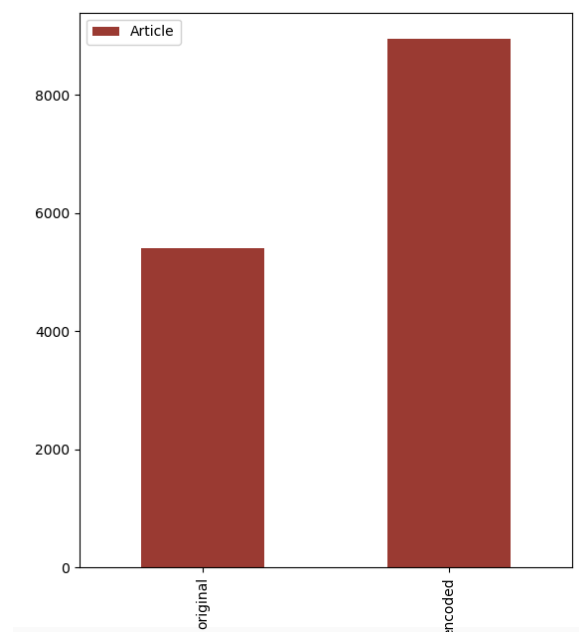
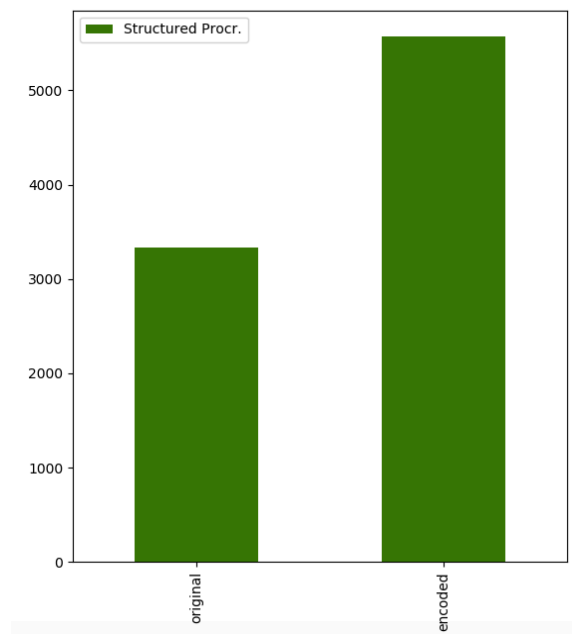
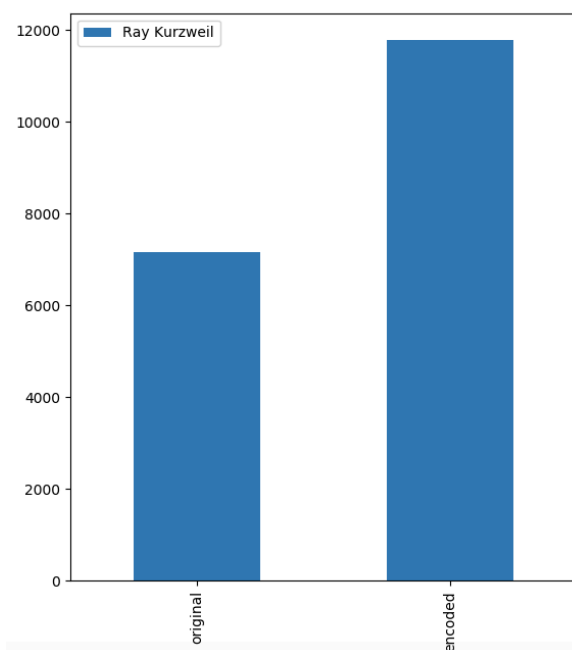
Результат сравнения:

[illegible]

Як можна побачити, єдине, чим відрізняються результати - це переходи на нову строку. Це зумовлено тим, що вивід написаної мною програми - сплошний текст, а результат на сайті розбитий на строки довжиною не більше 76 символів. Отже, програма кодує вірно.

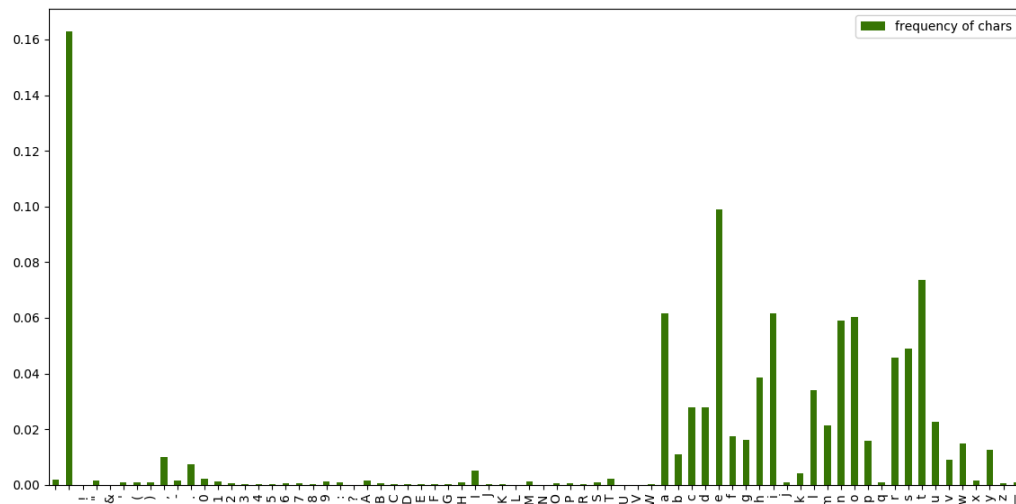
Кількість інформації в оригінальному та закодованому файлах:

	Original	Encoded
Ray Kurzweil	7170.248740	11778.310417
Structured Procr.	3332.215782	5567.081340
Article	5414.451134	8947.269930

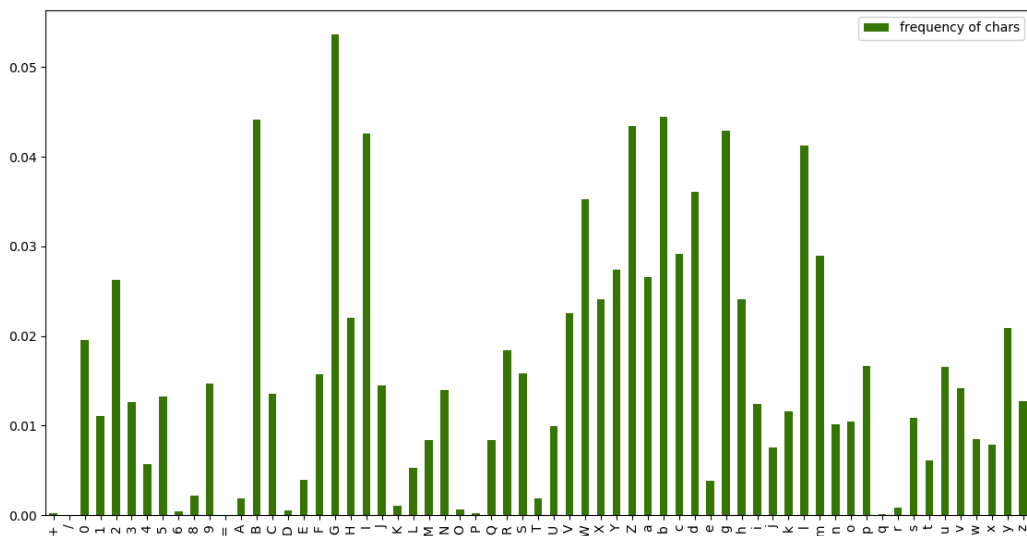


Отже, у кодуванні base64 кількість інформації збільшується, на прикладі текстів, що я використовував - приблизно в 1.65 разів. Це пов'язано з тим, що кількість символів, що міститься у закодованому файлі, більша, що впливає з опису кодування - грубо кажучи, кожним 8 оригінальним бітам відповідають 6 закодованих. Таким чином, збільшення кількості символів - 33%. Решта прирісту може бути зумовлена тим, що кількість символів закодованого алгоритму більша за кількість символів оригінального.

Порівняння розподілу символів вихідного та закодованого алфавітів (на прикладі витягу з книги Рея Курцвела, для інших текстів ситуація аналогічна):



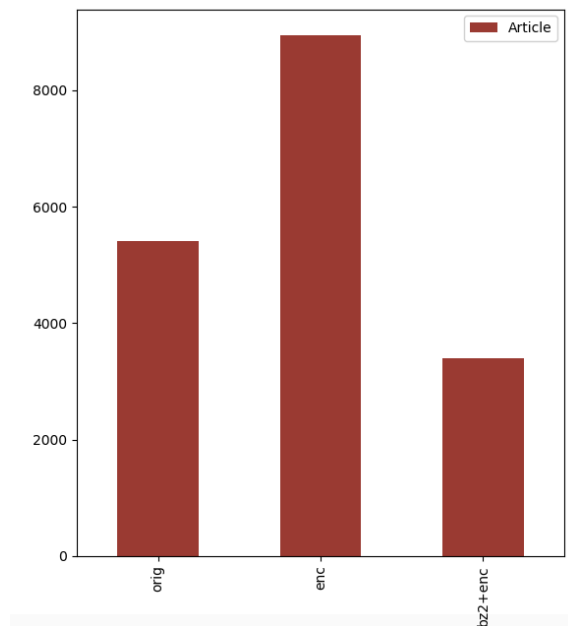
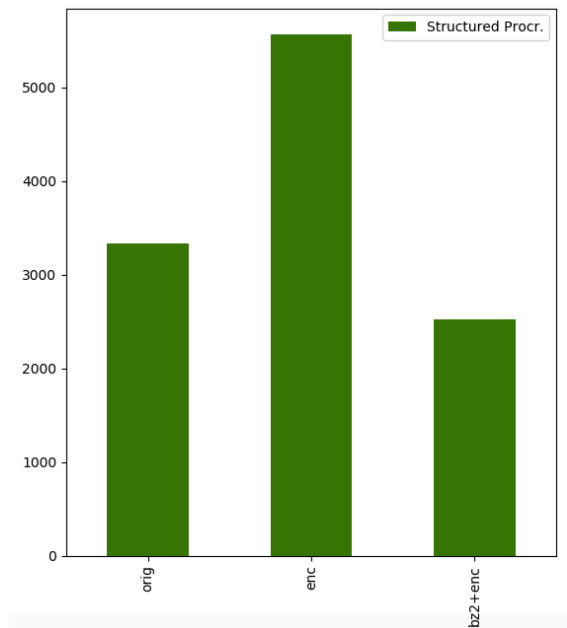
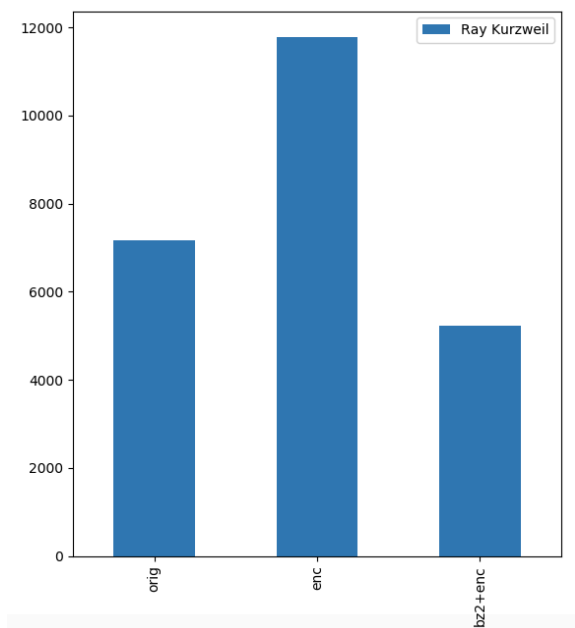
Розподіл символів алфавіту (оригінальний текст)

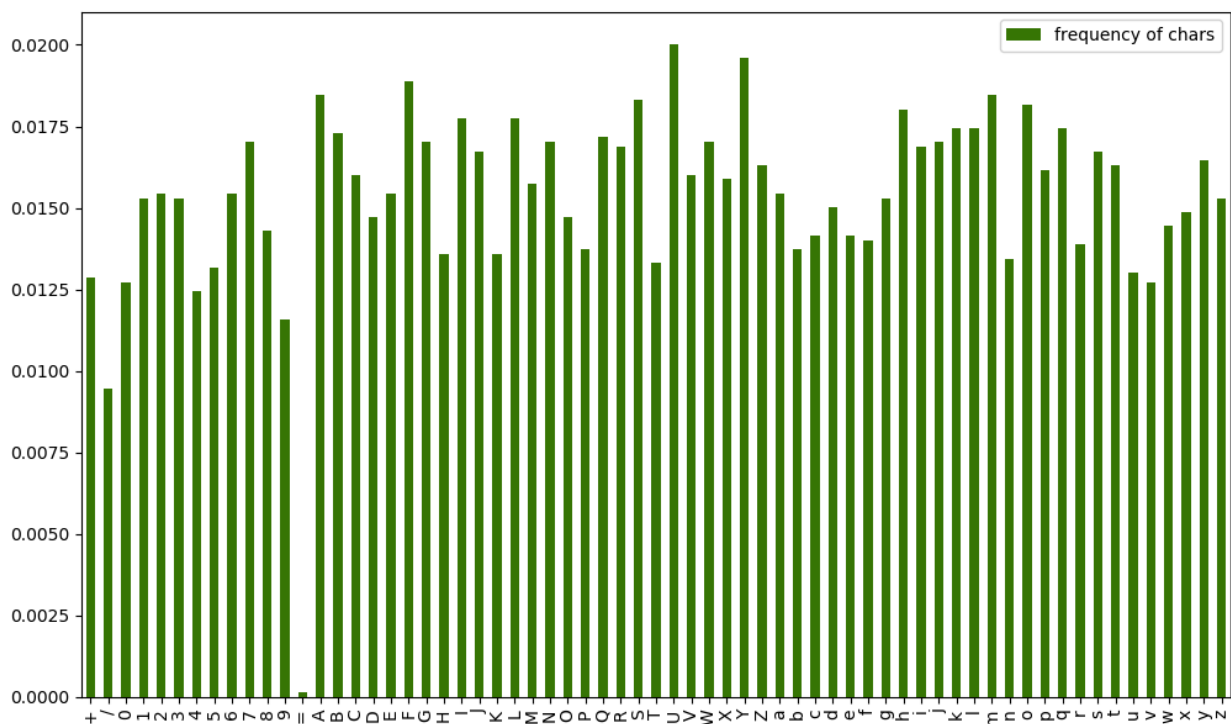


Закодований у base64 текст

Порівняння кількості інформації в оригінальному, закодованому та закодованому стисненому файлах:

	Original	Encoded	Compressed+Encoded
Ray Kurzweil	7170.248740	11778.310417	5234.775139
Structured Procr.	3332.215782	5567.081340	2522.673411
Article	5414.451134	8947.269930	3391.485013





Розподіл символів алфавіту (стиснений та закодований файл)

Отже, кількість інформації у стисненому та закодованому файлі становить ~ 0.7 від кількості інформації оригінального файлу. Зрозуміло, що кількість інформації повинна бути меншою \rightarrow bzr2 стискає розмір файлу більше ніж у 2,5 рази (судячи з графіків, отриманих у ч. 1 цієї лаби). Тоді виникає питання, чому ж тоді кількість інформації менша усього лише на ~ 0.3 .

Можна побачити, що в закодованому після стиснення файлі усі символи алфавіту майже рівноймовірні (символ “=” зустрічається лише 1 раз наприкінці файлу, і означає, що кількість біт до кодування не була кратною 24).

При нерівноймовірному алфавіті $H = \sum_i p_i \log_2 \frac{1}{p_i} \leq \log_2 n$

При рівноймовірному алфавіті $H = \log_2 n$

Таким чином, ентропія майже рівноймовірного алфавіту ближче до максимального значення, що також збільшило свій внесок у значення кількості інформації.

Висновки: за допомогою інструментів програмування можна обрахувати середню ентропію та кількість інформації в тексті та файлі. Ці величини пов’язані між собою наступним чином: кількість інформації дорівнює добутку середньої ентропії тексту на його кількість символів. При різних кодуваннях значення кількості інформації тексту буде відрізнятись; те саме можна сказати і про кодування після стискання.

Приклад роботи програми-кодувальника

[illegible]