

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

Лабораторна робота № 1

**Тема: « Дослідження кількості інформації при різних варіантах
кодування »**

Роботу виконав
студент 3 курсу
КІ - СА
Гулівець Владислав
Андрійович

Київ 2020

Мета: Дослідити імовірнісні параметри української мови для оцінки кількості інформації текстів. Дослідити вплив різних методів кодування інформації на її кількість.

Теоретичні відомості

Відносна частота появи символу - імовірність появи певного символу в певному місці тексту - відношення числа появи символу в тексті до загальної кількості символів.

Середня ентропія нерівноймовірного алфавіту:

$$H = \sum_{i=1}^m p_i \log_2 \frac{1}{p_i} = - \sum_{i=1}^m p_i \log_2 p_i$$

де m - кількість символів алфавіту, p - імовірність появи символу
Ентропія вимірюється в **БІТАХ** (як представлення кількості можливих варіантів).

Кількість інформації в тексті - середня ентропія вихідного алфавіту помножена на кількість символів тексту. (**HINT:** результат обрахунку для порівняння значення з розміром файлів треба перевести з бітів в байти)

Хід виконання роботи:

Дослідження кількості інформації в тексті

1. Оберіть 3 текстових файла різного тематичного та лінгвістичного спрямування (наприклад, вірш Тараса Шевченка "Мені тринадцятий минало", "Казка про репку" Леся Подерв'янського та специфікацію інтерфейсу PCI)
 - text1.txt – куплет пісні Скрибіна.
 - text2.txt – історія України за 1917 рік.
 - text3.txt – рецензія на фільм «Захар Беркут».
2. Створіть програму (будь-якою зручною для вас мовою), яка в якості вхідних даних приймає текстовий файл, та аналізуючи його вміст:
 - обраховує частоти (імовірності) появи символів в тексті
 - обраховує середню ентропію алфавіту для даного тексту
 - виходячи з ентропії визначає кількість інформації та порівнює її з розмірами файлів
 - виводить на екран значення частот, ентропії та кількості інформації
4. Проведіть стиснення кожного вхідного файлу за допомогою 5 різних алгоритмів стиснення (zip, rar, gzip, bzip2, xz, або будь-які інші на ваш вибір, можна використовувати готові програмні засоби для стиснення).

Analyze of *text3.txt*

```
Ця програма рахує в'дносну частоту появи символу в тексті.
В текст? було знайдено наступні букви:

Цей символ А зустрічається 14 раз з ймовірністю 0,03063457. Ентропія символу = 0,15405194
Цей символ Б зустрічається 3 раз з ймовірністю 0,00656455. Ентропія символу = 0,04760014
Цей символ В зустрічається 11 раз з ймовірністю 0,02407002. Ентропія символу = 0,12941533
Цей символ Г зустрічається 2 раз з ймовірністю 0,00437637. Ентропія символу = 0,03429344
Цей символ ? зустрічається 0 раз з ймовірністю 0,00000000. Ентропія символу = 0,00000000
Цей символ Д зустрічається 3 раз з ймовірністю 0,00656455. Ентропія символу = 0,04760014
Цей символ Е зустрічається 7 раз з ймовірністю 0,01531729. Ентропія символу = 0,09234326
Цей символ Є зустрічається 1 раз з ймовірністю 0,00218818. Ентропія символу = 0,01933490
Цей символ Ж зустрічається 4 раз з ймовірністю 0,00875274. Ентропія символу = 0,05983414
Цей символ З зустрічається 5 раз з ймовірністю 0,01094092. Ентропія символу = 0,07127048
Цей символ И зустрічається 15 раз з ймовірністю 0,03282276. Ентропія символу = 0,16178861
Цей символ ? зустрічається 15 раз з ймовірністю 0,03282276. Ентропія символу = 0,16178861
Цей символ І зустрічається 2 раз з ймовірністю 0,00437637. Ентропія символу = 0,03429344
Цей символ Й зустрічається 1 раз з ймовірністю 0,00218818. Ентропія символу = 0,01933490
Цей символ К зустрічається 9 раз з ймовірністю 0,01969365. Ентропія символу = 0,11158671
Цей символ Л зустрічається 8 раз з ймовірністю 0,01750547. Ентропія символу = 0,10216281
Цей символ М зустрічається 9 раз з ймовірністю 0,01969365. Ентропія символу = 0,11158671
Цей символ Н зустрічається 8 раз з ймовірністю 0,01750547. Ентропія символу = 0,10216281
Цей символ О зустрічається 14 раз з ймовірністю 0,03063457. Ентропія символу = 0,15405194
Цей символ П зустрічається 8 раз з ймовірністю 0,01750547. Ентропія символу = 0,10216281
Цей символ Р зустрічається 6 раз з ймовірністю 0,01312910. Ентропія символу = 0,08207118
Цей символ С зустрічається 10 раз з ймовірністю 0,02188184. Ентропія символу = 0,12065913
Цей символ Т зустрічається 7 раз з ймовірністю 0,01531729. Ентропія символу = 0,09234326
Цей символ У зустрічається 5 раз з ймовірністю 0,01094092. Ентропія символу = 0,07127048
Цей символ Ф зустрічається 0 раз з ймовірністю 0,00000000. Ентропія символу = 0,00000000
Цей символ Х зустрічається 2 раз з ймовірністю 0,00437637. Ентропія символу = 0,03429344
Цей символ Ц зустрічається 1 раз з ймовірністю 0,00218818. Ентропія символу = 0,01933490
Цей символ Ч зустрічається 1 раз з ймовірністю 0,00218818. Ентропія символу = 0,01933490
Цей символ Ш зустрічається 0 раз з ймовірністю 0,00000000. Ентропія символу = 0,00000000
Цей символ Щ зустрічається 0 раз з ймовірністю 0,00000000. Ентропія символу = 0,00000000
Цей символ Ъ зустрічається 4 раз з ймовірністю 0,00875274. Ентропія символу = 0,05983414
Цей символ Ю зустрічається 1 раз з ймовірністю 0,00218818. Ентропія символу = 0,01933490
Цей символ Я зустрічається 8 раз з ймовірністю 0,01750547. Ентропія символу = 0,10216281
Цей символ ? зустрічається 38 раз з ймовірністю 0,08315098. Ентропія символу = 0,29835595
Загальна ентропія: 2,63565820
HINT файлу: 1204,49579898
```

Entropy:

text1.txt – 4.25588104

text2.txt – 2.27625325

text3.txt – 2.63565820

Size of file (bytes):

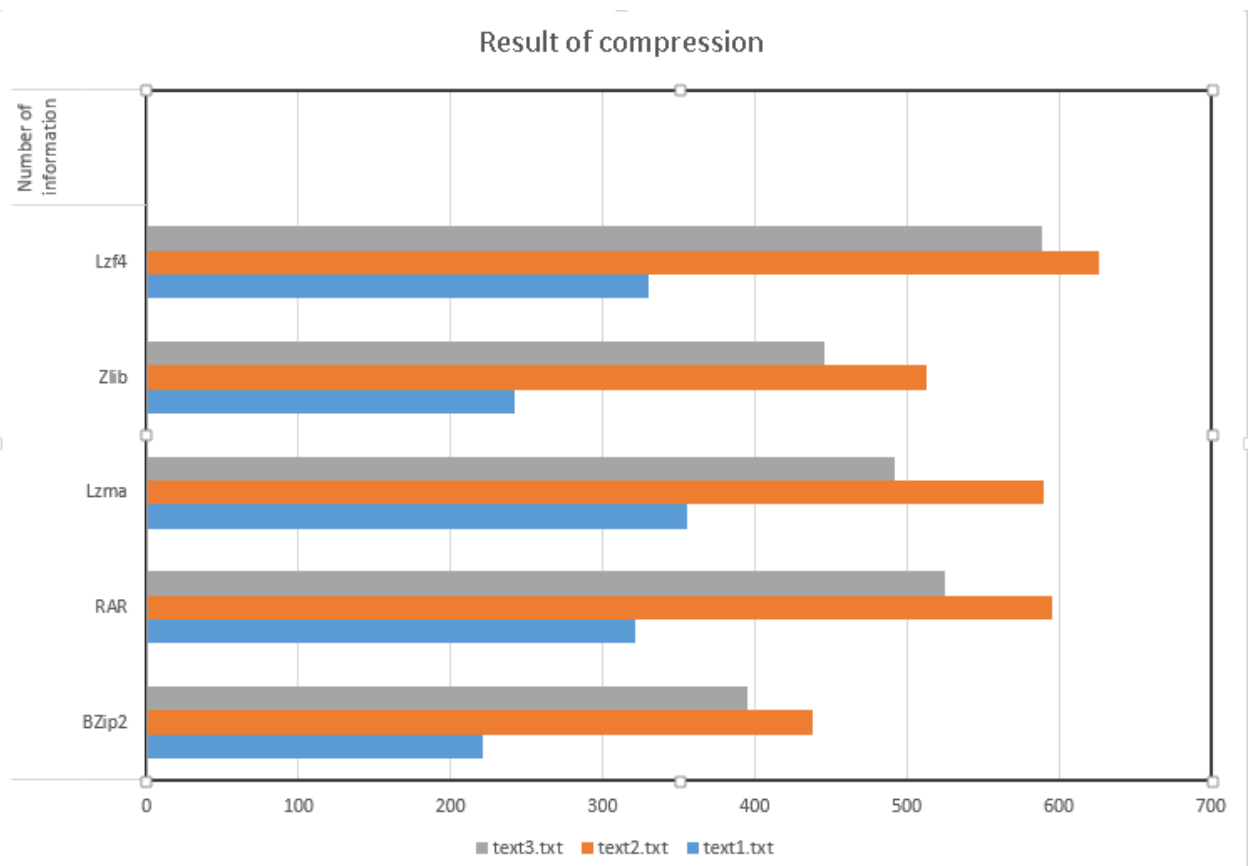
text1.txt – 987bit/8 = 123.375bytes

text2.txt – 1277bit/8 = 159.625bytes

text3.txt – 1204bit/8 = 150.5bytes

Result of compression

Назва файлу	BZip2	RAR	Lzma	Zlib	Lzf4	Number of information
text1.txt	221	321	356	242	330	123.375
text2.txt	438	595	590	513	626	159.625
text3.txt	395	525	492	446	589	150.5



Дослідження способів кодування інформації на прикладі Base64

1. Ознайомтесь зі стандартом [RFC4648](#)
2. Для практичного засвоєння методу кодування, створіть програму, що кодує довільний файл в Base64 (шляхом реалізації алгоритму вручну, а не виклику бібліотечної функції).
Перевірте коректність роботи програми, порівнявши результат з існуючими програмними засобами (наприклад, openssl enc -base64)
3. Закодуйте в Base64 обрані вами текстові файли
Обрахуйте кількість інформації в base64-закодованому варіанті файлу
Порівняйте отримане значення з кількістю інформації вихідного файлу
Зробіть висновки з отриманого результату
4. Закодуйте в Base64 стиснені кращим з алгоритмів текстові файли
Обрахуйте кількість інформації в base64-закодованому варіанті стисненого файлу
Порівняйте отримане значення з кількістю інформації вихідного файлу та base64-закодованого файлу
Зробіть висновки з отриманого результату

Хід виконання роботи:

Base64 compression of *text1.txt*

```
Original string:
Я завжди мр?яв написати п?сню про маму
Але р?зн? поети вс? слова вже сказали
? я не хот?в повторити когось ?з них
Я б?ля своєї мами буду завжди маленьким
? як т?льки покличе приб?жу скоренько
Тому я їй ? написав ц? слова як м?г

Original string length:
232

Encrypted string:
0K8g0LfQsNCy0LbQtNC4INC80YDRltGP0LIg0L3QsNC/0LjRgdCw0YlQuCDQv9Gw0YHQvdGOINC/0YDQvIDQvNCw0LzRgw0K0JD
NGC0Lgg0LrQvtCz0L7RgdGmINGW0Lcg0L3QuNGFDQrQr yDQsdGW0LvRjyDRgdCy0L7RlNGXINC80LDQvNC4INCx0YPQtNGDINC3
LQvtC80Ymg0Y8g0ZfQuSDRliDQvdCw0L/QuNGB0LDQsiDRhtGWINGB0LvQvtCy0LAg0Y/QuIDQvNGW0LM=

Encrypted string length:
556
```

Base64 compression of *text2.txt*

```
На початку березня 1917 року до Києва династ?я Романових пала. ?мператор Микола ?? зр?кся п
В умовах невизначеност? майбутнього пол?тичного життя професор Михайло Грушевський побачив
их пол?тичних сил та орган?зац?й, як л?берального, так ? соц?ал?стичного спрямування.

Original string length:
557

Encrypted string:
0J3QsCDQv9C+0YfQsNGC0LrRgyDQsdC10YDQtdC30L3RjyAxOTE3INGA0L7QutGDINC00L4g0JrQuNGU0LLQsCDQtNC
iDQktC70LDQtNGDINC/0LXRgNC10LnQvdGP0LvQuCDRgNC+0YHRltC50YHRjNC60ZYg0LvRltCx0LXRgNCw0LvRjNC9
PQstGW0LnRiNC70Lgg0LIg0ZbRgdGc0L7RgNGW0Y4g0Y/QuIDQm9G00YlQvdC10LLQsCDRgNC10LLQvtC70Y7RhTGW0
A0L7RhNC10YHQvtGAINCc0LjRhdcw0LnQu9C+INCT0YDRg9GI0LXQstGB0YzQutC40Lkg0L/QvtCx0LDRh9C40LIg0Y
gNCw0LvRjNC90Ymg0KDQsNC00YMuINCm0LXQuSDQvtGA0LPQsNC9INGD0YlQstC+0YDQuNC70Lgg0Ymg0LHQtdGA0LX
SDRgdC40Lsg0YlQsCDQvtGA0LPQsNC90ZbQt9Cw0YbRltC5LCDRj9C6INC70ZbQsdC10YDQsNC70YzQvdC+0LPQviwg

Encrypted string length:
1364
```

Base64 compression of *text3.txt*

```
Original string:
Друга екран?зац?я псевдо?сторичної пов?ст? ?вана Франка <Захар Беркут> (?снує ще ф?льм 1971 р.) - на
ржк?но), неабияк? плани просування ф?льму на зах?дний ринок. <Захар Беркут> ?з самого початку зн?має

Original string length:
457

Encrypted string:
0JTRgNGD0LPQsCDQtDC60YDQsNC90ZbQt9Cw0YbRltGPINC/0YHQtdCy0LTQvtGW0YHRgtC+0YDQuNGH0L3QvtGXINC/0L7QstGv
dCw0LnQsNC80LHRltGC0L3RltGI0LAg0YlQsCDQvdCw0LnQtNC+0YDQvtC20YfQsCDRgdGC0YDRltGH0LrQsCDRh9Cw0YHRltCyI
GI0YlQvtGA0LjRgSDQsiAXmTMsNSDQvNC70L0uINCz0YDQvS4gKDMwINC80LvQvS4g0L3QsNC00LDQvdC+INCu0LXRgNC20LrRlt
Cq9CX0LDRhdCw0YAg0JHQtdGA0LrRg9GCwrsG0ZbQtYDRgdCw0LzQvtCz0L4g0L/QvtGH0LDRgtC60Ymg0LfQvdGW0LzQsNCy0YH
0LLQvtGOLiDQktGW0L0g0ZYg0LLQuNCz0LvRj9C00LDRlCDRj9C6INCz0L7Qu9C70ZbQstG00LTRgdGM0LrQuNC5INGE0ZbQu9Gf

Encrypted string length:
1104
```

Size of file (bytes):

text1encrypted.txt = $1267\text{bit}/8 = 158.375$ bytes

text2encrypted.txt = $1554\text{bit}/8 = 194.25$ bytes

text3encrypted.txt = $1486\text{bit}/8 = 185.75$ bytes

Entropy:

text1.txt – 4.25588104

text2.txt – 2.27625325

text3.txt – 2.63565820

Size of file (bytes):

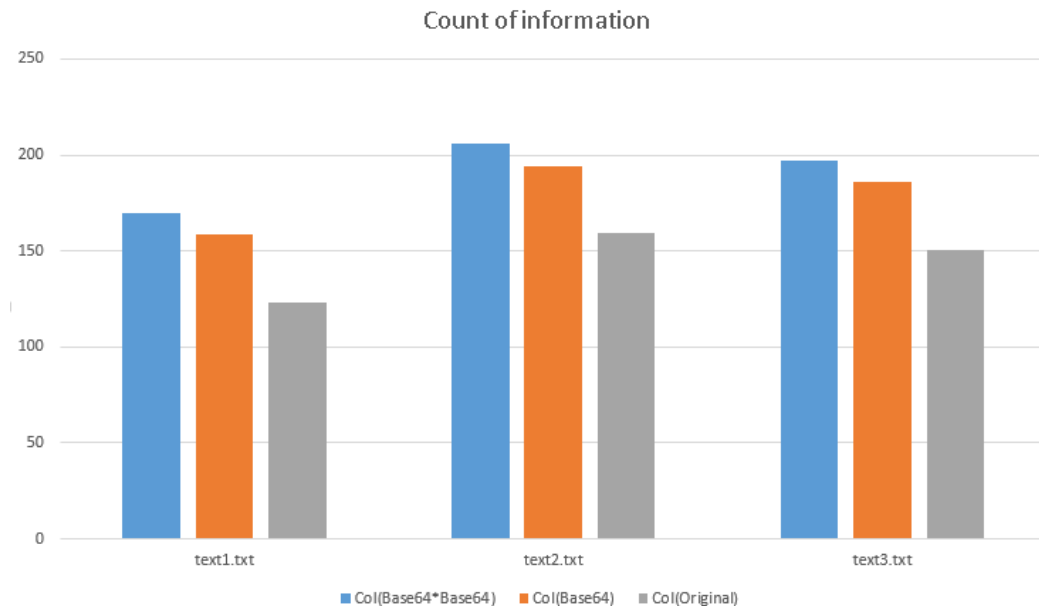
text1.txt – $987\text{bit}/8 = 123.375$ bytes

text2.txt – $1277\text{bit}/8 = 159.625$ bytes

text3.txt – $1204\text{bit}/8 = 150.5$ bytes

Назва файлу		Entropy		Length	Count of information	
text1encrypted.txt		2,2793		556	158.375	
text2encrypted.txt		1,1398		1364	194.25	
text3encrypted.txt		1,3469		1104	185.75	
Назва файлу	BZip2	RAR	Lzma	Zlib	Lzf4	
text1.txt	221	321	356	242	330	
text2.txt	438	595	590	513	603	
text3.txt	395	525	492	446	541	
Назва файлу		BZip(BASE64)		Count of information		Number of information
text1encrypted.txt		335		158.375		123.375
text2encrypted.txt		630		194.25		159.625
text3encrypted.txt		577		185.75		150.5

Назва файлу	CoI(BASE64) ²	Number of information
text1encrypted.txt	170	123.375
text2encrypted.txt	205.875	159.625
text3encrypted.txt	197.375	150.5



Висновок: В цій лабораторній роботі мною була зроблена програма що аналізує текст та обчислює ймовірності появи букви в тексті, ентропію тексту, загальний обсяг інформації. Мною було проаналізований обсяг зайнятого простору на диску файлів з різним типом стиснення. Як результат аналізу можна сказати, що обсяг інформації є значно меншим аніж обсяг який займає стиснутий цей же файл на диску. Було встановлено що найкращим з перевірених алгоритмом стиснення є VZip2. Об'єм його файлів є найближчим до кількості інформації. Також були опрацьовані навички в кодуванні **Base64**. Як результат можна сказати, що с кожним повторним кодуванням файлу обсяг інформації збільшується. Це зумовлено тим, що сама кількість символів в файлі збільшується, а отже й сам обсяг збільшується також.

My personal GitHub link: <https://github.com/gulivec84>

