

Київський національний університет імені Тараса Шевченка
Факультет радіофізики, електроніки та комп'ютерних систем
Навчальна дисципліна «Комп'ютерні системи»

Звіт з лабораторної роботи №1
на тему «Дослідження кількості інформації
при різних варіантах кодування»

Роботу виконав
Студент 3 курсу
КІ, група МА
Мормуль Олексій
Володимирович

Київ 2019

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

Хід роботи

Дослідження кількості інформації в тексті

1. Оберіть 3 текстових файла різного тематичного та лінгвістичного спрямування (файли також є у репозиторії):
 - [Об'єкт SCP-173](#)
 - [Гальмування двигуном](#)
 - [Перекладений текст пісні Rick Astley - Never Gonna Give You Up](#) (натисніть, щоб пійматися на рікролл українською)
2. Переконайтесь, що тексти, які ви використовуєте є унікальними і не повторюються у ваших колег! Використовуйте наявні електронні засоби зв'язку та документообігу, щоб уникнути дублювання! Вдруге аналіз того самого тексту не зараховується!
 - Заради цього створено мною [цей документ](#)
3. Код створеної програми міститься у репозиторії, посилання буде вкінці.
Результат роботи програми:

файл для аналізу: text1.enc

Загальна кількість символів файлу: 338

[illegible]

Відносна частота появи літери "щ" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "ь" у тексті = 0,0177514792899408 ; Літера присутня у тексті: 6 разів.
Відносна частота появи літери "ю" у тексті = 0,0236686390532544 ; Літера присутня у тексті: 8 разів.
Відносна частота появи літери "я" у тексті = 0,0177514792899408 ; Літера присутня у тексті: 6 разів.
Відносна частота появи літери "і" у тексті = 0,029585798816568 ; Літера присутня у тексті: 10 разів.
Відносна частота появи літери "ї" у тексті = 0,14792899408284 ; Літера присутня у тексті: 50 разів.
Відносна частота появи літери "(" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери ")" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "." у тексті = 0,00887573964497041 ; Літера присутня у тексті: 3 разів.
Відносна частота появи літери ":" у тексті = 0,00591715976331361 ; Літера присутня у тексті: 2 разів.
Відносна частота появи літери "-" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "''" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "0" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "1" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "2" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "3" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "4" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "5" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "6" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "7" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "8" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "9" у тексті = 0 ; Літера присутня у тексті: 0 разів.

Середня ентропія нерівномірного алфавіту у заданому тексті: 4,6186527830823
Кількість інформації у тексті: 195,138080085227

Загальна кількість символів файлу: 4694

[illegible]

Відносна частота появи літери "ш" у тексті = 0,00191734128674904 ; Літера присутня у тексті: 9 разів.
Відносна частота появи літери "б" у тексті = 0,0249254367277375 ; Літера присутня у тексті: 117 разів.
Відносна частота появи літери "ю" у тексті = 0,00511291009799744 ; Літера присутня у тексті: 24 разів.
Відносна частота появи літери "я" у тексті = 0,0157648061354921 ; Літера присутня у тексті: 74 разів.
Відносна частота появи літери "і" у тексті = 0,0136344269279932 ; Літера присутня у тексті: 64 разів.
Відносна частота появи літери "ї" у тексті = 0,137835534725181 ; Літера присутня у тексті: 647 разів.
Відносна частота появи літери "(" у тексті = 0,00489987217724755 ; Літера присутня у тексті: 23 разів.
Відносна частота появи літери ")" у тексті = 0,00489987217724755 ; Літера присутня у тексті: 23 разів.
Відносна частота появи літери ":" у тексті = 0,0129953131657435 ; Літера присутня у тексті: 61 разів.
Відносна частота появи літери ";" у тексті = 0,00063911376224968 ; Літера присутня у тексті: 3 разів.
Відносна частота появи літери "-" у тексті = 0,00106518960374947 ; Літера присутня у тексті: 5 разів.
Відносна частота появи літери "''" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "0" у тексті = 0,000852151682999574 ; Літера присутня у тексті: 4 разів.
Відносна частота появи літери "1" у тексті = 0,0034086067319983 ; Літера присутня у тексті: 16 разів.
Відносна частота появи літери "2" у тексті = 0,00127822752449936 ; Літера присутня у тексті: 6 разів.
Відносна частота появи літери "3" у тексті = 0,00063911376224968 ; Літера присутня у тексті: 3 разів.
Відносна частота появи літери "4" у тексті = 0,00127822752449936 ; Літера присутня у тексті: 6 разів.
Відносна частота появи літери "5" у тексті = 0,00149126544524925 ; Літера присутня у тексті: 7 разів.
Відносна частота появи літери "6" у тексті = 0,00063911376224968 ; Літера присутня у тексті: 3 разів.
Відносна частота появи літери "7" у тексті = 0,000852151682999574 ; Літера присутня у тексті: 4 разів.
Відносна частота появи літери "8" у тексті = 0,000852151682999574 ; Літера присутня у тексті: 4 разів.
Відносна частота появи літери "9" у тексті = 0,00106518960374947 ; Літера присутня у тексті: 5 разів.

Середня ентропія нерівномірного алфавіту у заданому тексті: 4,78265636880096

Кількість інформації у тексті: 2806,22362439397

Загальна кількість символів файлу: 3584

[illegible]

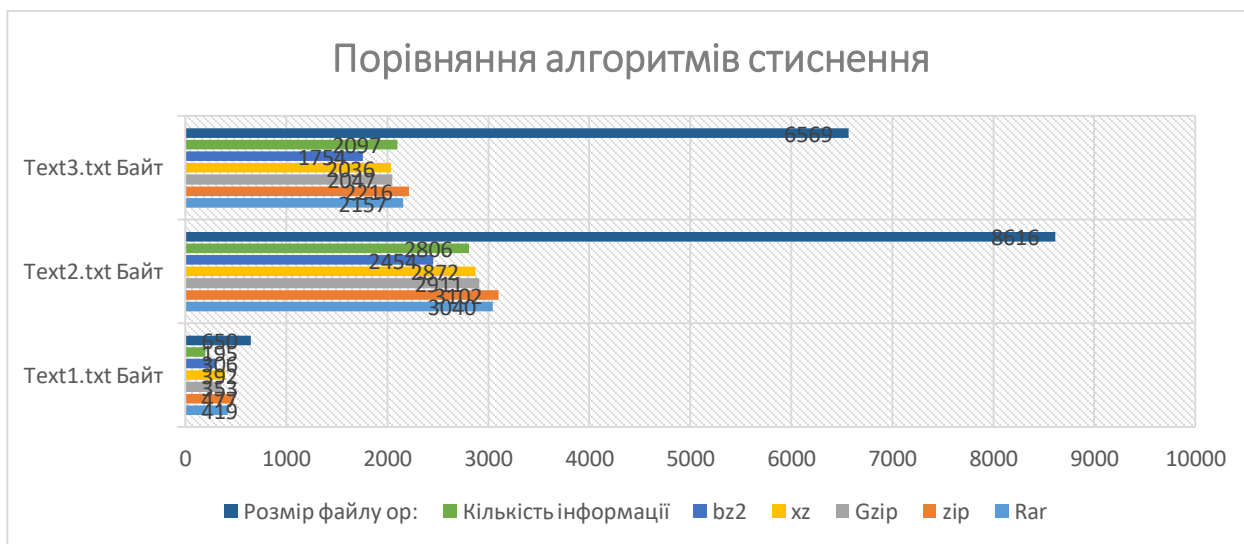
Відносна частота появи літери "ь" у тексті = 0,0139508928571429 ; Літера присутня у тексті: 50 разів.
Відносна частота появи літери "ю" у тексті = 0,00837053571428571 ; Літера присутня у тексті: 30 разів.
Відносна частота появи літери "я" у тексті = 0,0189732142857143 ; Літера присутня у тексті: 68 разів.
Відносна частота появи літери "і" у тексті = 0,00864955357142857 ; Літера присутня у тексті: 31 разів.
Відносна частота появи літери " " у тексті = 0,132254464285714 ; Літера присутня у тексті: 474 разів.
Відносна частота появи літери "(" у тексті = 0,001953125 ; Літера присутня у тексті: 7 разів.
Відносна частота появи літери ")" у тексті = 0,001953125 ; Літера присутня у тексті: 7 разів.
Відносна частота появи літери "." у тексті = 0,0108816964285714 ; Літера присутня у тексті: 39 разів.
Відносна частота появи літери ":" у тексті = 0,000837053571428571 ; Літера присутня у тексті: 3 разів.
Відносна частота появи літери "-" у тексті = 0,00139508928571429 ; Літера присутня у тексті: 5 разів.
Відносна частота появи літери "" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "0" у тексті = 0,00279017857142857 ; Літера присутня у тексті: 10 разів.
Відносна частота появи літери "1" у тексті = 0,000279017857142857 ; Літера присутня у тексті: 1 разів.
Відносна частота появи літери "2" у тексті = 0,00111607142857143 ; Літера присутня у тексті: 4 разів.
Відносна частота появи літери "3" у тексті = 0,00111607142857143 ; Літера присутня у тексті: 4 разів.
Відносна частота появи літери "4" у тексті = 0,000558035714285714 ; Літера присутня у тексті: 2 разів.
Відносна частота появи літери "5" у тексті = 0,000279017857142857 ; Літера присутня у тексті: 1 разів.
Відносна частота появи літери "6" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "7" у тексті = 0 ; Літера присутня у тексті: 0 разів.
Відносна частота появи літери "8" у тексті = 0,000279017857142857 ; Літера присутня у тексті: 1 разів.
Відносна частота появи літери "9" у тексті = 0 ; Літера присутня у тексті: 0 разів.

Середня ентропія нерівномірного алфавіту у заданому тексті: 4,68133589749239
Кількість інформації у тексті: 2097,23848207659

4. Проведіть стиснення кожного вхідного файлу за допомогою 5 різних алгоритмів стиснення.

Для цього завдання я використовував WinRAR та 7-Zip. Для кожного алгоритму я використовував рівень стиснення Normal.

Файл	Text1.txt Байт	Text2.txt Байт	Text3.txt Байт
Rar	419	3040	2157
Zip	477	3102	2216
Gzip	353	2911	2047
Xz	392	2872	2036
Bz2	306	2454	1754
Кількість інформації	195	2806	2097
Розмір файлу ор:	650	8616	6569



5. У результаті ідеального стиснення розмір файлу повинен бути рівним кількості інформації. Але у реальності розміри архівованих файлів у більшості випадків дещо більші за кількість інформації, окрім випадку з великим текстом.

Це відбувається тому, що алгоритми архіваторів побудовані таким чином аби використати повторювані частини тексту. Виходячи з цього, формула розрахунку кількості інформації, використана для програми, не є досконалою, бо вона не враховує передбачення наступного шматочку тексту.

- Хочу звернути уваги на алгоритм bzip2, який виявився найефективнішим у всіх випадках. А також навіть упорався «ідеальним стисненням» у випадку великого файлу, тобто стиснений архів має розмір навіть менший, ніж кі-сть інформації. У випадку з піснею, хоч файл і меншого розміру, ніж текст про SCP-173, але є багато частин повторюваного тексту (як-от приспів)

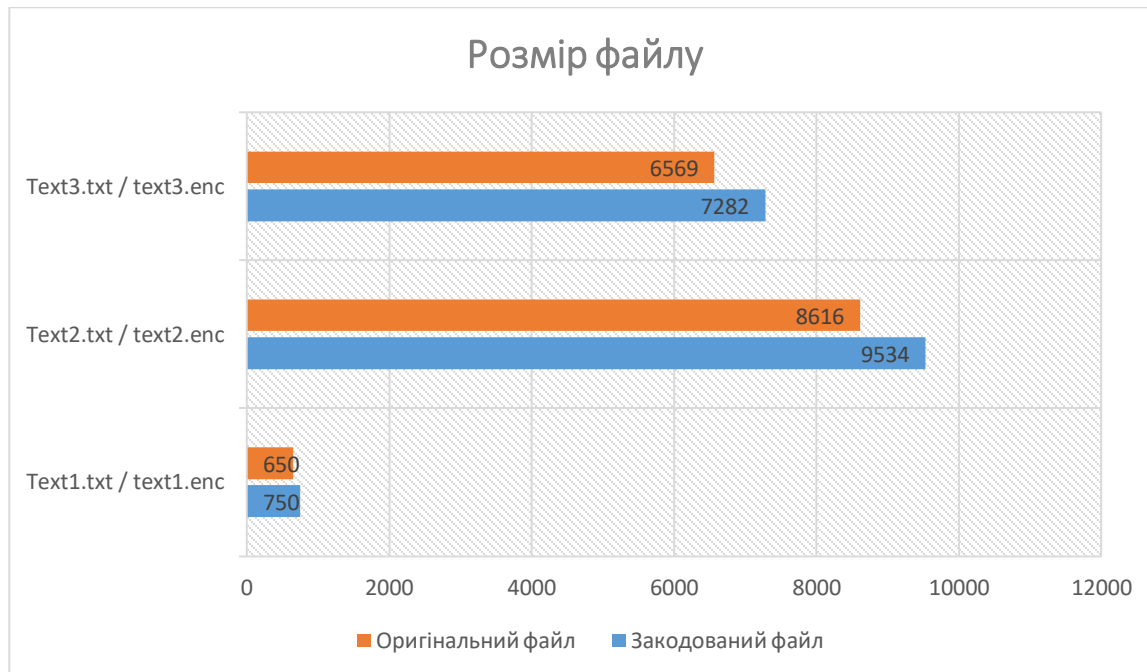
[illegible]

[illegible]

Кількість інформації

Файл	Оригінал	base64
Text3.txt / text3.enc	1662	2097
Text2.txt / text2.enc	2312	2806
Text1.txt / text1.enc	195	195

	Text1.txt / text1.enc	Text2.txt / text2.enc	Text3.txt / text3.enc
Base64	290	2806	2097
Оригінал	195	2312	1662



Можна помітити, що розмір закодованих файлах зріс. Це пов'язано з алгоритмом кодування base64 – перетворення, наприклад 3 октетів (по 8 біт) у 4 секстети (по 6 біт), що збільшує розмір на третину.

Найефективнішим алгоритмом стиснення виявився bzip2.

Висновок

У ході виконання лабораторної роботи ознайомився з поняттям ентропії інформації та пов'язаних з ентропією понять. Порівняв алгоритми стиснення – обрав кращий з них для випадків, коли треба буде зекономити місце на носії. Закріпив теоретичні знання практикою. Теоретично та практично ознайомився з алгоритмом кодування Base64, його перевагами та недоліками.

Код програм, звіт та текстові файли, використані у роботі містяться у репозиторії за посиланням: [\(тут\)](#)