

Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Сибирский государственный университет телекоммуникаций и
информатики»
(СибГУТИ)

Кафедра прикладной математики и кибернетики

Теория Информации
Лабораторная работа №4

Выполнил:
Студент IV курса ИВТ,
группы ИП-713
Михеев Никита Алексеевич

Работу проверил:
доцент кафедры ПМиК
Мачикина Е. П.

Новосибирск 2020 г.

1. Постановка задачи

Цель работы:

Экспериментальное изучение процесса сжатия текстового файла с помощью бинарного кодирования.

Задание:

1. Запрограммировать процедуру двоичного кодирования текстового файла. В качестве метода кодирования использовать или метод Шеннона, или метод Фано, или метод Хаффмана. Текстовые файлы использовать те же, что и в практических работах 1, 2, 3.
2. Проверить, что построенный код для каждого файла является префиксным. Вычислить среднюю длину кодового слова и оценить избыточность каждого построенного кода.
3. После кодирования текстового файла вычислить оценки энтропии выходной последовательности, используя частоты отдельных символов, пар символов и троек символов и заполнить таблицу.
4. Оформить отчет, загрузить отчет и файл с исходным кодом в электронную среду. Отчет обязательно должен содержать заполненную таблицу и анализ полученных результатов.

2. Ход работы

Для выполнения лабораторной работы была написана программа на языке Python версии 3.9, которая сначала считывает текстовый файл с исходным кодом, приводит его к удобному для работы программы виду. Для кодирования текста был выбран метод Шеннона.

Затем идет подсчет каждого символа (или их последовательности) и вычисляются вероятности их встречи относительно всего текста. Далее вычисляются длины кодовых слов, а также кумулятивные вероятности.

После этого в программе начинается вычисление кодов символов: берется двоичное представление кумулятивной вероятности встречи символа, из дробной части берется такое количество символов, сколько было подсчитано выше. Результат выводится на экран.

3. Результат работы

```
 : 0.1906 - 000
e: 0.0986 - 0011
t: 0.0731 - 0100
o: 0.0666 - 0101
a: 0.0635 - 0110
h: 0.0590 - 01111
n: 0.0548 - 10001
i: 0.0502 - 10011
r: 0.0501 - 10101
s: 0.0478 - 10110
d: 0.0376 - 11001
l: 0.0312 - 110010
u: 0.0241 - 110100
f: 0.0222 - 110110
w: 0.0186 - 110111
m: 0.0185 - 111000
c: 0.0181 - 111010
y: 0.0168 - 111011
g: 0.0151 - 1111101
b: 0.0142 - 1111010
p: 0.0116 - 1111110
v: 0.0080 - 1111101
k: 0.0071 - 11111101
x: 0.0008 - 11111111011
j: 0.0008 - 11111111100
z: 0.0004 - 111111111100
q: 0.0004 - 111111111110
```

Рис.1 – вывод на экран символа, его вероятности и кодовое слово.

```
Original text: 4.075871094438105
```

Рис.2 – энтропия источника

```
L average: 4.6009678019728275
Coded entropy for 1 symbols in a row: 0.9997165783152567
Coded entropy for 2 symbols in a row: 0.9996531013458817
Coded entropy for 3 symbols in a row: 0.9987596407383864
r = 0.5250967075347228
```

Рис.3 – средняя длина кодового слова, энтропия закодированного текста, избыточность кодирования.

Метод кодирования	Название текста	Оценка избыточности кодирования	Оценка энтропии выходной посл-ти (частоты символов)	Оценка энтропии выходной посл-ти (частоты пар символов)	Оценка энтропии выходной посл-ти (частоты троек символов)
Метод Шеннона	«The White Company»	0.5251	0.99971	0.99965	0.99876

Таблица 1 – полученные данные

Вывод:

При кодировании методом Шеннона был получен префиксный код, в котором используется избыточность сообщения (коды более частых символов состоят из коротких последовательностей, а коды более редких символов – из более длинных). По нашим данным выяснилось, что энтропия оригинального текста далека от реальных значений и поэтому избыточность получилась сравнительно небольшой. Но, если взять энтропию реально текста (около 2), то избыточность будет большой (~2.6), что показывает не оптимальность кодировки Шеннона.