

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

Кафедра прикладной математики и кибернетики

Теория Информации
Лабораторная работа №5

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 троичное 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 - 00
e: 0.0986 - 012
t: 0.0731 - 021
o: 0.0666 - 121
a: 0.0635 - 102
h: 0.0590 - 111
n: 0.0548 - 112
i: 0.0502 - 121
r: 0.0501 - 122
s: 0.0478 - 201
d: 0.0376 - 202
l: 0.0312 - 2101
u: 0.0241 - 2110
f: 0.0222 - 2112
w: 0.0186 - 2121
m: 0.0185 - 2122
c: 0.0181 - 2201
y: 0.0168 - 2202
g: 0.0151 - 2211
b: 0.0142 - 2212
p: 0.0116 - 22201
v: 0.0080 - 22211
k: 0.0071 - 22220
x: 0.0008 - 2222210
j: 0.0008 - 2222212
z: 0.0004 - 22222211
q: 0.0004 - 22222220
```

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

```
L average: 3.052205471803462
Coded entropy for 1 symbols in a row: 1.5691209415605152
Coded entropy for 2 symbols in a row: 1.5626590808958385
Coded entropy for 3 symbols in a row: 1.5480100302254158
```

Рис.2 – энтропия закодированного источника с различной длиной цепочек.

```
L average: 3.052205471803462
```

Рис.3 – средняя длина кодового слова.

Метод кодирования	Название текста	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»	1.5691	1.5627	1.5480

Таблица 1 – полученные данные

Вывод:

При кодировании методом Шеннона был получен префиксный код (коды более частых символов состоят из коротких последовательностей в троичном представлении, а коды более редких символов – из более длинных).

При сравнении с бинарным кодированием по методу Шеннона, средняя длина у троичного кода меньше почти на 1.6.

Из всего этого можно сделать вывод, что хоть код остается префиксным, энтропия близка к максимальной, но наиболее эффективно информация кодируется с помощью бинарного кода, так как при дальнейшем увеличении количества символов в коде идет уменьшение средней длины слов и увеличивается разрыв между закодированным источником и предельной энтропией.