



МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ

Федеральное государственное бюджетное образовательное

учреждение высшего образования

"МИРЭА - Российский технологический университет"

РТУ МИРЭА

Институт информационных технологий (ИТ)

Кафедра математического обеспечения и стандартизации информационных

технологий (МОСИТ)

ОТЧЕТ ПО ПРАКТИЧЕСКОЙ РАБОТЕ № 8.1

по дисциплине

«Структуры и алгоритмы обработки данных»

Тема: Алгоритмы кодирования и сжатия данных

Выполнил студент группы ИКБО-41-23

Гольд Д.В.

Принял старший преподаватель

Рысин М.Л.

Москва 2024

СОДЕРЖАНИЕ

ЗАДАНИЕ 1	3
ЦЕЛЬ:	3
1. Закодировать фразу методами Шеннона–Фано	3
2. Сжатие данных по методу Лемпеля–Зива LZ77, используя двухсимвольный алфавит (0, 1) закодировать следующую фразу:.....	4
3. Закодировать следующую фразу, используя код LZ78	5
ЗАДАНИЕ 2	6
ЦЕЛЬ:	6
Алгоритм Шеннона-Фано:	6
Код программы:.....	7
Тестирование:	9
Алгоритм Хаффмана:	10
Код программы:.....	11
Тестирование:	14
Практическая оценка сложности алгоритма Хаффмана	15
ВЫВОД	16

ЗАДАНИЕ 1

ЦЕЛЬ:

Исследовать такие алгоритмы сжатия, как Шеннон-Фано, LZ77, LZ78

1. Закодировать фразу методами Шеннона–Фано

Эни-бени рити-Фати. Дорба, дорба сентибрати. Дэл. Дэл. Кошка. Дэл. Фати!
(72 символа)

Незакодированная фраза = 576 бит

Закодированная фраза = 285 бит

Символ	Кол-во	Вероятность	Код	Кол-во бит
	9	0.125	001	27
и	8	0.111	000	24
а	6	0.083	010	18
.	6	0.083	0111	24
т	5	0.0694	0110	20
б	4	0.055	1100	16
р	4	0.055	1101	16
Д	4	0.055	1110	16
н	3	0.0416	11110	15
о	3	0.0416	11111	15
л	3	0.0416	10000	12
г	2	0.027	10010	10
е	2	0.027	10011	10
Ф	2	0.027	101011	12

Э	1	0.0138	101010	6
д	1	0.0138	101000	6
ь	1	0.0138	101001	6
ш	1	0.0138	101110	6
!	1	0.0138	101100	6
к	1	0.0138	101101	6
с	1	0.0138	1011110	7
К	1	0.0138	1011111	7

2. Сжатие данных по методу Лемпеля–Зива LZ77, используя двухсимвольный алфавит (0, 1) закодировать следующую фразу:

<u>Исходный текст</u>	0.00.10.01.011.001.000.100.01
<u>LZ-код</u>	0.00.010.001.1001.0101.0100.0110.0100
<u>R</u>	2 3 4
<u>Вводимые коды</u>	- 00 10 01 011 001 000 100

Таблица соответствия кодов и символов:

Символ	Код, R = 1	Код, R = 2	Код, R = 3	Код, R = 4
0	0	00	000	0000
1	1	01	001	0001
00		10	010	0010
10		11	011	0011
01			100	0100
011			101	0101
001			110	0110
000			111	0111
100				1000

3. Закодировать следующую фразу, используя код LZ78

klolklonkolonklonkl

Словарь	Считываемое содержимое	Код
	k	<0, k>
k = 1	l	<0, l>
k = 1, l = 2,	o	<0, o>
k = 1, l = 2, o = 3,	kl	<1, l>
k = 1, l = 2, o = 3, kl = 4	klo	<4, o>
k = 1, l = 2, o = 3, kl = 4, klo = 5	n	<0, n>
k = 1, l = 2, o = 3, kl = 4, klo = 5, n = 6,	ko	<1, o>
k = 1, l = 2, o = 3, kl = 4, klo = 5, n = 6, ko = 7	lo	<2, o>
k = 1, l = 2, o = 3, kl = 4, klo = 5, n = 6, ko = 7, lo = 8, nk = 9	nk	<6, k>
k = 1, l = 2, o = 3, kl = 4, klo = 5, n = 6, ko = 7, lo = 8, nk = 9	lon	<8, n>
k = 1, l = 2, o = 3, kl = 4, klo = 5, n = 6, ko = 7, lo = 8, nk = 9, lon = 10	kl	<4, EOF>

<0, k><0, l><0, o><1, l><4, o><0, n><1, o><2, o><6, k><8, n><4, EOF>

ЗАДАНИЕ 2

ЦЕЛЬ:

Разработать программы сжатия и восстановления текста методами Хаффмана и Шеннона–Фано.

Алгоритм Шеннона-Фано:

Формирование кодов:

- Символы сортируются по частоте вхождения в строку.
- Разделяются на группы с равной суммой вероятностей.
- Каждой группе присваиваются биты 0 и 1.
- Процесс повторяется до назначения кода каждому символу.

Код программы:

```
1  #include <iostream>
2  #include <vector>
3  #include <map>
4  #include <string>
5  #include <algorithm>
6  #include <locale>
7  #include <codecvt>
8
9  using namespace std;
10
11 struct Node { //узел дерева шеннона-фано
12     char symbol;
13     double symbol_frequency;
14     string character_code;
15 };
16
17 bool compare_frequency(Node tree_nodes_a, Node tree_nodes_b) { //сравниваем частот для сортировки
18     return tree_nodes_a.symbol_frequency > tree_nodes_b.symbol_frequency;
19 }
20
21 void generate_codes(vector<Node>& symbol_nodes, int start, int end, string character_code) { //рекурс функ - построение кодов
22     if (start == end) {
23         symbol_nodes[start].character_code = character_code;
24         return;
25     }
26
27     double sum_symbol_frequency_left = 0, sum_symbol_frequency_right = 0;
28     int left_froup_end_index = start;
29
30     for (int i = start; i <= end; ++i) sum_symbol_frequency_right += symbol_nodes[i].symbol_frequency;
31
32     for (int i = start; i <= end; ++i) {
33         sum_symbol_frequency_left += symbol_nodes[i].symbol_frequency;
34         sum_symbol_frequency_right -= symbol_nodes[i].symbol_frequency;
35
36         if (sum_symbol_frequency_left >= sum_symbol_frequency_right) {
37             sum_symbol_frequency_right -= symbol_nodes[i].symbol_frequency;
38
39             if (sum_symbol_frequency_left >= sum_symbol_frequency_right) {
40                 left_froup_end_index = i;
41                 break;
42             }
43         }
44
45         generate_codes(symbol_nodes, start, left_froup_end_index, character_code + "0");
46         generate_codes(symbol_nodes, left_froup_end_index + 1, end, character_code + "1");
47     }
48 }
49
50 int main() {
51     setlocale(LC_ALL, "");
52
53     //string input = "Эни-бени-рити-Фати. Дорба, дорба-сентибрати. Дэл. Дэл. Кошка. Дэл. Фати!";
54     //э = q ; Э = Q ; ш = s;
55
56     string input = "Qni-beni riti-Fati. Dorba, dorba centibrati. Dql. Dql. Koska. Dql. Fati!";
57
58     map<char, double> frequency_map; //частота симв
59     for (char c : input) {
60         frequency_map[c]++;
61     }
62
63     vector<Node> symbol_nodes;
64     for (auto pair : frequency_map) {
65         symbol_nodes.push_back({ pair.first, pair.second / input.size(), "" });
66     }
67
68     sort(symbol_nodes.begin(), symbol_nodes.end(), compare_frequency);
69
70     generate_codes(symbol_nodes, 0, symbol_nodes.size() - 1, ""); //строим дерево
71
72     map<char, string> codes;
73     for (auto node : symbol_nodes) {
```

```

69     map<char, string> codes;
70     for (auto node : symbol_nodes) {
71         codes[node.symbol] = node.character_code;
72         cout << node.symbol << " -> " << node.character_code << endl;
73     }
74
75     string encoded = ""; //кодирование
76     for (char c : input) {
77         encoded += codes[c];
78     }
79
80     cout << "encoded text: " << encoded << endl;
81
82     //сжатый текст
83     int original_size_bits = input.size() * 8; //исходный размер в битах
84     int compressed_size_bits = encoded.size(); //сжатый размер в битах
85
86     double compression_ratio = (double)compressed_size_bits / original_size_bits * 100; //процент сжатия
87
88     string decoded = ""; //декодирование
89     string temp = "";
90     for (char bit : encoded) {
91         temp += bit;
92         for (auto node : symbol_nodes) {
93             if (node.character_code == temp) {
94                 decoded += node.symbol;
95                 temp = "";
96                 break;
97             }
98         }
99     }
100
101     double average_code_length = 0.0; //вычисление средней длины кода
102     for (auto node : symbol_nodes) {
103         average_code_length += node.character_code.length() * node.symbol_frequency;
104     }
105
106     double variance = 0.0; //вычисление дисперсии длины кода
107     for (auto node : symbol_nodes) {
108         double length_diff = node.character_code.length() - average_code_length;
109         variance += pow(length_diff, 2) * node.symbol_frequency;
110     }
111
112     cout << "\ncompressed text in bits: " << compressed_size_bits << endl;
113     cout << "compression ratio: " << compression_ratio << "%\n";
114     cout << "\ndecoded text: " << decoded << endl;
115     cout << "\naverage code length: " << average_code_length << endl;
116     cout << "code length variance: " << variance << endl;
117
118     return 0;
119 }
120

```


Тестирование:

```
-> 0000
i -> 0001
. -> 001
a -> 0100
t -> 0101
D -> 011
b -> 10000
r -> 10001
l -> 1001
n -> 10100
o -> 10101
q -> 1011
- -> 11000
F -> 11001
e -> 11010
! -> 110110
, -> 110111
K -> 111000
Q -> 111001
c -> 111010
d -> 111011
k -> 11110
s -> 11111
encoded text: 1110011010000011100010000110101010000010000100010001010100011100011001010001010001001000001110101100011000
001001101110000111011101011000110000010000001110101101010100010100011000010001010001010001001000001110111001001000001110
11100100100001110001010111111111100100001000001110111001001000011001010001010001110110
compressed text in bits: 312
compression ratio: 54.1667%

decoded text: Qni-beni riti-Fati. Dorba, dorba centibrati. Dql. Dql. Koska. Dql. Fati!

average code length: 4.33333
code length variance: 0.666667
```

Алгоритм Хаффмана:

Составление списка узлов:

- Символы сортируются по частоте вхождения в строку.

Построение дерева:

- Найти два узла с наименьшими весами.
- Объединить их в новый узел с весом, равным сумме их частот.
- Новый узел становится родительским для этих двух.
- Повторять процесс, пока не останется один узел (корень дерева).

Назначение кодов:

- Пройти по дереву от корня к листьям.
- Каждой левой ветви присвоить 0, а правой — 1.
- Листовые узлы содержат символы с их уникальными бинарными кодами.

Код программы:

```
1  #include <iostream>
2  #include <string>
3  #include <map>
4  #include <queue>
5  #include <vector>
6  #include <algorithm>
7
8  using namespace std;
9
10 struct huffman_node { //узел дерева хаффмана
11     char symbol;
12     double symbol_frequency;
13     huffman_node* left_child;
14     huffman_node* right_child;
15 };
16
17 struct compare { //компаратор для приоритетной очереди
18     bool operator()(huffman_node* tree_nodes_a, huffman_node* tree_nodes_b) {
19         return tree_nodes_a->symbol_frequency > tree_nodes_b->symbol_frequency;
20     }
21 };
22
23 void generate_huffman_codes(huffman_node* node, string current_code, map<char, string>& huffman_codes) {
24     if (!node) return;
25     //рекурсивная функция для генерации кодов символов
26
27     if (!node->left_child && !node->right_child) { //если узел листовой - добавляем символ и его код
28         huffman_codes[node->symbol] = current_code;
29     }
30
31     //рекурсивно вызываем для левого и правого поддеревя
32     generate_huffman_codes(node->left_child, current_code + "0", huffman_codes);
33     generate_huffman_codes(node->right_child, current_code + "1", huffman_codes);
34 }
35
36 void delete_huffman_tree(huffman_node* node) { //рекурсивное удаление дерева
37     if (!node) return;
38     delete_huffman_tree(node->left_child);
39     delete_huffman_tree(node->right_child);
40 }
```

```

36 void delete_huffman_tree(huffman_node* node) { //рекурсивное удаление дерева
37     if (!node) return;
38     delete_huffman_tree(node->left_child);
39     delete_huffman_tree(node->right_child);
40     delete node;
41 }
42
43 int main() {
44     string input_text = "Gold Daniil Vladimirovich";
45
46     map<char, double> frequency_map;
47     for (char c : input_text) {
48         frequency_map[c]++;
49     }
50
51     //создание приоритетной очереди
52     priority_queue<huffman_node*, vector<huffman_node*>, compare> priority_queue_huffman_nodes;
53     for (auto pair : frequency_map) {
54         huffman_node* new_huffman_node = new huffman_node{ pair.first, pair.second, nullptr, nullptr };
55         priority_queue_huffman_nodes.push(new_huffman_node);
56     }
57
58     //построение дерева хатфмана
59     while (priority_queue_huffman_nodes.size() > 1) {
60         huffman_node* left_child = priority_queue_huffman_nodes.top(); priority_queue_huffman_nodes.pop();
61         huffman_node* right_child = priority_queue_huffman_nodes.top(); priority_queue_huffman_nodes.pop();
62
63         huffman_node* merged = new huffman_node{ '\\0', left_child->symbol_frequency + right_child->symbol_frequency,
64             left_child, right_child };
65         priority_queue_huffman_nodes.push(merged);
66     }
67
68     huffman_node* root = priority_queue_huffman_nodes.top(); //корень дерева хатфмана
69     map<char, string> huffman_codes; //генерация кодов символов
70     generate_huffman_codes(root, "", huffman_codes);
71

```

```

69     map<char, string> huffman_codes; //генерация кодов символов
70     generate_huffman_codes(root, "", huffman_codes);
71
72     cout << "symbols codes:\n"; //отображение кодов символов
73     for (auto pair : huffman_codes) {
74         cout << pair.first << " -> " << pair.second << endl;
75     }
76
77     string encoded_text = "";
78     for (char c : input_text) {
79         encoded_text += huffman_codes[c];
80     }
81     cout << "\nencoded text: " << encoded_text << endl;
82
83     //сжатый текст
84     int original_size_bits = input_text.size() * 8; //исходный размер в битах (по 8 бит на символ)
85     int compressed_size_bits = encoded_text.size(); //сжатый размер в битах
86
87     //коэффициенты сжатия
88     double compression_ratio_ascii = (double)original_size_bits / compressed_size_bits;
89     double average_code_length = 0.0;
90     double variance = 0.0;
91
92     for (auto pair : huffman_codes) { //расчет средней длины кода и дисперсии
93         int code_length = pair.second.size();
94         average_code_length += code_length * (pair.second.size() / (double)input_text.size());
95     }
96
97     for (auto pair : huffman_codes) { //вычисление дисперсии длины кодов
98         int code_length = pair.second.size();
99         variance += (code_length - average_code_length) * (code_length - average_code_length) *
100             (pair.second.size() / (double)input_text.size());
101     }
102
103     //коэффициент сжатия относительно равномерного кода
104     int uniform_code_size = (int)ceil(log2(frequency_map.size())) * input_text.size(); //равномерный код
105     double compression_ratio_uniform = (double)original_size_bits / uniform_code_size;

```

```

103 //коэффициент сжатия относительно равномерного кода
104 int uniform_code_size = (int)ceil(log2(frequency_map.size())) * input_text.size(); //равномерный код
105 double compression_ratio_uniform = (double)original_size_bits / uniform_code_size;
106
107 cout << "\ncompressed size in bits: " << compressed_size_bits << endl;
108 cout << "compression ratio (ASCII): " << compression_ratio_ascii << endl;
109 cout << "compression ratio (uniform code): " << compression_ratio_uniform << endl;
110 cout << "average code length: " << average_code_length << " bits\n";
111 cout << "variance code length: " << variance << endl;
112
113 string decoded_text = ""; //декодирование текста
114 huffman_node* current_node = root;
115 for (char bit : encoded_text) {
116     if (bit == '0') {
117         current_node = current_node->left_child;
118     }
119     else {
120         current_node = current_node->right_child;
121     }
122
123     if (!current_node->left_child && !current_node->right_child) {
124         decoded_text += current_node->symbol;
125         current_node = root;
126     }
127 }
128
129 cout << "\ndecoded text: " << decoded_text << endl;
130
131 delete_huffman_tree(root);
132
133 return 0;
134 }
135

```

Тестирование:

```

symbols codes:
-> 1011
D -> 11100
G -> 11101
V -> 11010
a -> 1111
c -> 0000
d -> 1001
h -> 1000
i -> 01
l -> 001
m -> 11011
n -> 0001
o -> 1010
r -> 11000
v -> 11001

encoded text: 111011010001100110111110011110001010100110111101000111111001011101101110001010110010100001000

compressed size in bits: 93
compression ratio (ASCII): 2.15054
compression ratio (uniform code): 2
average code length: 11 bits
variance code length: 112.24

decoded text: Gold Daniil Vladimirovich

```

Таблица частот:

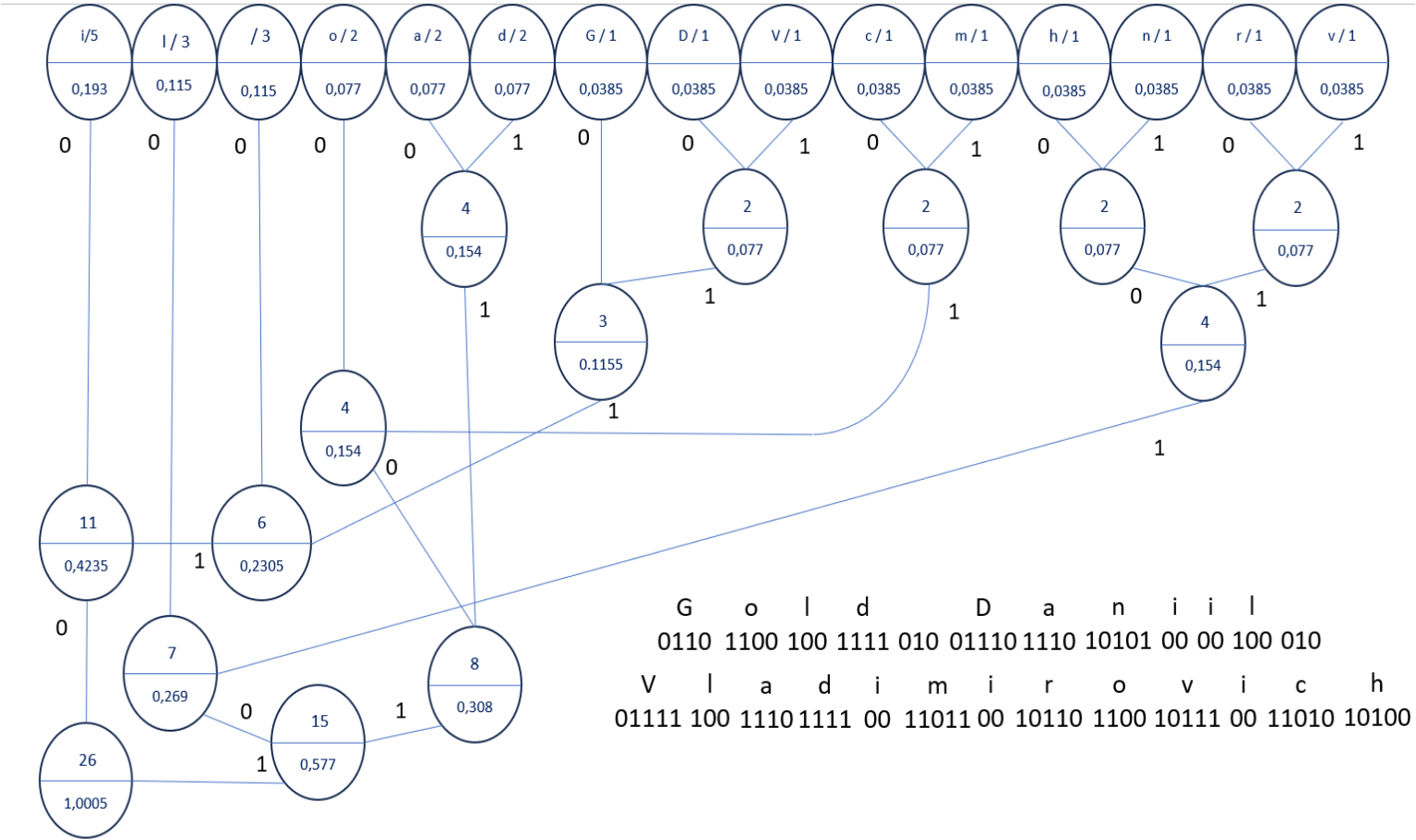
Алфавит	G	o	l	d		D	a	n	i	V
Кол.вх.	1	2	3	2	3	1	2	1	5	1
Вероятн.	0,0385	0,077	0,115	0,077	0,115	0,0385	0,077	0,0385	0,193	0,0385

Алфавит	m	r	v	c	h
Кол.вх.	1	1	1	1	1
Вероятн.	0,0385	0,0385	0,0385	0,0385	0,0385

Таблица отсортированных частот:

Алфавит	i	l		o	a	d	G	D	V	c
Кол.вх.	5	3	3	2	2	2	1	1	1	1
Вероятн.	0,193	0,115	0,115	0,077	0,077	0,077	0,0385	0,0385	0,0385	0,0385

Алфавит	m	h	n	r	v
Кол.вх.	1	1	1	1	1
Вероятн.	0,0385	0,0385	0,0385	0,0385	0,0385



Практическая оценка сложности алгоритма Хаффмана

Алгоритм Хаффмана состоит из нескольких этапов, каждый из которых имеет свою сложность:

- Подсчёт частоты символов: $O(m)$

Где m — длина исходного текста (проход по тексту для подсчёта частот).

- Построение дерева Хаффмана: $O(n \log(n))$

Где n — количество уникальных символов во входных данных

- Кодирование текста: $O(m)$

Замена каждого символа текста его кодом.

- Общая сложность: $O(m+n \log n)$

ВЫВОД

В ходе исследования были реализованы и протестированы алгоритмы сжатия текста: Шеннона-Фано, LZ77, LZ78 и Хаффмана.