

Теория кодирования и сжатия информации

Лабораторная работа №1

Гущин Андрей, 431 группа, 1 подгруппа

2022 г.

1 Задача

Разработать программу осуществляющую архивацию и разархивацию текстового файла используя алгоритм Хаффмана. Программы архивации и разархивации должны быть представлены отдельно и работать независимо друг от друга. Определить для данного шифра характеристики 1 (коэффициент сжатия) и 2 (скорость сжатия). К работе необходимо прикрепить отчет и программный проект.

2 Алгоритм

Алгоритм Хаффмана заключается в создании нового кода для каждого встречающегося в тексте символа на основе частоты встречи этого символа.

Алгоритм состоит из следующих шагов:

1. Вычислить частоты всех встретившихся символов;
2. Построить очередь с приоритетом на основе полученных частот;
3. Построить дерево с помощью очереди с приоритетом;
4. На основе дерева вычислить код для каждого символа.

3 Тестирование

Для проверки программы были использованы тестовые тексты 1 (рис. 1) и 6 (рис. 2). Можно заметить, что после распаковки архива полученный файл совпадает с исходным (проверка с помощью утилиты diff). Также можно заметить, что для файлов малого размера архив увеличивает их размер за счёт метаданных.

```
lab1 -zsh
$ ./lab1-rs
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab1-rs -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab1-rs -i Tectr_1.txt -o test1.gsch1 --compress
$ ./lab1-rs -i test1.gsch1 -o test.txt --decompress
$ diff Tectr_1.txt test.txt --report-identical-files
Files Tectr_1.txt and test.txt are identical
$ ls -l
total 1840
-rwxr-xr-x  1 aguschin  staff   926709 Oct 20 16:05 lab1-rs
-rw-r--r--  1 aguschin  staff         2 Oct 20 16:07 test.txt
-rw-r--r--  1 aguschin  staff         7 Oct 20 16:07 test1.gsch1
-rw-----@ 1 aguschin  staff         2 Oct 20 16:07 Tectr_1.txt
$
```

Рис. 1: Сжатие текста Тест_1.txt

```
lab1 -zsh
$ ./lab1-rs
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab1-rs -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab1-rs -i Tectr_6.txt -o test6.gsch1 --compress
$ ./lab1-rs -i test6.gsch1 -o test6.txt --decompress
$ diff Tectr_6.txt test6.txt --report-identical-files
Files Tectr_6.txt and test6.txt are identical
$ ls -l
total 1864
-rwxr-xr-x  1 aguschin  staff   926709 Oct 20 16:05 lab1-rs
-rw-r--r--  1 aguschin  staff    7843 Oct 20 16:08 test6.gsch1
-rw-r--r--  1 aguschin  staff    7958 Oct 20 16:08 test6.txt
-rw-----@ 1 aguschin  staff    7958 Oct 20 16:08 Tectr_6.txt
$
```

Рис. 2: Сжатие текста Тест_6.txt

4 Вычисленные характеристики

4.1 Характеристика 1 (Коэффициент сжатия)

Результаты применения программы к каждому из тестовых текстовых файлов занесены в таблицу 1. Можно заметить, что в некоторых случаях архивация с помощью алгоритма Хаффмана является неэффективной.

Название	Исходный размер, байт	Сжатый размер, байт	Коэффициент
Тест_1.txt	2	7	0.28571
Тест_2.txt	33	92	0.3587
Тест_3.txt	2739	1998	1.37087
Тест_4.txt	330	48	6.875
Тест_5.txt	59	168	0.35119
Тест_6.txt	7958	7843	1.01466
Тест_7.txt	138245	120488	1.14738
Тест_8.txt	574426	516945	1.11119
Тест_9.txt	2752	349	7.88539
Тест_10.txt	2814	368	7.64674

Таблица 1: результаты тестирования

4.2 Характеристика 2 (Скорость сжатия)

Для тестирования скорости сжатия использовался произвольный двоичный файл размера 76450438 байт (≈ 73 мегабайта). В результате пяти последовательных запусков, среднее время запаковки файла составило 4.2 секунды, среднее время распаковки составило 4.7 секунды.

Таким образом, средняя скорость сжатия составила 17.35924 Мбайт в секунду, а средняя скорость разжатия составила 15.51251 Мбайт в секунду.

5 Реализация

Программа реализована на языке программирования Rust с использованием библиотеки clap для чтения параметров командной строки. Сборка производится с помощью программы cargo, поставляющейся вместе с языком.

5.1 Содержимое файла priority_queue.rs

```
1 use std::cmp::Ordering;
2 use std::collections::BinaryHeap;
3
4 #[derive(Debug)]
5 pub struct Weighted<T> {
6     pub weight: u32,
7     pub value: T,
8 }
9
10 impl<T> PartialEq for Weighted<T> {
11     fn eq(&self, other: &Self) -> bool {
```

```

12         self.weight == other.weight
13     }
14 }
15
16 impl<T> Eq for Weighted<T> {}
17
18 impl<T> PartialOrd for Weighted<T> {
19     fn partial_cmp(&self, other: &Self) -> Option<Ordering> {
20         Some(self.cmp(other))
21     }
22 }
23
24 impl<T> Ord for Weighted<T> {
25     fn cmp(&self, other: &Self) -> Ordering {
26         self.weight.cmp(&other.weight).reverse()
27     }
28 }
29
30 #[derive(Debug)]
31 pub struct PriorityQueue<T> {
32     heap: BinaryHeap<Weighted<T>>,
33 }
34
35 impl<T> PriorityQueue<T> {
36     pub fn new() -> Self {
37         PriorityQueue {
38             heap: BinaryHeap::new(),
39         }
40     }
41
42     pub fn insert(&mut self, priority: u32, value: T) {
43         self.heap.push(Weighted {
44             weight: priority,
45             value,
46         });
47     }
48
49     pub fn len(&self) -> usize {
50         self.heap.len()
51     }
52
53     pub fn pop(&mut self) -> Option<Weighted<T>> {
54         Some(self.heap.pop()?)
55     }
56 }

```

5.2 Содержимое файла huffman.rs

```

1 use crate::priority_queue::PriorityQueue;
2 use std::collections::HashMap;
3
4 #[derive(Clone, Copy, Debug)]
5 enum Bit {
6     Zero,

```

```

7     One,
8 }
9 impl Bit {
10     fn from_u8(x: u8) -> Bit {
11         if x == 0 {
12             Bit::Zero
13         } else {
14             Bit::One
15         }
16     }
17
18     fn to_u8(x: &Bit) -> u8 {
19         match *x {
20             Bit::Zero => 0,
21             Bit::One  => 1,
22         }
23     }
24 }
25
26 #[derive(Debug)]
27 enum HuffmanTree<T>
28 where
29     T: std::cmp::Eq + std::hash::Hash + Copy,
30 {
31     Leaf(T),
32     Node(Box<HuffmanTree<T>>, Box<HuffmanTree<T>>),
33 }
34
35 impl<T: std::cmp::Eq + std::hash::Hash + Copy> HuffmanTree<T> {
36     fn from_queue(mut queue: PriorityQueue<T>) -> Self {
37         let mut huf_queue = PriorityQueue::new();
38         while let Some(p) = queue.pop() {
39             huf_queue.insert(p.weight,
↳ Box::new(HuffmanTree::Leaf(p.value)));
40         }
41         assert!(huf_queue.len() > 0);
42         while huf_queue.len() > 1 {
43             let v1 = huf_queue.pop().unwrap();
44             let v2 = huf_queue.pop().unwrap();
45             huf_queue.insert(
46                 v1.weight + v2.weight,
47                 Box::new(HuffmanTree::Node(v1.value, v2.value)),
48             );
49         }
50         return *huf_queue.pop().unwrap().value;
51     }
52
53     fn dfs(tree: &HuffmanTree<T>, mut acc: Vec<Bit>, code: &mut
↳ HashMap<T, Vec<Bit>>) -> Vec<Bit> {
54         match tree {
55             HuffmanTree::Leaf(val) => {
56                 code.insert(*val, acc.clone());
57             }
58             HuffmanTree::Node(left, right) => {

```

```

59         acc.push(Bit::Zero);
60         acc = HuffmanTree::dfs(left, acc, code);
61         acc.pop();
62         acc.push(Bit::One);
63         acc = HuffmanTree::dfs(right, acc, code);
64         acc.pop();
65     }
66 };
67     return acc;
68 }
69
70 fn get_code(&self) -> HashMap<T, Vec<Bit>> {
71     let mut code = HashMap::new();
72     if let HuffmanTree::Leaf(val) = &self {
73         code.insert(*val, vec![Bit::One]);
74     } else {
75         HuffmanTree::dfs(self, Vec::new(), &mut code);
76     }
77     return code;
78 }
79 }
80
81 fn get_weights(data: &Vec<u8>) -> HashMap<u8, u32> {
82     let mut freq = HashMap::new();
83     data.iter().for_each(|byte| {
84         *freq.entry(*byte).or_insert(1) += 1;
85     });
86     let mut freq: Vec<(u8, u32)> = freq.drain().collect();
87     freq.sort();
88     let mut weights = HashMap::new();
89     freq.iter().enumerate().for_each(|(i, (byte, _))| {
90         weights.insert(*byte, i as u32);
91     });
92     return weights;
93 }
94
95 #[derive(Debug)]
96 struct Metadata {
97     weights: HashMap<u8, u32>,
98     tree: HuffmanTree<u8>,
99     code: HashMap<u8, Vec<Bit>>,
100     remainder: u8,
101 }
102
103 impl Metadata {
104     fn compute(data: &Vec<u8>) -> Self {
105         let weights = get_weights(data);
106         let mut queue = PriorityQueue::new();
107         for (value, priority) in weights.iter() {
108             queue.insert(*priority, *value);
109         }
110         let tree = HuffmanTree::from_queue(queue);
111         let code = tree.get_code();
112

```

```

113         return Self {
114             weights,
115             tree,
116             code,
117             remainder: 0,
118         };
119     }
120
121     fn load(data: &Vec<u8>) -> Self {
122         let remainder = data[0];
123         let dict_len = data[1] as usize + 1;
124         let mut weights = HashMap::new();
125         for i in 0..dict_len {
126             weights.insert(data[2 + 2 * i], data[2 + 2 * i + 1] as u32);
127         }
128
129         let mut queue = PriorityQueue::new();
130         for (value, priority) in weights.iter() {
131             queue.insert(*priority, *value);
132         }
133         let tree = HuffmanTree::from_queue(queue);
134         let code = tree.get_code();
135
136         return Self {
137             weights,
138             tree,
139             code,
140             remainder,
141         };
142     }
143
144     fn dump(&self) -> Vec<u8> {
145         let mut result = Vec::new();
146
147         result.push(self.remainder);
148         result.push((self.weights.len() - 1) as u8);
149         for (byte, weight) in &self.weights {
150             result.push(*byte);
151             result.push(*weight as u8);
152         }
153
154         return result;
155     }
156 }
157
158 struct BitWriter {
159     buffer: Vec<Bit>,
160     remainder: u8,
161     result: Vec<u8>,
162 }
163
164 impl BitWriter {
165     fn new() -> Self {
166         Self {

```

```

167         buffer: Vec::new(),
168         remainder: 0,
169         result: Vec::new(),
170     }
171 }
172
173 fn dump_byte(&mut self) {
174     let mut byte = 0_u8;
175     self.buffer
176         .iter()
177         .map(Bit::to_u8)
178         .enumerate()
179         .for_each(|(i, bit)| byte |= bit << i);
180     self.result.push(byte);
181     self.buffer.clear();
182 }
183
184 fn write_bit(&mut self, bit: Bit) {
185     self.buffer.push(bit);
186     if self.buffer.len() == 8 {
187         self.dump_byte();
188     }
189 }
190
191 fn write_bits(&mut self, bits: &Vec<Bit>) {
192     bits.iter().for_each(|bit| self.write_bit(*bit));
193 }
194
195 fn finish(&mut self) {
196     let remainder = 8 - self.buffer.len() as u8 % 8;
197     if remainder != 0 {
198         self.dump_byte();
199     }
200     self.remainder = remainder;
201 }
202 }
203
204 pub fn compress(data: &Vec<u8>) -> Vec<u8> {
205     let mut result = Vec::new();
206     let mut metadata = Metadata::compute(data);
207
208     let mut writer = BitWriter::new();
209     data.iter().for_each(|byte| {
210         let bits = metadata.code.get(byte).unwrap();
211         writer.write_bits(&bits);
212     });
213     writer.finish();
214     metadata.remainder = writer.remainder;
215
216     let md_dump = metadata.dump();
217     md_dump.iter().for_each(|byte| result.push(*byte));
218     writer.result.iter().for_each(|byte| result.push(*byte));
219
220     return result;

```



```

221 }
222
223 struct BitReader<'a> {
224     data: &'a Vec<u8>,
225     metadata: &'a Metadata,
226     buffer: Vec<Bit>,
227     ptr: usize,
228 }
229
230 impl<'a> BitReader<'a> {
231     fn new(data: &'a Vec<u8>, metadata: &'a Metadata) -> Self {
232         Self {
233             data,
234             metadata,
235             buffer: Vec::new(),
236             ptr: 0,
237         }
238     }
239
240     fn read_byte(&mut self) {
241         if let Some(byte) = self.data.get(self.ptr) {
242             for i in 0..=7 {
243                 let bit = byte & (1 << i);
244                 self.buffer.push(Bit::from_u8(bit));
245             }
246             if self.ptr == self.data.len() - 1 {
247                 for _ in 0..self.metadata.remainder {
248                     self.buffer.pop();
249                 }
250             }
251             self.buffer.reverse();
252             self.ptr += 1;
253         }
254     }
255
256     fn read_bit(&mut self) -> Option<Bit> {
257         if self.buffer.len() == 0 {
258             self.read_byte();
259         }
260         self.buffer.pop()
261     }
262 }
263
264 pub fn decompress(archive: &Vec<u8>) -> Vec<u8> {
265     let mut result = Vec::new();
266     let metadata = Metadata::load(archive);
267
268     let data = archive[2 + metadata.weights.len() * 2..].to_vec();
269     let mut reader = BitReader::new(&data, &metadata);
270
271     let mut state = &metadata.tree;
272     while let Some(bit) = reader.read_bit() {
273         if let HuffmanTree::Node(left, right) = state {
274             match bit {

```

```

275         Bit::Zero => state = left,
276         Bit::One => state = right,
277     }
278 }
279
280     if let HuffmanTree::Leaf(byte) = state {
281         result.push(*byte);
282         state = &metadata.tree;
283     }
284 }
285
286     return result;
287 }

```

5.3 Содержимое файла main.rs

```

1  mod huffman;
2  mod priority_queue;
3  use clap::Parser;
4  use std::fs::File;
5  use std::io::{Error, ErrorKind, Read, Write};
6  use std::path::PathBuf;
7
8  #[derive(Parser)]
9  struct Cli {
10     #[arg(short)]
11     input_file: PathBuf,
12
13     #[arg(short)]
14     output_file: PathBuf,
15
16     #[arg(long, default_value_t = true)]
17     compress: bool,
18
19     #[arg(long, default_value_t = false)]
20     decompress: bool,
21 }
22
23 fn run_decompressor(cli: &Cli) -> Result<(), Error> {
24     let mut input_f = File::open(cli.input_file.to_str().unwrap())?;
25     let mut archive = Vec::new();
26     input_f.read_to_end(&mut archive)?;
27     let data = huffman::decompress(&archive);
28
29     let mut output_f = File::create(cli.output_file.to_str().unwrap())?;
30     output_f.write_all(&data)?;
31
32     return Ok(());
33 }
34
35 fn run_compressor(cli: &Cli) -> Result<(), Error> {
36     let mut input_f = File::open(cli.input_file.to_str().unwrap())?;
37     let mut data = Vec::new();
38     input_f.read_to_end(&mut data)?;

```

```

39     let archive = huffman::compress(&data);
40
41     let mut output_f = File::create(cli.output_file.to_str().unwrap())?;
42     output_f.write_all(&archive)?;
43
44     return Ok(());
45 }
46
47 fn main() -> std::io::Result<()> {
48     let cli = Cli::parse();
49
50     let result = if cli.decompress {
51         run_decompressor(&cli)
52     } else {
53         run_compressor(&cli)
54     };
55
56     if let Err(error) = result {
57         match error.kind() {
58             ErrorKind::NotFound => println!("Указанный файл не найден"),
59             ErrorKind::AlreadyExists => println!("Указанный файл уже
↪ существует"),
60             _ => println!("Произошла непредвиденная ошибка"),
61         };
62     }
63
64     return Ok(());
65 }

```

5.4 Содержимое файла Cargo.toml

```

1 [package]
2 name = "lab1-rs"
3 version = "0.1.0"
4 edition = "2021"
5
6 [dependencies]
7 clap = { version = "4.0.17", features = ["derive"] }

```