

МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ  
(НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ)

Институт №8 «Информационные технологии и прикладная математика»  
Кафедра 806 «Вычислительная математика и программирование»

**Курсовой проект  
по курсу «Дискретный анализ»**

Студент: Д. Д. Наумов  
Преподаватель: Н. А. Зацепин  
Группа: М8О-406Б-17  
Дата:  
Оценка:  
Подпись:

**Москва, 2020**

# Тема: Архиватор

## Описание

Задача состоит в реализации архиватора, на основе алгоритма арифметического сжатия. Программа состоит из двух частей: кодера и декодера, которые соответственно сжимают и разжимают входные данные. В качестве входных данных декодеру принимаются текстовые файлы размеров меньше чем  $2^{31}$  бит = 2,6 ГБ. На выходе получается сжатый файл, который подается на вход декодеру.

## Описание алгоритма

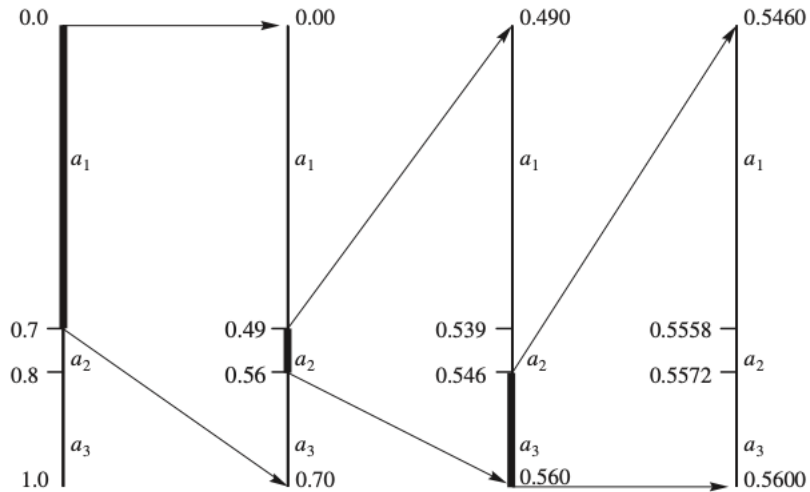
При арифметическом кодировании для кодируемой последовательности создается уникальный идентификатор или тег. Этот тег соответствует двоичной дроби, которая становится двоичным кодом последовательности. Уникальный арифметический код может быть сгенерирован для последовательности длиной  $m$  без необходимости генерировать кодовые слова для всех последовательностей длины  $m$ . Это не похоже на ситуацию с кодами Хаффмана. Чтобы сгенерировать код Хаффмана для последовательности длины  $m$ , где код не является конкатенацией кодовых слов для отдельных символов, нам необходимо получить коды Хаффмана для всех последовательностей длины  $m$ .

Чтобы отличить последовательность символов от другой последовательности символов, нам нужно пометить ее уникальным идентификатором. Один из возможных наборов тегов для представления последовательностей символов - это числа в единичном интервале  $[0, 1)$ . Поскольку количество чисел в единичном интервале бесконечно, должно быть возможно присвоить уникальный тег каждой отдельной последовательности символов.

Процедура создания тега работает за счет уменьшения размера интервала, в котором находится тег, по мере получения все большего количества элементов последовательности.

## Пример

Рассмотрим трехбуквенный алфавит  $A = \{a_1, a_2, a_3\}$  с  $P(a_1) = 0,7$ ,  $P(a_2) = 0,1$  и  $P(a_3) = 0,2$ . Используя отображение  $F_X(1) = 0,7$ ,  $F_X(2) = 0,8$  и  $F_X(3) = 1$  разбиваем единичный интервал, как показано на рисунке ниже.



Для последовательности, начинающейся с  $a_1, a_2, a_3, \dots$ , к тому времени, когда будет получен третий символ  $a_3$ , тег будет ограничен субынтервалом  $[0, 546, 0, 56)$ . Если бы третьим символом был  $a_1$  вместо  $a_3$ , тег находился бы в подынтервале  $[0, 49, 0, 539)$ , который не пересекается с подынтервалом  $[0, 546, 0, 56)$ . Даже если с этого момента две последовательности идентичны (одна начинается с  $a_1, a_2, a_3$ , а другая начинается с  $a_1, a_2, a_1$ ), интервал тегов для двух последовательностей всегда будет непересекающимся.

## Описание структуры программы

Программа разделена на следующие фрагменты:

- **BitIoStream** - поток побитового ввода/вывода
- **FrequencyTable** - таблица с частотами символов
- **ArithmeticCoder** - функции для обновления границ тега
- **ArithmeticCompress** - основная программа для сжатия
- **ArithmeticDecompress** - программа для декодирования

## Пример работы программы

```
1 | (base) MacBook:src dandachok$ ./encoder ../tests/input/1mb.t ../tests/encode/1mb.az
2 | (base) MacBook:src dandachok$ du -h ../tests/input/1mb.t
3 | 980K ../tests/input/1mb.t
4 | (base) MacBook:src dandachok$ du -h ../tests/encode/1mb.az
5 | 568K ../tests/encode/1mb.az
6 | (base) MacBook:src dandachok$ ./decoder ../tests/encode/1mb.az ../tests/output/1mb.t
7 | (base) MacBook:src dandachok$ du -h ../tests/output/1mb.t
8 | 980K ../tests/output/1mb.t
9 | (base) MacBook:src dandachok$ diff ../tests/input/1mb.t ../tests/output/1mb.t
10| (base) MacBook:src dandachok$
```

## Результаты

Файла	Размер файла	Алгоритм	Время сжатия	Время деком-прессии	Размер сжатого файла	Коэффициент сжатия
world95.txt	3005020 Б	Арифметика	0m2.306s	0m2.725s	1920022 Б	1.56
world95.txt	3005020 Б	gzip	0m0.309s	0m4.359s	868440 Б	3.4
enwik8	100000000 Б	Арифметика	1m15.385s	1m29.597s	63502180 Б	1.57
enwik8	100000000 Б	gzip	0m8.411s	0m2.545s	36475811 Б	2.75

- **Процессор** 1,8 GHz 2-ядерный процессор Intel Core i5
- **Память** 4 ГБ 1600 MHz DDR3
- **Графика** Intel HD Graphics 4000 1536 МБ

## Листинг программы

### ArithmeticCoder.hpp

```
1  #pragma once
2
3  #include <algorithm>
4  #include <cstdint>
5  #include "BitInputStream.hpp"
6  #include "FrequencyTable.hpp"
7
8  class ArithmeticCoderBase {
9  public:
10     explicit ArithmeticCoderBase(int numBits);
11     virtual ~ArithmeticCoderBase() = 0;
12
13 protected:
14
15     int numStateBits;
16     std::uint64_t fullRange;
17     std::uint64_t halfRange;
18     std::uint64_t quarterRange;
19     std::uint64_t minimumRange;
20     std::uint64_t maximumTotal;
21     std::uint64_t stateMask;
22     std::uint64_t low;
23     std::uint64_t high;
24
25     virtual void update(const FrequencyTable &freqs, std::uint32_t symbol);
26     virtual void shift() = 0;
27     virtual void underflow() = 0;
28 };
29
30
31 class ArithmeticDecoder final : private ArithmeticCoderBase {
32 public:
33     explicit ArithmeticDecoder(int numBits, BitInputStream &in);
34     std::uint32_t read(const FrequencyTable &freqs);
35
36 protected:
37     void shift() override;
38     void underflow() override;
39
40 private:
41     BitInputStream &input;
42     std::uint64_t code;
43
44     int readCodeBit();
45 };
46
```

```

47 |
48 | class ArithmeticEncoder final : private ArithmeticCoderBase {
49 |
50 | public:
51 |     explicit ArithmeticEncoder(int numBits, BitOutputStream &out);
52 |     void write(const FrequencyTable &freqs, std::uint32_t symbol);
53 |     void finish();
54 |
55 | protected:
56 |     void shift() override;
57 |     void underflow() override;
58 |
59 | private:
60 |     BitOutputStream &output;
61 |     unsigned long numUnderflow;
62 | };

```

### ArithmeticCoder.cpp

```

1 | #include <limits>
2 | #include <stdexcept>
3 | #include "ArithmeticCoder.hpp"
4 |
5 | using std::uint32_t;
6 | using std::uint64_t;
7 |
8 |
9 | ArithmeticCoderBase::ArithmeticCoderBase(int numBits) {
10 |     if (numBits < 1 || numBits > 63)
11 |         throw std::domain_error("State size out of range");
12 |     numStateBits = numBits;
13 |     fullRange = static_cast<decltype(fullRange)>(1) << numStateBits;
14 |     halfRange = fullRange >> 1;
15 |     quarterRange = halfRange >> 1;
16 |     minimumRange = quarterRange + 2;
17 |     maximumTotal = std::min(std::numeric_limits<decltype(fullRange)>::max() / fullRange,
18 |                             minimumRange);
19 |     stateMask = fullRange - 1;
20 |     low = 0;
21 |     high = stateMask;
22 | }
23 |
24 | ArithmeticCoderBase::~~ArithmeticCoderBase() {}
25 |
26 | void ArithmeticCoderBase::update(const FrequencyTable &freqs, uint32_t symbol) {
27 |     if (low >= high || (low & stateMask) != low || (high & stateMask) != high)
28 |         throw std::logic_error("Assertion error: Low or high out of range");
29 |     uint64_t range = high - low + 1;
30 |     if (range < minimumRange || range > fullRange)
31 |         throw std::logic_error("Assertion error: Range out of range");

```



```

31
32 uint32_t total = freqs.getTotal();
33 uint32_t symLow = freqs.getLow(symbol);
34 uint32_t symHigh = freqs.getHigh(symbol);
35 if (symLow == symHigh)
36     throw std::invalid_argument("Symbol has zero frequency");
37 if (total > maximumTotal)
38     throw std::invalid_argument("Cannot code symbol because total is too large");
39
40 uint64_t newLow = low + symLow * range / total;
41 uint64_t newHigh = low + symHigh * range / total - 1;
42 low = newLow;
43 high = newHigh;
44
45 while (((low ^ high) & halfRange) == 0) {
46     shift();
47     low = ((low << 1) & stateMask);
48     high = ((high << 1) & stateMask) | 1;
49 }
50
51 while ((low & ~high & quarterRange) != 0) {
52     underflow();
53     low = (low << 1) ^ halfRange;
54     high = ((high ^ halfRange) << 1) | halfRange | 1;
55 }
56 }
57
58
59 ArithmeticDecoder::ArithmeticDecoder(int numBits, BitInputStream &in) :
60     ArithmeticCoderBase(numBits),
61     input(in),
62     code(0) {
63     for (int i = 0; i < numStateBits; i++)
64         code = code << 1 | readCodeBit();
65 }
66
67 uint32_t ArithmeticDecoder::read(const FrequencyTable &freqs) {
68     uint32_t total = freqs.getTotal();
69     if (total > maximumTotal)
70         throw std::invalid_argument("Cannot decode symbol because total is too large");
71     uint64_t range = high - low + 1;
72     uint64_t offset = code - low;
73     uint64_t value = ((offset + 1) * total - 1) / range;
74     if (value * range / total > offset)
75         throw std::logic_error("Assertion error");
76     if (value >= total)
77         throw std::logic_error("Assertion error");
78
79     uint32_t start = 0;

```

```

80 | uint32_t end = freqs.getSymbolLimit();
81 | while (end - start > 1) {
82 |     uint32_t middle = (start + end) >> 1;
83 |     if (freqs.getLow(middle) > value)
84 |         end = middle;
85 |     else
86 |         start = middle;
87 | }
88 | if (start + 1 != end)
89 |     throw std::logic_error("Assertion error");
90 |
91 | uint32_t symbol = start;
92 | if (offset < freqs.getLow(symbol) * range / total || freqs.getHigh(symbol) * range /
    |     total <= offset)
93 |     throw std::logic_error("Assertion error");
94 | update(freqs, symbol);
95 | if (code < low || code > high)
96 |     throw std::logic_error("Assertion error: Code out of range");
97 | return symbol;
98 | }
99 |
100 | void ArithmeticDecoder::shift() {
101 |     code = ((code << 1) & stateMask) | readCodeBit();
102 | }
103 |
104 | void ArithmeticDecoder::underflow() {
105 |     code = (code & halfRange) | ((code << 1) & (stateMask >> 1)) | readCodeBit();
106 | }
107 |
108 | int ArithmeticDecoder::readCodeBit() {
109 |     int temp = input.read();
110 |     if (temp == -1)
111 |         temp = 0;
112 |     return temp;
113 | }
114 |
115 |
116 | ArithmeticEncoder::ArithmeticEncoder(int numBits, BitOutputStream &out) :
117 |     ArithmeticCoderBase(numBits),
118 |     output(out),
119 |     numUnderflow(0) {}
120 |
121 | void ArithmeticEncoder::write(const FrequencyTable &freqs, uint32_t symbol) {
122 |     update(freqs, symbol);
123 | }
124 |
125 | void ArithmeticEncoder::finish() {
126 |     output.write(1);
127 | }

```

```

128
129 void ArithmeticEncoder::shift() {
130     int bit = static_cast<int>(low >> (numStateBits - 1));
131     output.write(bit);
132
133     // Write out the saved underflow bits
134     for (; numUnderflow > 0; numUnderflow--)
135         output.write(bit ^ 1);
136 }
137
138 void ArithmeticEncoder::underflow() {
139     if (numUnderflow == std::numeric_limits<decltype(numUnderflow)>::max())
140         throw std::overflow_error("Maximum underflow reached");
141     numUnderflow++;
142 }

```

### ArithmeticCompress.cpp

```

1  #include <stdint>
2  #include <stdlib>
3  #include <fstream>
4  #include <iostream>
5  #include <stdexcept>
6  #include <vector>
7  #include "ArithmeticCoder.hpp"
8  #include "BitIoStream.hpp"
9  #include "FrequencyTable.hpp"
10
11 using std::uint32_t;
12
13
14 int main(int argc, char *argv[]) {
15     if (argc != 3) {
16         std::cerr << "Usage: " << argv[0] << " InputFile OutputFile" << std::endl;
17         return EXIT_FAILURE;
18     }
19     const char *inputFile = argv[1];
20     const char *outputFile = argv[2];
21
22     std::ifstream in(inputFile, std::ios::binary);
23     SimpleFrequencyTable freqs(std::vector<uint32_t>(257, 0));
24     freqs.increment(256);
25     while (true) {
26         int b = in.get();
27         if (b == EOF)
28             break;
29         if (b < 0 || b > 255)
30             throw std::logic_error("Assertion error");
31         freqs.increment(static_cast<uint32_t>(b));
32     }

```

```

33
34 in.clear();
35 in.seekg(0);
36 std::ofstream out(outputFile, std::ios::binary);
37 BitOutputStream bout(out);
38 try {
39
40     for (uint32_t i = 0; i < 256; i++) {
41         uint32_t freq = freqs.get(i);
42         for (int j = 31; j >= 0; j--)
43             bout.write(static_cast<int>((freq >> j) & 1)); // Big endian
44     }
45
46     ArithmeticEncoder enc(32, bout);
47     while (true) {
48         int symbol = in.get();
49         if (symbol == EOF)
50             break;
51         if (symbol < 0 || symbol > 255)
52             throw std::logic_error("Assertion error");
53         enc.write(freqs, static_cast<uint32_t>(symbol));
54     }
55
56     enc.write(freqs, 256);
57     enc.finish();
58     bout.finish();
59     return EXIT_SUCCESS;
60
61 } catch (const char *msg) {
62     std::cerr << msg << std::endl;
63     return EXIT_FAILURE;
64 }
65 }

```

### ArithmeticDecompress.cpp

```

1 #include <stdint>
2 #include <cstdlib>
3 #include <fstream>
4 #include <iostream>
5 #include <limits>
6 #include <vector>
7 #include "ArithmeticCoder.hpp"
8 #include "BitIoStream.hpp"
9 #include "FrequencyTable.hpp"
10
11 using std::uint32_t;
12
13
14 int main(int argc, char *argv[]) {

```

```

15 | if (argc != 3) {
16 |     std::cerr << "Usage: " << argv[0] << " InputFile OutputFile" << std::endl;
17 |     return EXIT_FAILURE;
18 | }
19 | const char *inputFile = argv[1];
20 | const char *outputFile = argv[2];
21 |
22 | std::ifstream in(inputFile, std::ios::binary);
23 | std::ofstream out(outputFile, std::ios::binary);
24 | BitInputStream bin(in);
25 | try {
26 |
27 |     SimpleFrequencyTable freqs(std::vector<uint32_t>(257, 0));
28 |     for (uint32_t i = 0; i < 256; i++) {
29 |         uint32_t freq = 0;
30 |         for (int j = 0; j < 32; j++)
31 |             freq = (freq << 1) | bin.readNoEof(); // Big endian
32 |         freqs.set(i, freq);
33 |     }
34 |     freqs.increment(256);
35 |
36 |     ArithmeticDecoder dec(32, bin);
37 |     while (true) {
38 |         uint32_t symbol = dec.read(freqs);
39 |         if (symbol == 256)
40 |             break;
41 |         int b = static_cast<int>(symbol);
42 |         if (std::numeric_limits<char>::is_signed)
43 |             b -= (b >> 7) << 8;
44 |         out.put(static_cast<char>(b));
45 |     }
46 |     return EXIT_SUCCESS;
47 |
48 | } catch (const char *msg) {
49 |     std::cerr << msg << std::endl;
50 |     return EXIT_FAILURE;
51 | }
52 | }

```

## Выводы

Идея арифметического сжатия довольно простая, однако при реализации возникает много нюансов. Из-за ограничений точности типов с плавающей запятой, невозможно просто пересчитывать границы интервала. Необходимо переходить от вещественных значений к целым. Так же при реализации возникают ограничения на максимальный размер входного файла.