EE 248583

Advisor: 電子所 高立人 副教授

**106368002 張昌祺 Justin, Chang-Qi Zhang**

justin840727@gmail.com

Due Date: November 12 2018

## Problem 1 Entropy

Let X be a random variable with an alphabet $H = \{1, 2, 3, 4, 5\}$. Please determine $H(X)$ for the following three cases of probability mass function $p(i) = prob[X = i]$. (15%)

(a) $P(1) = P(2) = 1/2$:

**Ans**

$$
\begin{aligned}
H(X) &= -(P(1)\log_2 P(1) + P(2)\log_2 P(2)) \\
&= -(0.5\log_2(0.5) + 0.5\log_2(0.5)) \\
&= -(-0.5 - 0.5) \\
&= 1 \ bits/symbol
\end{aligned}
$$

(b) $P(i) = 1/4, for\ i = 1, 2, 3, and\ p(4) = p(5) = 1/8$:

**Ans**

$$
\begin{aligned}
H(X) &= -(3 \times P(1)\log_2 P(1) + P(4)\log_2 P(4) + P(5)\log_2 P(5)) \\
&= -(3 \times 0.25\log_2(0.25) + 2 \times 0.125\log_2(0.125)) \\
&= -(-1.5 - 0.75) \\
&= 2.25 \ bits/symbol
\end{aligned}
$$

(c) $P(i) = 2^{-i}, for\ i = 1, 2, 3, 4, and\ p(5) = 1/16$:

**Ans**

$$
\begin{aligned}
H(X) &= -(\sum_{i=1}^{4} 2^{-i}\log_2 2^{-i} + \frac{1}{16}\log_2\frac{1}{16}) \\
&= -(0.5 \times (-1) + 0.25 \times (-2) + 0.125 \times (-3) + 0.0625 \times (-4) + 0.0625 \times (-4)) \\
&= 1.875 \ bits/symbol
\end{aligned}
$$

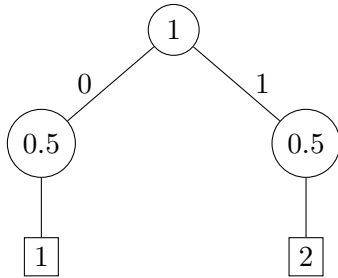# Problem 2 Huffman Code

Design a Huffman code C for the source in Problem 1. (15%)

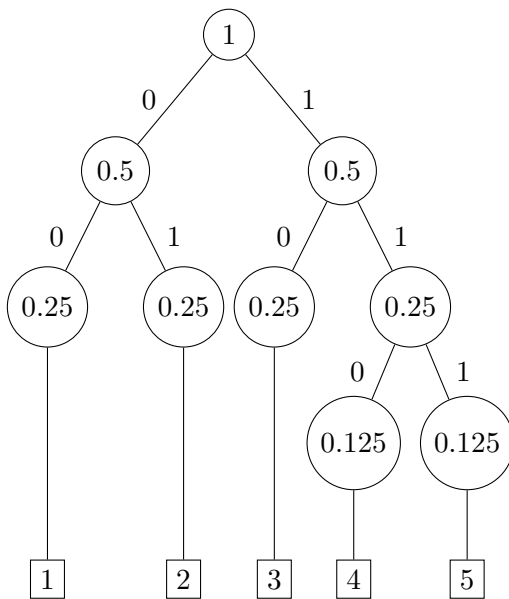(a) Specify your codewords for individual pmf model in Problem 1.
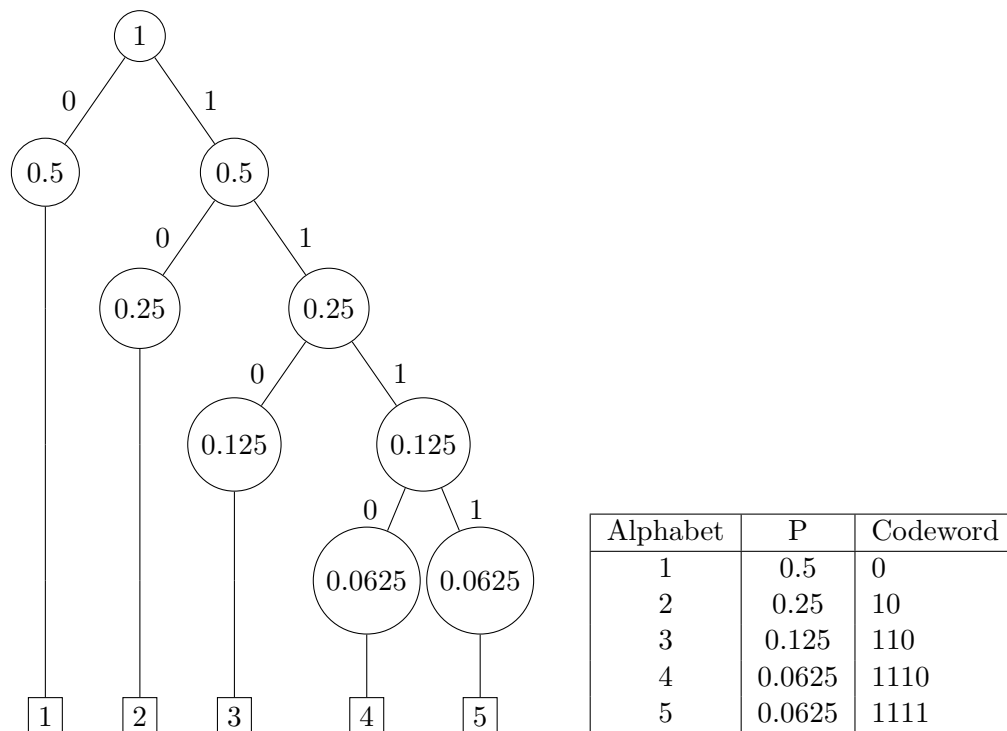
**Ans**

**1.(a)**



| Alphabet | P | Codeword |
|----------|-----|----------|
| 1 | 0.5 | 0 |
| 2 | 0.5 | 1 |

**1.(b)**



| Alphabet | P | Codeword |
|----------|-------|----------|
| 1 | 0.25 | 00 |
| 2 | 0.25 | 01 |
| 3 | 0.25 | 10 |
| 4 | 0.125 | 110 |
| 5 | 0.125 | 111 |

**1.(c)**



| Alphabet | P | Codeword |
|----------|--------|----------|
| 1 | 0.5 | 0 |
| 2 | 0.25 | 10 |
| 3 | 0.125 | 110 |
| 4 | 0.0625 | 1110 |
| 5 | 0.0625 | 1111 |

(b) Compute the expected codeword length and compare with the entropy for your codes in (a).

**Ans**

**1.(b)**

$$Expected\ codeword\ length = 0.5 \times 1 + 0.5 \times 1$$
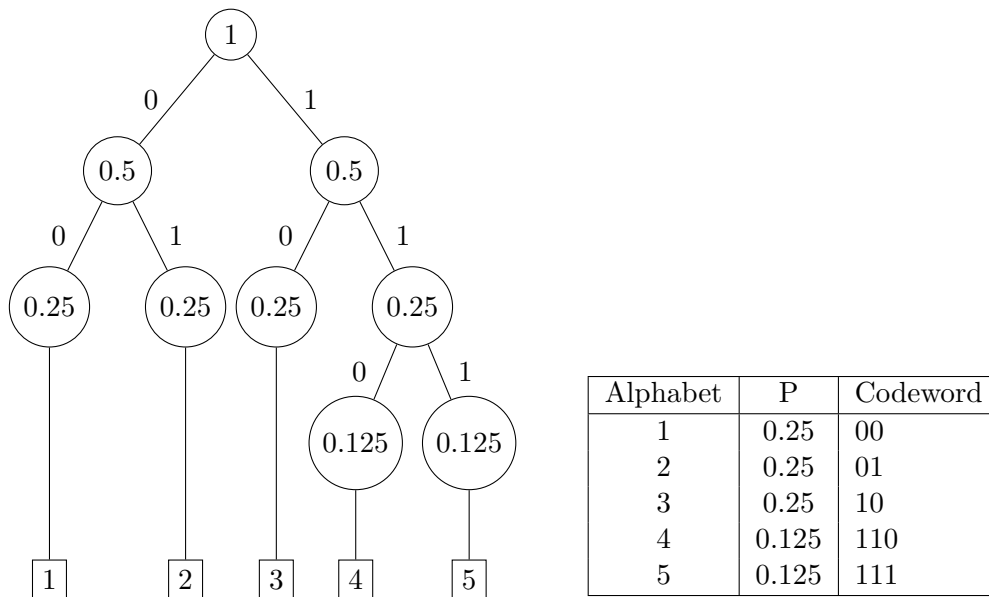$$= 1\ bits/symbol\ \textbf{(Equal Entropy)}$$

**1.(b)**

$$Expected\ codeword\ length = 0.25 \times 2 + 0.25 \times 2 + 0.25 \times 2 + 0.125 \times 3 + 0.125 \times 3$$
$$= 2.25\ bits/symbol\ \textbf{(Equal Entropy)}$$

**1.(c)**

$$Expected\ codeword\ length = 0.5 \times 1 + 0.25 \times 2 + 0.125 \times 3 + 0.0626 \times 4 + 0.0625 \times 4$$
$$= 4.125\ bits/symbol\ \textbf{(NOT Equal Entropy)}$$

(c) Design a code with minimum codeword length variance for the pmf model in Problem 1.(b)

**Ans**



| Alphabet | P | Codeword |
|----------|-------|----------|
| 1 | 0.25 | 00 |
| 2 | 0.25 | 01 |
| 3 | 0.25 | 10 |
| 4 | 0.125 | 110 |
| 5 | 0.125 | 111 |

# Problem 3 Empirical Distribution C++

Empirical distribution. In the case a probability model is not known, it can be estimated from empirical data. Let's say the alphabet is $H = \{1, 2, 3, \dots, m\}$. Given a set of observations of length $N$, the empirical distribution is given by $p = total\ number\ of\ symbol\ 1/N,\ for\ i = 1, 2, 3, ..., m$. Please determine the empirical distribution for **santaclaus.txt**, which is an ASCII file with only lower-cased English letters (i.e., $a \sim z$), space and CR (carriage return), totally 28 symbols. The file can be found on the class web site. Compute the entropy. (14%)

**Ans**

The source code for this problem are available at `https://github.com/justin-changqi/2018_fall_data_compression.git`. Please check README.md to know how to execute the code. After I executed the program the entropy is 4.12 bits/symbol. Empirical distribution shows in Figure 2



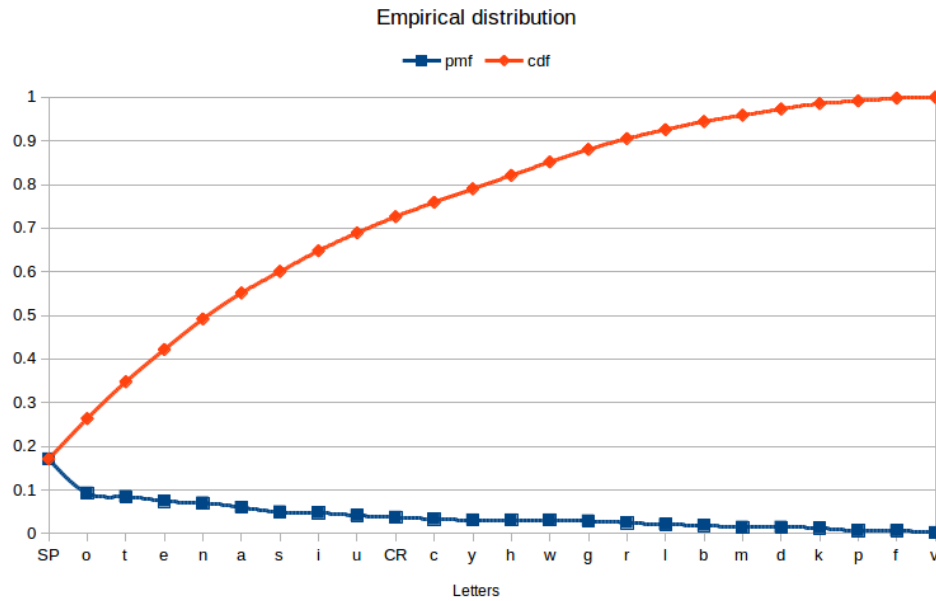Figure 1: Statistics result for **santaclaus.txt**

Figure 2: Empirical distribution for **santaclaus.txt**

# Problem 4 Huffman Code Encode C++

Write a program that designs a Huffman code for the given distribution in Problem 3. (14

**Ans**

The program for this problem was wrote together with Problem 3. The execute print the Huffman encode result as Figure 3



```
================== Codeword Table ==================

    SP:       111          t:      1101
     l:    110011          b:    110010
    CR:     11000          e:      1011
     n:      1010          c:     10011
     y:     10010          h:     10001
     w:     10000          a:      0111
     g:     01101          m:     011001
     d:    011000          f:   01011111
     v:  01011110          p:    0101110
     k:    010110          r:      01010
     s:      0100          o:       001
     i:      0001          u:      0000
```

Figure 3: Huffman encode result for **santaclaus.txt**

5

# Problem 5 Adaptive Huffman Tree

Let X be a random variable with an alphabet $H$ , i.e., the 26 lower-case letters. Use adaptive Huffman tree to find the binary code for the sequence **a a b b a**. (24%)

You are asked to use the following 5 bits fixed-length binary code as the initial codewords for the 26 letters. That is

a: 00000

b: 00001

$\vdots$

z: 11001

**Note**: Show the Huffman tree during your coding process.

**Ans**

1. Initial step:

   $Total\ nodes = 2m - 1 = 26 \times 2 - 1 = 51$

   $\boxed{\text{NYT}}$

   51

2. **a** encoded:

   51

   (1)

   NYT $\boxed{0}$      $\boxed{1}$ a

   49          50

   00000
   a

3. **a a** encoded:

   51

   (2)

   NYT $\boxed{0}$      $\boxed{2}$ a

   49          50

   00000   1
   a       a

6

4. **a a b** encoded:

51
(3)
49    50
(1)    [2] a
47    48
NYT [0]    [1] b

```
00000   1    0    00001
  a     a   NYT     b
```

5. **a a b b** encoded:

51
(4)
49    50
(2)    [2] a
47    48
NYT [0]    [2] b

```
00000   1    0    00001   01
  a     a   NYT     b      b
```

6. **a a b b a** encoded:

51
(5)
49    50
(2)    [3] a
47    48
NYT [0]    [2] b

```
00000   1    0    00001   01   1
  a     a   NYT     b      b    a
```

# Problem 6 Golomb Encoding and Decoding.

(a) Find the Golomb code of n=21 when m=4.

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^2 - 4 = 0$$
$$encoded\ 21 = 21/4 = 5 \ldots 1 = 111110\ 01$$

(b) Find the Golomb code of n=14 when m=4.

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^2 - 4 = 0$$
$$encoded\ 14 = 14/4 = 3 \ldots 2 = 1110\ 10$$

(c) Find the Golomb code of n=21 when m=5.

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^3 - 5 = 3$$
$$encoded\ 21 = 21/5 = 2 \ldots 1 = 110\ 01$$

(d) Find the Golomb code of n=14 when m=5.

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^3 - 5 = 3$$
$$encoded\ 14 = 14/5 = 2 \ldots 4 = 110\ 111$$

(e) A two-integer sequence is encoded by Golomb code with m=4 to get the bitstream 11101111000. What's the decoded two-integer sequence?

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^2 - 4 = 0$$

$$\underline{1110} \quad \underline{11} \quad \underline{110} \quad \underline{00}$$
$$\phantom{11}3 \qquad 3 \qquad 2 \qquad 0$$
$$15 \qquad\qquad 8$$
$$\text{sequence: } 15,\ 8$$

(f) A two-integer sequence is encoded by Golomb code with m=5 to get the bitstream 11101111000 (the same bitstream as that in (e)). What's the decoded two-integer sequence?

**Hint**: The unary code for a positive integer q is simply q 1s followed by a 0.

**Ans**

$$2^{\lceil \log_2^m \rceil} - m = 2^3 - 5 = 3$$

$$\underline{1110} \quad \underline{111} \quad \underline{10} \quad \underline{00}$$
$$3 \qquad 7\text{-}3 = 4 \quad 2 \qquad 0$$
$$18 \qquad\qquad 10$$
$$\text{sequence: } 19,\ 10$$

## Source code for Problem 3 & 4

*huffman_code.hpp*

```cpp
1  #include <iostream>
2  #include <fstream>
3  #include <algorithm>     // std::find
4  #include <vector>        // std::vector
5  #include <iomanip>
6  #include <sstream>
7  #include <string>
8  #include <numeric>
9  #include <cstring>
10
11 class Node
12 {
13   public:
14   std::string letter;
15   char symbol;
16   int cnt;
17   Node *child_r;
18   Node *child_l;
19   Node(std::string letter, char symbol, int cnt);
20   Node(int cnt, Node *child_l, Node *child_r);
21   bool isLeaf();
22   bool operator<(const  Node & other)
23   {
24       return cnt < other.cnt;
25   }
26   bool operator>=(const  Node & other)
27   {
28       return cnt >= other.cnt;
29   }
30   bool operator==(const  Node & other)
31   {
32       return cnt == other.cnt;
33   }
34   int operator+(const  Node & other)
35   {
36       return cnt + other.cnt;
37   }
38 };
39
40 class HuffmanCode {
```

```cpp
41   public:
42     std::vector<Node> nodes;
43     Node *root;
44     std::vector< std::pair <char, std::string> > code_list;
45     double total_letters;
46     HuffmanCode(std::string file, std::string csv_file);
47     void getPmfCdf(std::vector<Node> node_list,
48                    std::vector<double> &pmf,
49                    std::vector<double> &cdf);
50     double getEntropy(std::vector<double> &pmf);
51     void printTable(std::vector<Node> node_list,
52                     std::vector<double> &pmf,
53                     std::vector<double> &cdf);
54     void writeToCsv(std::string file_path,
55                     std::vector<Node> node_list,
56                     std::vector<double> &pmf,
57                     std::vector<double> &cdf);
58     std::string getSymbol(char c);
59     void initNodes(std::vector<char> alphabas, std::vector<int> counts);
60     void buildTree();
61     void mergeNodesToList(std::vector<Node> &list, int start_indx, int end_indx);
62     void printTree(Node *root, int spaces);
63     void encodeData(Node *root, std::string code);
64     bool isLeaf(Node *root);
65     void printCodeWord();
66 };
```

*huffman_code.cpp*

```cpp
1 #include "huffman_code.hpp"
2
3 Node::Node(std::string letter, char symbol, int cnt)
4 {
5    this->symbol = symbol;
6    this->letter = letter;
7    this->cnt = cnt;
8    this->child_l = NULL;
9    this->child_r = NULL;
10 }
11
12 Node::Node(int cnt, Node *child_l, Node *child_r)
13 {
14    this->cnt = cnt;
15    // this->parent = NULL;
16    this->child_l = child_l;
17    this->child_r = child_r;
18 }
19
20 bool Node::isLeaf()
21 {
22    if (this->child_l == NULL && this->child_r == NULL)
23    {
24      return true;
25    }
26    else
27    {
28      return false;
29    }
30 }
31
32 HuffmanCode::HuffmanCode( std::string file, std::string csv_file)
```

```cpp
{
  std::vector<char> alphabas;
  std::vector<int> counts;
  std::ifstream infile;
  char letter[0];
  infile.open (file, std::ios::app);
  while ( infile.peek()  != EOF ) {
    infile.read (letter, 1);
    // counting letter
    std::vector<char>::iterator it;
    it = std::find (alphabas.begin(), alphabas.end(), *letter);
    if (it != alphabas.end())
    {
      size_t index = it - alphabas.begin();
      counts[index] += 1;
    }
    else
    {
      if(letter[0] != 0x0a)
      {
        alphabas.push_back(letter[0]);
        counts.push_back(1);
      }
    }
  }
  infile.close();
  this->initNodes(alphabas, counts);
  // Get reverse list
  std::vector<Node> nodes_r;
  for (int i = this->nodes.size()-1; i >= 0; i--)
  {
    nodes_r.push_back(this->nodes[i]);
  }
  std::vector <double> pmf, cdf;
  this->getPmfCdf(nodes_r, pmf, cdf);
  this->printTable(nodes_r, pmf, cdf);
  std::cout << "\nEntropy: " << this->getEntropy(pmf) << " bits/symbol" << std::
    endl;
  this->writeToCsv(csv_file, nodes_r, pmf, cdf);
  this->buildTree();
}

void HuffmanCode::getPmfCdf(std::vector<Node> node_list,
                            std::vector<double> &pmf,
                            std::vector<double> &cdf)
{
  for (int i = 0; i < node_list.size(); i++)
  {
    pmf.push_back(node_list[i].cnt / this->total_letters);
    if (i == 0)
    {
      cdf.push_back(pmf[0]);
    }
    else
    {
      cdf.push_back(cdf[i-1]+pmf[i]);
    }
  }
}

```

```
92   double HuffmanCode::getEntropy(std::vector<double> &pmf)
93   {
94     double entropy = 0;
95     for (int i = 0; i < pmf.size(); i++)
96     {
97       entropy += pmf[i] * log2(pmf[i]);
98     }
99     return -entropy;
100  }
101
102  void HuffmanCode::printTable(std::vector<Node> node_list,
103                               std::vector<double> &pmf,
104                               std::vector<double> &cdf)
105  {
106    std::cout << std::setw(15) << "Alphaba";
107    for (int i = 0; i < node_list.size() / 2; i++)
108    {
109      std::cout << std::setw(10) << node_list[i].letter;
110    }
111    std::cout << std::endl;
112    std::cout << std::setw(15) << "Total Number";
113    for (int i = 0; i < node_list.size() / 2; i++)
114    {
115      std::cout << std::setw(10) << node_list[i].cnt;
116    }
117    std::cout << std::endl;
118    std::cout << std::setw(15) << "PMF";
119    for (int i = 0; i < node_list.size() / 2; i++)
120    {
121      std::cout << std::setw(10) << std::setprecision (3)<< pmf[i];
122    }
123    std::cout << std::endl;
124    std::cout << std::setw(15) << "CDF";
125    for (int i = 0; i < node_list.size() / 2; i++)
126    {
127      std::cout << std::setw(10) << std::setprecision (3)<< cdf[i];
128    }
129    std::cout << std::endl << std::endl;;
130    std::cout << std::setw(15) << "Alphaba";
131    for (int i = node_list.size() / 2; i < node_list.size(); i++)
132    {
133      std::cout << std::setw(10) << node_list[i].letter;
134    }
135    std::cout << std::endl;
136    std::cout << std::setw(15) << "Total Number";
137    for (int i = node_list.size() / 2; i < node_list.size(); i++)
138    {
139      std::cout << std::setw(10) << node_list[i].cnt;
140    }
141    std::cout << std::endl;
142    std::cout << std::setw(15) << "PMF";
143    for (int i = node_list.size() / 2; i < node_list.size(); i++)
144    {
145      std::cout << std::setw(10) << std::setprecision (3)<< pmf[i];
146    }
147    std::cout << std::endl;
148    std::cout << std::setw(15) << "CDF";
149    for (int i = node_list.size() / 2; i < node_list.size(); i++)
150    {
151      std::cout << std::setw(10) << std::setprecision (3)<< cdf[i];
```

```cpp
152    }
153    std::cout << std::endl;
154 }
155
156 void HuffmanCode::writeToCsv( std::string file_path,
157                                std::vector<Node> node_list,
158                                std::vector<double> &pmf,
159                                std::vector<double> &cdf)
160 {
161    std::ofstream myfile(file_path);
162    myfile << "letter,pmf,cdf" << std::endl;
163    for (int i = 0; i < node_list.size(); i++)
164    {
165      myfile << node_list[i].letter << "," << pmf[i] << "," <<  cdf[i] << std::endl
      ;
166    }
167    myfile.close();
168 }
169
170 std::string HuffmanCode::getSymbol(char c)
171 {
172    switch(c)
173    {
174      case 0x0d:
175        return "CR";
176        break;
177      case 0x20:
178        return "SP";
179        break;
180      default:
181        std::string s;
182        s += c;
183        return s;
184        break;
185    }
186 }
187
188 void HuffmanCode::initNodes(std::vector<char> alphabas, std::vector<int> counts)
189 {
190    for (int i = 0; i < alphabas.size(); i++)
191    {
192      Node n(this->getSymbol(alphabas[i]), alphabas[i], counts[i]);
193      this->nodes.push_back(n);
194    }
195    std::sort(this->nodes.begin(), this->nodes.end());
196    this->total_letters = 0;
197    for (auto& n : this->nodes)
198      this->total_letters += n.cnt;
199 }
200
201 void HuffmanCode::buildTree()
202 {
203    int start_index = 0;
204    while(this->nodes.size()-1 > start_index)
205    {
206      int node_merge_cnt = 1;
207      Node &node = this->nodes[start_index];
208      for (int i = start_index+1; i < this->nodes.size(); i++)
209      {
210        Node node_next =  this->nodes[i];
```

```cpp
211        if (node == node_next)
212        {
213          node_merge_cnt++;
214        }
215        else
216        {
217          break;
218        }
219      }
220      if (node_merge_cnt == 1)
221      {
222        node_merge_cnt = 2;
223      }
224      else
225      {
226        node_merge_cnt = (node_merge_cnt / 2) * 2;
227      }
228      // Merge Node then put into node list
229      this->mergeNodesToList(this->nodes, start_index, start_index+node_merge_cnt);
230      start_index += node_merge_cnt;
231    }
232    this->root = &this->nodes.back();
233  }
234
235  void HuffmanCode::mergeNodesToList(std::vector<Node> &list,
236                                     int start_indx, int end_indx)
237  {
238    for (int i = start_indx; i < end_indx; i = i + 2)
239    {
240      Node *nodes_ptr = this->nodes.data();
241      Node *node_i = nodes_ptr+i;
242      Node *node_i_1 = nodes_ptr+i+1;
243      Node n(list[i]+list[i+1], node_i, node_i_1);
244      if (list.size() > 1)
245      {
246        bool inserted = false;
247        for (int j = 0; j < list.size(); j++)
248        {
249          if (list[j] >= n)
250          {
251            list.insert(list.begin()+j, n);
252            inserted = true;
253            break;
254          }
255        }
256        if (!inserted)
257        {
258          this->nodes.push_back(n);
259        }
260      }
261      else
262      {
263        this->nodes.push_back(n);
264      }
265    }
266  }
267
268  void HuffmanCode::printTree(Node *root, int spaces)
269  {
270    if(root != NULL)
```

```cpp
271    {
272      this->printTree(root->child_r, spaces + 5);
273      for(int i = 0; i < spaces; i++)
274        std::cout << ' ';
275      std::cout << "    " << root->cnt << std::endl;
276      this->printTree(root->child_l, spaces + 5);
277    }
278    else
279    {
280      return;
281    }
282  }
283
284  void HuffmanCode::encodeData(Node *root, std::string code)
285  {
286    if(root != NULL)
287    {
288      this->encodeData(root->child_r, code+"1");
289      this->encodeData(root->child_l, code+"0");
290      if (root->isLeaf())
291      {
292        std::pair <char, std::string> codeword(root->symbol, code);
293        this->code_list.push_back(codeword);
294      }
295    }
296    else
297    {
298      return;
299    }
300  }
301
302  void HuffmanCode::printCodeWord()
303  {
304    int cols = 2;
305    for (int i = 0; i < code_list.size(); i=i+cols)
306    {
307      for (int j = i; j < i + cols; j++)
308      {
309        if (j <  code_list.size())
310        {
311          std::cout << std::setw(10) << this->getSymbol(this->code_list[j].first)
      << ": "
312                    << std::setw(10) << this->code_list[j].second;
313        }
314        else
315        {
316          break;
317        }
318      }
319      std::cout << std::endl;
320    }
321  }
322
323  int main(int argc, char const *argv[])
324  {
325    HuffmanCode huffman_code("../santaclaus.txt", "../statistic_result.csv");
326    // HuffmanCode huffman_code("../test.txt",  "../statistic_result.csv");
327    // HuffmanCode huffman_code("./hw#1_entropy_huffman_golomb/test.txt",  "../hw#1
      _entropy_huffman_golomb/statistic_result.csv");
328    std::cout << std::endl << "═════════════════ Huffman Tree ═════════════════"
```

```
            << std::endl;
329     // HuffmanCode huffman_code("./hw#1_entropy_huffman_golomb/test.txt");
330     huffman_code.printTree(huffman_code.root, 1);
331     std::cout << std::endl << "══════════════════ Codeword Table ══════════════"
            << std::endl << std::endl;
332     huffman_code.encodeData(huffman_code.root, "");
333     huffman_code.printCodeWord();
334     std::cout << std::endl;
335     return 0;
336 }
```