# HUAZHONG UNIVERSITY OF SCIENCE AND TECHNOLOGY

### ADVANCED CLASS 1501

### PROJECTS FOR INFORMATION THEORY

# Implementation of Huffman Coding and Shannon Coding

*Author:*
Haoran YOU

*Supervisor:*
Yayu GAO

June 10, 2017

# Contents

# 1   Introduction

Here we will talk with two Coding Methods: *Huffman Coding* and *Shannon Coding* for the document *Steve Jobs Speech.doc*. It was presented by *Shannon*[1] the first time that the concept of information entropy and the *Shannon's First theorem*, which is called *zero-error coding theorem* for foreseeing the average code length per symbol of an instantaneous code can be made as closes to the entropy, but never be smaller.

# 2   Restatement of the Problem

The objective of this programming project is to deeply understand the Huffman coding and Shannon coding method. Here are the assignment:

- Collect the statistics of the letters, punctuation, space in the document *"Steve Jobs Speech.doc"*.

- Compute the entropy of *"Steve Jobs Speech.doc"*.

- Apply the *Huffman coding method* and *Shannon coding method* for *Steve Jobs Speech.doc*. Output the letters/punctuation/space and their Huffman codewords and Shannon codes. respectively.

- Decode the code which generated by the *Huffman coding method* and *Shannon coding method*, respectively.

- Compute the average code length of *Steve Jobs Speech.doc* using Huffman codes and Shannon codes, respectively.

# 3   Detailed requirements

## 3.1   Collect the statistics of *Steve Jobs Speech.doc*

First we should get the content of this document, the corresponding module in my code is *"Load Data"*. And then we can vote for different symbols in this document to get the frequency of occurrence for each symbol, which corresponding to the module *"Text Statistics"* in my code. The result can be seen from Table1: the number of different characteristics are $71 - symbols$, and the number of total characteristics in this document are $11838 - symbols$.

Table1: Symbols statistics in *Steve Jobs Speech.doc*.

---

[1]who published ≪A Mathematical Theory of Communication≫ in Bell System Technical Journal.1948.

| symbols | Times | Probability | symbols | Times | Probability |
|---|---|---|---|---|---|
|  | 28 | 2.365264e-03 | P | 4 | 3.378949e-04 |
|  | 2228 | 1.882075e-01 | R | 5 | 4.223686e-04 |
| ” | 12 | 1.013685e-03 | S | 22 | 1.858422e-03 |
| $ | 1 | 8.447373e-05 | T | 22 | 1.858422e-03 |
| ’ | 40 | 3.378949e-03 | W | 12 | 1.013685e-03 |
| , | 101 | 8.531847e-03 | X | 3 | 2.534212e-04 |
| - | 9 | 7.602636e-04 | Y | 5 | 4.223686e-04 |
| . | 142 | 1.199527e-02 | a | 741 | 6.259503e-02 |
| 0 | 11 | 9.292110e-04 | b | 121 | 1.022132e-02 |
| 1 | 7 | 5.913161e-04 | c | 215 | 1.816185e-02 |
| 2 | 2 | 1.689475e-04 | d | 408 | 3.446528e-02 |
| 3 | 6 | 5.068424e-04 | e | 1074 | 9.072478e-02 |
| 4 | 1 | 8.447373e-05 | f | 203 | 1.714817e-02 |
| 5 | 2 | 1.689475e-04 | g | 205 | 1.731711e-02 |
| 6 | 2 | 1.689475e-04 | h | 434 | 3.666160e-02 |
| 7 | 5 | 4.223686e-04 | i | 526 | 4.443318e-02 |
| 8 | 1 | 8.447373e-05 | j | 8 | 6.757898e-04 |
| 9 | 2 | 1.689475e-04 | k | 63 | 5.321845e-03 |
| : | 9 | 7.602636e-04 | l | 397 | 3.353607e-02 |
| ; | 2 | 1.689475e-04 | m | 205 | 1.731711e-02 |
| ? | 4 | 3.378949e-04 | n | 588 | 4.967055e-02 |
| A | 31 | 2.618686e-03 | o | 768 | 6.487582e-02 |
| B | 11 | 9.292110e-04 | p | 179 | 1.512080e-02 |
| C | 4 | 3.378949e-04 | q | 4 | 3.378949e-04 |
| D | 9 | 7.602636e-04 | r | 493 | 4.164555e-02 |
| E | 5 | 4.223686e-04 | s | 489 | 4.130765e-02 |
| F | 3 | 2.534212e-04 | t | 906 | 7.653320e-02 |
| G | 2 | 1.689475e-04 | u | 283 | 2.390607e-02 |
| H | 6 | 5.068424e-04 | v | 115 | 9.714479e-03 |
| I | 116 | 9.798953e-03 | w | 233 | 1.968238e-02 |
| J | 1 | 8.447373e-05 | x | 13 | 1.098158e-03 |
| K | 2 | 1.689475e-04 | y | 252 | 2.128738e-02 |
| L | 5 | 4.223686e-04 | z | 4 | 3.378949e-04 |
| M | 11 | 9.292110e-04 | — | 3 | 2.534212e-04 |
| N | 9 | 7.602636e-04 | ¢ | 1 | 8.447373e-05 |
| O | 4 | 3.378949e-04 | 71 | 11838 | 1 |

## 3.2 Compute the entropy of document

The definition formula of Information entropy can be drived as:

$$H(X) = \sum_{i=1}^{n} p_i \log \frac{1}{p_i} \tag{3.1}$$

Where $X$ is a random variable, $p_i$ is its probability distribution. If we view the document as a random variable, we can get it's probability distribution as Table1. So the entropy of this document can be calculated as:

$$H_{Steve-Jobs-Speech} \approx 4.387752 \tag{3.2}$$

## 3.3 Coding for the document

### 3.3.1 Huffman Coding

The fundamental thought of *Huffman Coding* is make the code symbols whose occurrence probability is large as short as possible. So we should sort the symbols according to its occurrence probability and link the two least likely outcome symbols as a new symbol whose occurrence probability equals to the sum of that probabilities of two original symbols.

Recur the symbols until we get two symbols only, so it's easy for us to code them as $'0'$ and $'1'$. Then we find their l-child and r-child, coding them as $'x0'$ and $'x1'$, respectively. Finally we can get the instantaneous code which has the average code length near to information entropy of this document, called *Huffman code*.

Here is our implementation of *Huffman Coding*:

```matlab
 1  function w = Huffman(r, P, S)
 2  %
 3  % Function : Huffman coding
 4  % input    : r --- the number of the source symbols
 5  %            P --- the probability distribution of source symbols
 6  %            S --- the number of source symbols si
 7  % output   : w --- the Huffman codewords wi corresponding to si
 8  %
 9  format long;
10  if(r == 2)
11      w{S(1)} = '0';
12      w{S(2)} = '1';
13  else
14      [P_descend, idx] = sort(P, 'descend');    % sort {Pi} in descending order
15      S = S(idx);                               % sort {Si}
16      P = [P_descend(1:r-2)  (P_descend(r-1) + P_descend(r))];    % update P
17      lchild = S(r - 1);                        % left child
18      rchild = S(r);                            % right child
19      S = S(1:r-1);                             % update S
20      w = Huffman(r-1, P, S);                   % recursion
21      w{rchild} = [w{lchild} '1'];
22      w{lchild} = [w{lchild} '0'];
23  end
```

Code1: Huffman Coding

The result of *Huffman Coding* can be seen as:

Table2: Huffman code for *Steve Jobs Speech.doc*

| Symbols | Pi | Code | Symbols | Pi | Code |
|---------|------|------|---------|------|------|
|  | 2.365264e-03 | 110011010 | 1 | 5.913161e-04 | 11001100110 |
|  | 1.882075e-01 | 01 | 2 | 1.689475e-04 | 1100110010110 |
| " | 1.013685e-03 | 1100110111 | 3 | 5.068424e-04 | 11001100111 |
| $ | 8.447373e-05 | 0001000001111 | 4 | 8.447373e-05 | 00010000011100 |
| ' | 3.378949e-03 | 00010011 | 5 | 1.689475e-04 | 1100110010111 |
| , | 8.531847e-03 | 1100111 | 6 | 1.689475e-04 | 1100110010100 |
| - | 7.602636e-04 | 10010100010 | 7 | 4.223686e-04 | 00010001111 |
| . | 1.199527e-02 | 000101 | 8 | 8.447373e-05 | 00010000011101 |
| 0 | 9.292110e-04 | 0001001000 | 9 | 1.689475e-04 | 1100110010101 |

| Symbols | Pi | Code | Symbols | Pi | Code |
|---|---|---|---|---|---|
| : | 7.602636e-04 | 10010100011 | c | 1.816185e-02 | 110010 |
| ; | 1.689475e-04 | 1001010011010 | d | 3.446528e-02 | 11011 |
| ? | 3.378949e-04 | 100101001010 | e | 9.072478e-02 | 1000 |
| A | 2.618686e-03 | 110011000 | f | 1.714817e-02 | 111010 |
| B | 9.292110e-04 | 0001001001 | g | 1.731711e-02 | 110100 |
| C | 3.378949e-04 | 100101001011 | h | 3.666160e-02 | 11000 |
| D | 7.602636e-04 | 10010100000 | i | 4.443318e-02 | 0011 |
| E | 4.223686e-04 | 00010010110 | j | 6.757898e-04 | 11001100100 |
| F | 2.534212e-04 | 000100011100 | k | 5.321845e-03 | 10010101 |
| G | 1.689475e-04 | 1001010011011 | l | 3.353607e-02 | 11100 |
| H | 5.068424e-04 | 00010000010 | m | 1.731711e-02 | 110101 |
| I | 9.798953e-03 | 1010110 | n | 4.967055e-02 | 0010 |
| J | 8.447373e-05 | 10010100110010 | o | 6.487582e-02 | 1111 |
| K | 1.689475e-04 | 1001010011000 | p | 1.512080e-02 | 111011 |
| L | 4.223686e-04 | 00010010111 | q | 3.378949e-04 | 100101001110 |
| M | 9.292110e-04 | 0001000110 | r | 4.164555e-02 | 10011 |
| N | 7.602636e-04 | 10010100001 | s | 4.130765e-02 | 10100 |
| O | 3.378949e-04 | 100101001000 | t | 7.653320e-02 | 1011 |
| P | 3.378949e-04 | 100101001001 | u | 2.390607e-02 | 00011 |
| R | 4.223686e-04 | 00010010100 | v | 9.714479e-03 | 1010111 |
| S | 1.858422e-03 | 000100001 | w | 1.968238e-02 | 101010 |
| T | 1.858422e-03 | 000100010 | x | 1.098158e-03 | 1100110110 |
| W | 1.013685e-03 | 0001000000 | y | 2.128738e-02 | 100100 |
| X | 2.534212e-04 | 000100011101 | z | 3.378949e-04 | 100101001111 |
| Y | 4.223686e-04 | 00010010101 | – | 2.534212e-04 | 000100000110 |
| a | 6.259503e-02 | 0000 | ¢ | 8.447373e-05 | 10010100110011 |
| b | 1.022132e-02 | 1001011 | | | |

### 3.3.2 Shannon Coding

The construction of *Shannon Code* is very ingenious, it use the conversion of decimal and binary to get instantaneous coding method.

The basic steps of our Shannon Coding Algorithm are:

**Step 1** Sort probability distribution $P(i)$ in descending order.

**Step 2** For i, changed from $1$ to $r$, where $r$ is total number of our symbols, we can calculate the cumulative distribution $F(i)$, and change it to binary description.

**Step 3** Then we should calculate the shortest code length for the probability $P(i)$, denoted as $l(i)$.

**Step 4** Take $l(i)$ digits after the dot of binary description $F(i)$ as the codeword for the source symbol $S(i)$.

Here is our implementation of *Shannon Coding*:

```
1  function w = Shannon(r, P, S)
2  %
3  % Function : Shannon coding
4  % input    : r --- the number of the source symbols
5  %            P --- the probability distribution of source symbols
6  %            S --- the number of source symbols si
7  % output   : w --- the Huffman codewords wi corresponding to si
8  %
9  format long;
10 [P_descend, idx] = sort(P, 'descend');       % sort {Pi} in descending order
11 F = zeros(1,r);                              % initialize {Fi}
12 l = zeros(1,r);                              % initialize {li}
13 for i = 1:r
14     F(i) = sum(P_descend(1:i-1));            % cumulative distribution function
15     l(i) = ceil(log2(1 / P_descend(i)));     % codelength
16     F_binary = dec_to_bin(F(i), l(i));       % convert decimal to binary
17
18     % take li digit after the dot as the codeword for the source symbol idx(i)
19     w{idx(i)} = num2str(F_binary);           % convert to string
20     w{idx(i)} = strrep(w{idx(i)}, ' ',''); % eliminate space
21 end
```

Code2: Shannon Coding

The result of *Shannon Coding* can be seen as:

Table2: Shannon code for *Steve Jobs Speech.doc*

| Symbols | Pi | Code | Symbols | Pi | Code |
| --- | --- | --- | --- | --- | --- |
| | 2.365264e-03 | 111110011 | P | 3.378949e-04 | 111111110010 |
| | 1.882075e-01 | 000 | R | 4.223686e-04 | 111111101010 |
| ” | 1.013685e-03 | 1111101110 | S | 1.858422e-03 | 1111101010 |
| $ | 8.447373e-05 | 11111111111001 | T | 1.858422e-03 | 1111101011 |
| ’ | 3.378949e-03 | 111110000 | W | 1.013685e-03 | 1111101111 |
| , | 8.531847e-03 | 1111010 | X | 2.534212e-04 | 111111110111 |
| - | 7.602636e-04 | 11111100111 | Y | 4.223686e-04 | 111111101100 |
| . | 1.199527e-02 | 1110101 | a | 6.259503e-02 | 0110 |
| 0 | 9.292110e-04 | 11111100010 | b | 1.022132e-02 | 1110110 |
| 1 | 5.913161e-04 | 11111101111 | c | 1.816185e-02 | 110101 |
| 2 | 1.689475e-04 | 1111111110010 | d | 3.446528e-02 | 10110 |
| 3 | 5.068424e-04 | 11111110000 | e | 9.072478e-02 | 0011 |
| 4 | 8.447373e-05 | 11111111111010 | f | 1.714817e-02 | 111000 |
| 5 | 1.689475e-04 | 1111111110100 | g | 1.731711e-02 | 110110 |
| 6 | 1.689475e-04 | 1111111110101 | h | 3.666160e-02 | 10101 |
| 7 | 4.223686e-04 | 111111100101 | i | 4.443318e-02 | 10001 |
| 8 | 8.447373e-05 | 11111111111011 | j | 6.757898e-04 | 11111101110 |
| 9 | 1.689475e-04 | 1111111110111 | k | 5.321845e-03 | 11110111 |
| : | 7.602636e-04 | 11111101001 | l | 3.353607e-02 | 10111 |
| ; | 1.689475e-04 | 1111111111000 | m | 1.731711e-02 | 110111 |
| ? | 3.378949e-04 | 111111101110 | n | 4.967055e-02 | 01111 |
| A | 2.618686e-03 | 111110010 | o | 6.487582e-02 | 0101 |
| B | 9.292110e-04 | 11111100011 | p | 1.512080e-02 | 1110011 |
| C | 3.378949e-04 | 111111101111 | q | 3.378949e-04 | 111111110011 |
| D | 7.602636e-04 | 11111101010 | r | 4.164555e-02 | 10010 |
| E | 4.223686e-04 | 111111100111 | s | 4.130765e-02 | 10011 |
| F | 2.534212e-04 | 111111110110 | t | 7.653320e-02 | 0100 |
| G | 1.689475e-04 | 1111111111001 | u | 2.390607e-02 | 110000 |
| H | 5.068424e-04 | 11111110001 | v | 9.714479e-03 | 1111001 |
| I | 9.798953e-03 | 1110111 | w | 1.968238e-02 | 110011 |
| J | 8.447373e-05 | 11111111111101 | x | 1.098158e-03 | 1111101101 |
| K | 1.689475e-04 | 1111111111011 | y | 2.128738e-02 | 110010 |
| L | 4.223686e-04 | 111111101000 | z | 3.378949e-04 | 111111110100 |
| M | 9.292110e-04 | 11111100101 | — | 2.534212e-04 | 111111111000 |
| N | 7.602636e-04 | 11111101100 | ¢ | 8.447373e-05 | 11111111111110 |
| O | 3.378949e-04 | 111111110000 | | | |

## 3.4 Decoding for the code of document

### 3.4.1 Decoding Huffman Code

After coding the document *Steve Jobs Speech.doc*, we must decode it to verify the correctness of our *Huffman Coding*. Here we load a string from the *'coding Huffman.txt'*, which is the result of *Huffman Coding*, and compare the part of string, form start to next characteristic, to that $71 - symbols$ codes, translate it to the corresponding symbols till the end of string, and find that we translate it to *'Decoding Huffman.txt'*, which is the same as original document!

Here is our implementation of *Decoding huffman code*.

```matlab
1  fids = fopen('Coding_Huffman.txt', 'r');              % get code
2  [A,COUNT] = fscanf(fids, '%c');
3
4  fidss = fopen('Decoding_Huffman.txt', 'wt');          % Decoding Huffman code
5  now = 1;
6  next = 1;
7  while (next ≤ COUNT)
8      idx = 0;
9      for i = 1:num_char
10         if(strcmp(A(now:next), w_H{i}) == 1)          % compare to the codewords
11             idx = i;
12         end
13     end
14     if(idx ≠ 0)
15         fprintf(fidss, '%s', char_doc{idx});          % translate to symbols
16         now = next + 1;
17         next = now + 1;
18     else
19         next = next + 1;                              % +1
20     end
21  end
22  fclose(fids);
```

Code3: Decoding Huffman code

### 3.4.2 Decoding Shannon Code

The process of decoding *Shannon code* is the same as decoding *Huffman code*. The only difference is we must compare our part of string to the *Shannon Codeword* instead of *Huffman Codeword*. The results turns to be alright with our methods.

Here is our implementation of *Decoding Shannon code*.

```matlab
1 fids = fopen('Coding_Shannon.txt', 'r');              % get code
2 [A,COUNT] = fscanf(fids, '%c');
3
4 fidss = fopen('Decoding_Shannon.txt', 'wt');           % Decoding Shannon code
5 now = 1;
6 next = 1;
7 while (next ≤ COUNT)
8     idx = 0;
9     for i = 1:num_char
10        if(strcmp(A(now:next), w_S{i}) == 1)           % compare to the codewords
11            idx = i;
12        end
13    end
14    if(idx ≠ 0)
15        fprintf(fidss, '%s', char_doc{idx});           % translate to symbols
16        now = next + 1;
17        next = now + 1;
18    else
19        next = next + 1;                                % + 1
20    end
21 end
22 fclose(fids);
```

Code4: Decoding Shannon code

## 3.5 Average code length

### 3.5.1 Average codelength for Huffman coding

The average code length for this document with *Huffman Coding* is:

$$l_{Huffman-Coding} \approx 4.428873 \tag{3.3}$$

### 3.5.2 Average codelength for Shannon coding

The average code length for this document with *Shannon Coding* is:

$$l_{Shannon-Coding} \approx 4.774286 \tag{3.4}$$

# 4   Conclusion

From the results of two coding methods, we could find that *Huffman Coding* is more efficient compared to *Shannon Coding*. But neither of them can be less than or equal to the information entropy of this document, which can be seen as the replication experiment for the correctness of *Shannon's First Theorem*!

Through this programming project, I get a deeper understanding of information entropy and the two coding methods: *Huffman Coding* and *Shannon Coding*. This programming experiment can promote my learning of *Information Theorem Course* and the ability for implementing algorithm. Thanks to this programming project, I could implement coding process by myself without extra help except the algorithm provided by teacher.

*Wish I could go further in the information road!*

# References

[1]  Thomas M.Cover , Joy A.Thomas . Elements of Information Theory. Wiley-Blackwell Press.