# INFORMATION THEORY PROJECT REPORT

## Source Coding for Game of Thrones

Name: Zhang Licheng & Sun Binghui

UESTC ID: 2017200602011 & 2017200602023

UoG ID: 2357767 & 2357779

Date: 28. 11. 2019

# Abstract

The purpose of this project is to encode English letters and the space in the English novel Game of Thrones. This experiment was performed by using Huffman coding technique and programming in Python. In this project, counting the frequency of each symbol, creating the Huffman tree, converting the text into binary code and building the executable file were all completed successfully. It concluded that Huffman coding could shorten the number of bits takes to encode a text effectively.
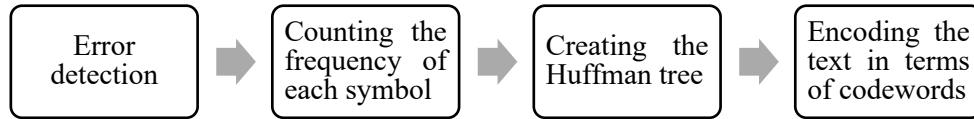
# Preliminary

As a predominant coding method, Huffman coding is a way to encode information using variable-length strings to represent symbols depending on how frequently a particular symbol occurring. By making symbols which are used more frequently be shorter while symbols which appear more rarely be longer, it is proved that Huffman coding always shorten the number of bits takes to encode a given text effectively.

The Huffman Coding Steps:

1) List the frequency of each symbol;

2) Select two symbols with the lowest frequencies;

3) Create a binary tree out of these two symbols, labeling one branch with a "1" and the other with a "0";

4) Add the probabilities of the two symbols to obtain the probability of the new subtree;

5) Remove the symbols from the list and add the subtree to the list;

6) Go back through the list and take the two symbols/subtrees with the smallest probabilities and combine those into a new subtree;

7) Remove the original symbols/subtrees from the list, and add the new subtree to the list;

8) Repeat until all of them are combined.

# Methodology

We divided the whole source coding process into four main stages as below.

| Error detection | Counting the frequency of each symbol | Creating the Huffman tree | Encoding the text in terms of codewords |
|---|---|---|---|

**Stage1: Error detection**

In our project, two kinds of common error are considered, Function *'try…except'* *is utilized to detect* FileNotFound Error and Permission Error.

**Stage2: Counting the frequency of each symbol**

Since each symbol maps to a particular ASCII code, so we set union in terms of ASCII codes to select symbols needed (26 letters and space) and utilized for loop to count the frequency of each symbol. Function 'findall' is used to match each string in the list of strings with chr(i).

**Stage3: Creating the Huffman tree**

Several classes and functions are used in the program to construct the Huffman code tree. They are *'TreeNode'*, *'createnodeQueue'*, *'addQueue'*, *'nodeQueue'*, *'freChar'*, *'createHuffmanTree'* and *'HuffmanCodeDic'*.

- Class *'TreeNode'* defines five built-in variables which represent the letter itself, frequency of occurrence, the 01 code of it, left child and right child during the encoded process;

- Function *'createnodeQueue'* uses *'for..in'* and *'. append'* to make the input list all Class *'nodeQueue'* form;

- Function *'addQueue'* adds the new merged tree nodes into the queue and sorts them;

- Class *'nodeQueue'* has two functions. Function *'addNode'* can add new tree node in sequence and function *'popNode'* deletes the smallest node in the tree;

- Function *'freChar'* inputs a two-dimensional list and returns to a new list which is rearranged according to the second dimension of the original one from big to small;

- Function *'createHuffmanTree'* is used to create Huffman binary tree and it makes a traversal while recording necessary information;
- Function *'HuffmanCodeDic'* is used to store information got from above.

**Stage4: Encoding the text in terms of codewords**

- Function *'TransEncode'* inputs the source novel and translates it into code from;
- Function *'PrintFunction'* prints the code form n words per row.

# Result

To complete this source coding project, we tried several programming platforms including MATLAB, Python and C++, and determined to utilize Python finally. As mentioned above, we divided the whole source coding process into four main parts: error detection, counting the frequency of each symbol, creating the Huffman tree, and encoding the text. After that, we programmed in terms of these four modules in Python, and the result is obtained as shown in *Table1*.

| Symbol | Weight | Codeword | Symbol | Weight | Codeword |
|--------|--------|----------|--------|--------|----------|
| space | 281965 | 111 | a | 97696 | 1001 |
| A | 3000 | 101110100 | b | 18740 | 110 |
| B | 1998 | 100100 | c | 21754 | 111 |
| C | 1455 | 1011100011 | d | 66066 | 11011 |
| D | 1794 | 1011101100 | e | 157340 | 10 |
| E | 978 | 10111011110 | f | 25057 | 101100 |
| F | 814 | 10111010110 | g | 26469 | 101111 |
| G | 1005 | 1001100 | h | 85523 | 111 |
| H | 3974 | 10001 | i | 69083 | 0 |
| I | 5104 | 10111 | j | 762 | 10111010101 |
| J | 1832 | 1011101101 | k | 14042 | 1100111 |
| K | 903 | 10111011100 | l | 52851 | 11000 |

| | | | | | |
|---|---|---|---|---|---|
| **L** | 2780 | 101110000 | **m** | 27948 | 110010 |
| **M** | 1959 | 100000 | **n** | 81760 | 110 |
| **N** | 1982 | 100001 | **o** | 94190 | 1000 |
| **O** | 934 | 10111011101 | **p** | 13911 | 1100110 |
| **P** | 759 | 10111010100 | **q** | 978 | 10111011111 |
| **Q** | 128 | 10111010111001 | **r** | 78829 | 11 |
| **R** | 2096 | 100101 | **s** | 78177 | 10 |
| **S** | 4774 | 10110 | **t** | 98608 | 1010 |
| **T** | 6094 | 10111001 | **u** | 30489 | 110100 |
| **U** | 186 | 1011101011101 | **v** | 8828 | 1010 |
| **V** | 512 | 101110101111 | **w** | 30936 | 110101 |
| **W** | 2436 | 100111 | **x** | 605 | 10011010 |
| **X** | 6 | 101110101110001 | **y** | 25269 | 101101 |
| **Y** | 1404 | 1011100010 | **z** | 609 | 10011011 |
| **Z** | 0 | 101110101110000 | | | |

*Table1. Result of Huffman Coding*

The results in the executable file are shown as follows:



*Figure1. Result of Error Detection in Executable File*

*Figure2. Result of Weight Counting in Executable File*



*Figure3. Result of Huffman Coding in Executable File*

*Figure4. Partial Result of Text Encoding in Executable File*

## Conclusion

During this experiment, we have done further research about Huffman coding and became more familiar with the Python programming language. By utilizing a variety of functions and classes in Python, we completed counting the frequency of each symbol, creating the Huffman tree, and converting the text into binary code successfully. From our test result, it can be seen that Huffman coding technique shortens the transmission bit length effectively.

## Reference

[1] Raymond W. Yeung, *Information Theory and Network Coding*, Springer, 2008;

[2] https://www.cnblogs.com/future-dream/p/10801934.html.

# Appendix

The project code in Python:

```
1.  #Created by Zhang Licheng(2017200602011) & Sun Binghui(2017200602023)
2.
3.  import re
4.
5.  #Error detection
6.  try:
7.      f1=open(r"A_Game_of_Thrones.txt",encoding='utf-8')
8.      s=f1.readlines()  # s is a list of strings
9.      f1.seek(0,0)
10.     f1.close()
11. except FileNotFoundError:
12.     print ("File is not found")
13.     input('Please press ENTER to quit')
14. except PermissionError:
15.     print("You don't have permission")
16.     input('Please press ENTER to quit')
17.
18. a=range(97,123) #a-z
19. b=range(65,91) #A-Z
20. c=range(32,33) #space
21. d=set(a) | set(b) #union
22. e=set(d) | set(c) #union
23.
24. char2=[]
25. count2=[]
26.
27. print ("Symbol\tWeight")
28. for i in e:  # i is int, mapping to ASCII code of each symbol
29.     count = 0
30.     for y in s:
31.         l=re.findall(chr(i),y)  # Utilize findall function to match each string in the list of strings with chr(i)
32.         count += len(l)
33.     count2.append(count)
34.     char2.append(chr(i))
35.     print(chr(i),'\t%d'%count)
36.
37. array=list(zip(char2,count2))
38.
39.
```

```
40. # Create tree node
41. class TreeNode(object):
42.     def __init__(self, data):
43.         self.val = data[0]
44.         self.priority = data[1]
45.         self.leftChild = None
46.         self.rightChild = None
47.         self.code = ""
48.
49. # Create node queue function
50. def createnodeQueue(codes):
51.     q = []
52.     for code in codes:
53.         q.append(TreeNode(code))
54.     return q
55.
56. # Add node to the queue and create priority queue to make sure it is sorted
    from large to small
57. def addQueue(queue, nodeNew):
58.     if len(queue) == 0:
59.         return [nodeNew]
60.     for i in range(len(queue)):
61.         if queue[i].priority >= nodeNew.priority:
62.             return queue[:i] + [nodeNew] + queue[i:]
63.     return queue + [nodeNew]
64.
65. # Define sequence
66. class nodeQueue(object):
67.
68.     def __init__(self, code):
69.         self.que = createnodeQueue(code)
70.         self.size = len(self.que)
71.
72.     def addNode(self, node):
73.         self.que = addQueue(self.que, node)
74.         self.size += 1
75.
76.     def popNode(self):
77.         self.size -= 1
78.         return self.que.pop(0)
79.
80. # Recieve the frequence of letters and space
81. def freChar(array):
82.
```

```python
83.     return sorted(array, key=lambda x: x[1])
84.
85. # Create Huffman Tree
86. def createHuffmanTree(nodeQ):
87.     while nodeQ.size != 1:
88.         node1 = nodeQ.popNode()
89.         node2 = nodeQ.popNode()
90.         r = TreeNode([None, node1.priority + node2.priority])
91.         r.leftChild = node1
92.         r.rightChild = node2
93.         nodeQ.addNode(r)
94.     return nodeQ.popNode()
95.
96. # Obtain Huffman Table in terms of Huffman Tree
97. def HuffmanCodeDic(head, x):
98.     global codeDic, codeList
99.     if head:
100.        HuffmanCodeDic(head.leftChild, x + '0')
101.        head.code += x
102.        if head.val:
103.            codeDic1[head.val] = head.code
104.        HuffmanCodeDic(head.rightChild, x + '1')
105.
106. # Encode string
107. def TransEncode(string):
108.     m=0
109.     global codeDic1
110.     transcode = ""
111.     for ch in string:
112.         if ch in char2:
113.             transcode += codeDic1[ch]
114.     return transcode
115.
116. # Print  200 words per row
117. def PrintFunction(results):
118.     n = 200  # print  200 words per row
119.     for q in range(len(results)):
120.         print(results[q], end=' ')
121.         if (q + 1) % n == 0:
122.             print(' ')
123.
124. codeDic1 = {}
125.
126. t = nodeQueue(freChar(array))
```

```
127. tree = createHuffmanTree(t)
128. HuffmanCodeDic(tree, '')
129.
130. list1= sorted(codeDic1.items(),key=lambda x:x[0])
131.
132. print ("\n")
133. print ("Symbol\tHuffman code")
134.
135. for head, x in list1:
136.     print(head,'\t',x)
137.
138. print('The encoded novel is:')
139.
140. test='a4 b+;;;'
141. list3=''.join(s)
142. PrintFunction(TransEncode(list3))
143. print(' ')
144. input('Please press ENTER to quit')
```