

# EE3980 Algorithms

## Homework 7. Huffman Code

**Due: Apr. 29, 2018**

The American Standard Code for Information Interchange (ASCII) is a standard encoding method for English letters and symbols. It is a fixed length encoding scheme, which means each letter takes a fixed number of bits to represent the letter. Since we know that some of the letters are not used as often as others, this fixed length encoding scheme may not be the most efficient in information storage or exchange.

To increase the storage efficiency, the variable length encoding scheme can be adopted. In this scheme, a less frequently used letter can be encoded with more bits, while more frequently used letters with fewer number of bits. In the end, the total number of bits for information storage should be minimized.

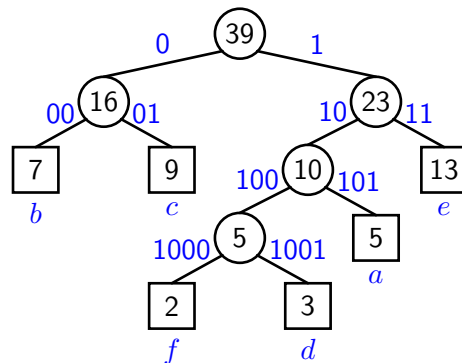
Let  $\mathcal{A} = \{a_1, a_2, \dots, a_n\}$  be the set of letters used, and  $\mathcal{F} = \{f_1, f_2, \dots, f_n\}$  be the corresponding frequency of the letters. In ASCII encoding, each letter takes 8 bits. Thus the total number of bits to store the information is

$$S_{ASCII} = \sum_{i=1}^n 8 \times f_i = 8 \sum_{i=1}^n f_i. \quad (7.1)$$

In a variable length encoding scheme, we assume the number of bits for each letters are  $\mathcal{L} = \{\ell_1, \ell_2, \dots, \ell_n\}$ , then the total number of bits is

$$S_{VL} = \sum_{i=1}^n \ell_i \times f_i. \quad (7.2)$$

The Huffman encoding scheme has been shown to achieve minimum storage requirement for any given set of letters,  $\mathcal{A}$ , and the corresponding frequencies,  $\mathcal{F}$ . As an example, suppose 6 letters,  $\mathcal{A} = \{a, b, c, d, e, f\}$  are used to store some information and the corresponding frequencies are also known,  $\mathcal{F} = \{5, 7, 9, 3, 13, 2\}$ . Then the Binary Merge Tree algorithm (5.3.4) constructs the following binary tree such that  $\sum_{i=1}^n d_i f_i$  is minimum, where  $d_i$  is the depth of the leave node and  $f_i$  is the frequency.



From the tree, we get the following Huffman encoding table.

letter	code
a	101
b	00
c	01
d	1001
e	11
f	1000

The goal of this homework is to construct the Huffman code given a set of words and to show the ratio between the Huffman encoded storage and the standard ASCII storage. The 9 wordlist files of hw06 should be used as test cases. (Note that the first line of each file, the number of words in that file, needs not be considered for this homework). The output of your program taking w101.dat as input is shown at the end of this file.

#### Notes.

1. One executable and error-free **C** source file should be turned in. This source file should be named as hw07.c.
2. A pdf file is also needed. This report file should be named as hw07a.pdf.
3. Submit your hw07.c and hw07a.pdf on EE workstations using the following command:  

```
$ ~ee3980/bin/submit hw07 hw07.c hw07a.pdf
```

where hw07 indicates homework 7.

4. Your report should be clearly written such that I can understand it. The writing, including English grammar, is part of the grading criteria.

Example program output:

---

```
$/a.out < w101.dat
Number of words: 40
Number of characters: 376
Huffman coding
c: 0000
t: 0001
g: 00100
d: 00101
u: 0011
m: 01100
v: 011010
b: 0110110
```

k: 0110111  
n: 0111  
\n: 010  
s: 1000  
o: 1001  
r: 1010  
h: 10110  
p: 10111  
i: 1100  
y: 110100  
f: 1101010  
q: 110101100  
x: 110101101  
w: 11010111  
l: 11011  
e: 1110  
a: 1111

Number of encoded bits: 1600

Ratio: 53.1915%

