**EE3980 Algorithms**

**HW7 Huffman Code**

**104061212  馮立俞**

**2018/4/26**

## Introduction

ASCII character encoding is one commonly-used English encoding system which requires 8 bits (i.e. 1 Byte) for each character. However, it's not very efficient in space usage since not every character appear equally often. In fact, we could compress space usage by encoding frequently-used characters with shorter bits, and less-frequently-used with longer bits. Such encoding method is called Huffman Coding.

## Approach

The target of Huffman Coding is minimizing total bit usage, $\sum_{i=1}^{n} b_i * f_i$, where $b_i$ and $f_i$ are bits required to encode and the usage frequency of that character. From Algorithm class we know that a binary merge tree could minimize $\sum_{i=1}^{n} d_i * f_i$, where $d_i$ and $f_i$ are depth of a node and the node's frequency. Therefore, we could achieve Huffman coding by building a binary merge tree by relating bit length with its level in tree.

## Building Binary Merge Tree

```
1.  Algorithm Tree(n, list) // Generate binary merge tree from list of n files.
2.  {
3.         for i: = 1 to(n- 1) do {
4.             pt: = new node;pt - > lchild: = Least(list);
5.                 // Find and remove min from list.
6.             pt - > rchild: = Least(list);
7.             pt - > w: =
8.               (pt - > lchild) - > w + (pt - > rchild) - > w;Insert(list, pt);
9.                 }
10.        return Least(list);
11. }
```
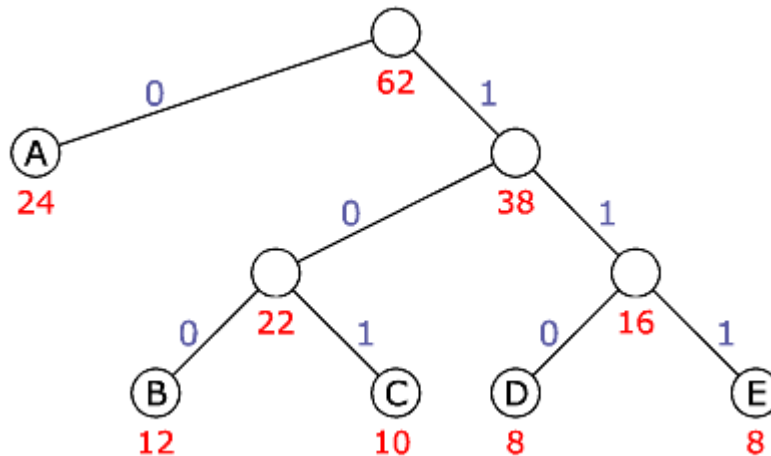
To build a binary merge tree, we keep finding the smallest element in `list,` merging them to form a subtree, then insert the subtree back to `list,` the procedure continues till there's only one element (i.e. root) in `list`.

Depending on how we achieve `Least` and `Insert` function, the resulting tree could be different. I adopted a Min-Heap in this assignment. Due to its unstable sorting property, the nodes with same frequency may switch places, compared to stable sorts. However, the total bit usage isn't affected by how we perform `Least` and `Insert` function. After all, their frequencies are the same.

`Least` and `Insert` function would also affect the complexity of this algorithm. Here because I used Heap, the overall time complexity is $O(nlogn)$, space complexity being $O(n)$ for a Heap and a Tree.

## Encoding

After the tree is built, we encode the left child node with additional '0', and the

right child with additional '1', like the following figure shows. This operation can be

done with recursion call. In practice, I also record bit usage when encoding.



```
1.  Algorithm int Encode(Node * node, char * str) {
2.      bitCount: = 0;node - > HMcode = str;
3.      if (is leaf node) { //print leaf nodes
4.          print( node - > c, node - > HMcode);
5.          bitCount = strlen(node - > HMcode) * node - > freq; //record freq
6.      }
7.      if (node - > lchild != NULL) {
8.          temp = str + '0';
9.          bitCount += Encode(node - > lchild, temp);
10.     }
11.     if (node - > rchild != NULL) {
12.         temp = str + '1';
13.         bitCount += Encode(node - > rchild, temp);
14.     }
15.     return bitCount;
16. }
```

The above recursion call needs to traverse all nodes and edges, whose

time complexity is thus $O(n + e)$, $n, e$ are the number of nodes and

edges in the tree, respectively. Additionally, the recursive nature would require $O(l)$ space complexity when calling themselves.

## Results and analysis

In the given test case, Huffman Code proved to compress bit usage down to around 53% of ASCII version. It seems that the test cases have similar character distribution.

## Observations and Conclusion

**1.Huffman Code compresses space while preserving some readability.**

Intuitively, one might consider another encoding method. That is, '0' for the most common character, '1' for the second, '00' for the third, '01' for the fourth…etc.

Huffman is not optimal?

Though this will save more bits than Huffman code, but it will cause some difficulty decoding it. For example, is '00' representing one or two character, or it's just a fraction of an encoded character? Yet to decode Huffman code, one only need to follow the binary merge tree. And the confusion above can be avoided if there's no noise during message transmission.
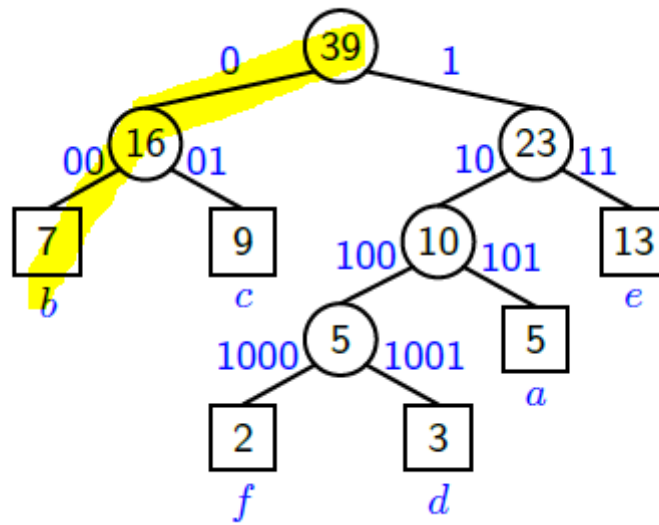
Figure. One only need to follow the tree to decode, say '00'

**2.Depending on the adopted sorting method, the resulting Huffman code may not**

**be unique.**

```
 1 /**************************************
 2   EE3980 Algorithms HW07 Huffman Code
 3   Li-Yu Feng 104061212
 4   Date:2018/4/23
 5 **************************************/
 6
 7 #include <stdio.h>
 8 #include <stdlib.h>
 9 #include <string.h>
10 #include <sys/time.h>
11 #include <stdbool.h>
12 #include <math.h>
13
14 #define LEN 14
15
16 typedef struct node{                //node in MinHeap & Tree
17     struct node *lchild,*rchild;
18     char c, *HMcode;
19     int freq;
20 }Node;
21
22
23 double GetTime(void);
24 void Heapify(Node *list, int i, int n);     // build a heap for nlgn speed-
25 void MinHeap(Node *list,int n);             // finding minimal member
26 Node *HeapPop(Node *list, int n);           //return minimum
27 void HeapInsert(Node *list, int n, Node *new);  //insert new member to MinHeap
28 Node *Tree(Node *list, int n);      //build minimum cost tree
29 int Encode(Node *node, char *str);  // Huffman encoder, print using recursion
30
31 double GetTime(void)
32 {
33     struct timeval tv;
34     gettimeofday(&tv,NULL);
35     return tv.tv_sec+1e-6*tv.tv_usec;
36 }
37
38 void Heapify(Node *list, int i, int n){
39     int j = i*2;
40     Node temp = list[i-1];
41     bool done = false;
42
43     while(j<=n && !done){
44         if(j<n && list[j-1].freq > list[ j+1 -1].freq ) j++;        //list[j
   -1].freq > list[ j+1 -1].freq
45         if(temp.freq < list[j-1].freq ) done = true;
46         else{
47             list[j/2-1] = list[j-1];
48             j *= 2;
49         }
```

1

```
50          //printf("%d\n",j);
51      }
52      list[j/2-1] = temp;
53  }
54
55  void MinHeap(Node *list,int n){
56      Node temp;
57      int i;
58
59      for( i = n/2 ; i>0 ; i--)
60          {Heapify(list,i,n);}
61  }
62
63  Node *HeapPop(Node *list, int n){
64      Node *min;
65
66      min = malloc(sizeof(Node));
67      *min = list[0];
68      list[0] = list[n-1];
69      Heapify(list,1, n-1);    //maintain MinHeap
70
71      return min;
72  }
73
74  void HeapInsert(Node *list, int n, Node *new){        //n = current member count
75      int j = n+1;
76
77      while( j > 1 && new->freq < list[ j/2 -1].freq ){                  //new->fr
    eq < list[ j/2 -1].freq
78          list[j-1] = list[j/2 -1];
79          j /= 2 ;
80      }
81      list[j-1] = *new;
82
83  }
84
85  Node *Tree(Node *list, int n){
86      Node *pt,*temp;
87      int i;
88
89      for( i = n; i > 1; ){
90
91          pt = malloc(sizeof(Node));
92          pt->c = '.';
93
94          for(temp = HeapPop(list, i--);                        //ignore node with

95              temp->freq ==0; temp = HeapPop(list, i--) );   //zero freq
96          pt->lchild = temp;
97          pt->rchild = HeapPop(list, i--);
```

```
 98
 99         pt->freq = pt->lchild->freq + pt->rchild->freq;
100
101         HeapInsert(list, i++, pt);
102         //printf("Take %c(%d) %c(%d)\n",pt->lchild->c,pt->lchild->freq,pt->rchi
    ld->c,pt->rchild->freq );
103         //printf("I am %d\n",pt->freq );
104     }
105     return HeapPop(list,1);
106 }
107
108 int Encode(Node *node, char *str){
109     char *temp;
110     int bitCount = 0;
111
112     node->HMcode = malloc( sizeof(str) );
113     temp = malloc( sizeof(str) +1 );
114     node->HMcode = str;
115     if(node->c != '.' && node->freq != 0){       //print leaf nodes
116         if(node->c == '\0') printf("  \\n:  %s\n",node->HMcode );
117         else    printf("  %c:  %s\n",node->c, node->HMcode );
118         bitCount = strlen(node->HMcode) * node->freq;
119     }
120
121     strcpy(temp,str);
122
123     if (node->lchild != NULL){
124         temp[strlen(temp)] = '0';
125         bitCount += Encode(node->lchild, temp);
126     }
127
128     if (node->rchild != NULL){
129         temp[strlen(temp)-1] = '1';
130         bitCount += Encode(node->rchild, temp);
131     }
132     return bitCount;
133 }
134
135 void InsertionSort(Node *list,int n){
136     int i,j;
137     Node temp;
138
139     for(j = 1; j < n; j++){
140         temp = list[j];
141         i = j-1;
142         while((i>=0) && temp.freq > list[i].freq ){
143             list[i+1] = list[i];
144             i--;
145         }
146         list[i+1] = temp;
```

```
147     }
148 }
149
150
151
152
153 int main(){
154     int i,j,k;
155     int Nwords,Nchar;
156     double t;
157     char **words;
158     Node *list,*temp;
159     bool done;
160
161     int count;
162
163     scanf("%d", &Nwords);
164
165
166     words = (char**)malloc(Nwords * sizeof(char*));          //
167     for(i = 0; i < Nwords; i++)                              //
168         words[i] = (char *)malloc( (LEN+1) * sizeof(char)); //
169
170
171     for(i = 0; i < Nwords ; i++){                            //
172         scanf("%s", words[i]);                               //scan words
173     }
174
175     list = (Node *)malloc( 27 * sizeof(Node));
176     list[0].c = '\0';
177     list[0].freq = 0;
178
179     for (i = 1; i < 27; i++){        //init list
180         list[i].c = 'a'-1 +i;
181         list[i].freq = 0;
182     }
183     Nchar = 0;
184     list[0].freq = Nwords;
185     for (i = 0; i < Nwords ; i++){                //count char freq
186         for(j = 0; words[i][j] != '\0'; j++)
187             list[ words[i][j]-'a'+1 ].freq++;
188         Nchar += j+1;
189     }
190
191     printf("Number of words: %d\nNumber of characters: %d\n",Nwords,Nchar );
192
193     /*InsertionSort(list,27);
194     for ( i = 0,j = 1, k = 2, count = 0; i < 27; ++i)
195     {
196         if(i+1 <= k)   count += list[i].freq * j;
```

4

```
197        else{
198            j++; k *= 2;
199
200                count += list[i].freq * j;
201        }
202    }
203    printf("dummy encoding bit count:%d\n",count );*/
204
205
206
207
208    printf("Huffman coding\n");
209
210    MinHeap(list,27);
211
212    temp = Tree(list,27);
213    i = Encode(temp,"");
214    printf("Number of encoded bits: %d\n", i );
215    printf("Ratio: %.4f%\n", i / (Nchar*8.0)*100.0 );
216    return 0;
217
218 }
```

Score: 60

[Compression] ratios should be listed in your report.

[Segmentation] fault may appear when running your compiled program.
    - node->lchild, node->rchild not properly initialized?

[Compiler warning] line 215.

[Program] can have more comments.

[Algorithms] need to be explained more clearly.

[Report] can be improved.