

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. }
```

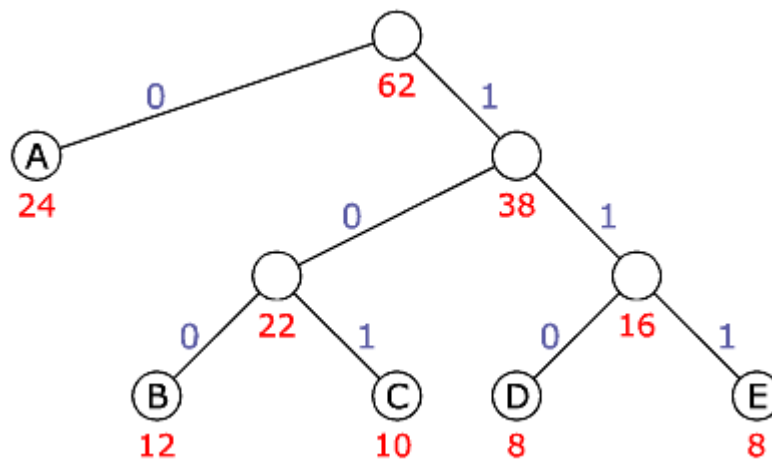
To build a binary merge tree, we keep finding the smallest element in `list`, merging them to form a subtree, then insert the subtree back to `list`, the procedure continues till there's only one element (i.e. root) in `list`.

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 $O(n \log n)$, 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. 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.

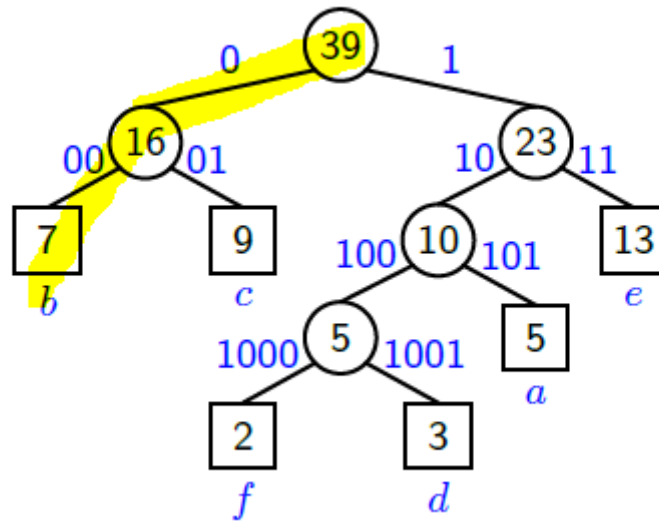


Figure. One only need to follow the tree to decode, say '00'

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