# EE 3980 Algorithms

# Homework 9. Encoding ASCII Texts

105061110 周柏宇

2020/5/16

## 1. Introduction

In this homework, we construct the Huffman code given the paragraph and show the codeword of each character and the ratio of the number of bytes required to store the paragraph using the Huffman code and standard ASCII code.

## 2. Analysis & Implementation

### 2.1 Overview

To construct the Huffman encoding table, what we need to do is first obtain the frequency for every character appearing in the paragraph and construct the binary merge tree. Once we obtain the binary merge tree, we can traverse the tree to retrieve the code.

### 2.2 Obtain Character Frequencies

In this phase, we will go through all characters in a paragraph and accumulate their frequencies, which involves a lot of searching the frequency of a character and increasing them by one. Therefore, it will be good to use a binary search tree to store the characters using their ASCII

code as a comparison key. Here's the data structure for this task.

```c
1. // data structure to store the char and its frequency
2. struct node {
3.     int ch; // character
4.     int n_ch; // frequency
5.     struct node *l; // left child
6.     struct node *r; // right child
7. };
```

The algorithm used to get the character frequencies is as below.

```
1.  // Calculate character frequency
2.  // Input: a paragraph P, binary search tree root bst
3.  // Output: characters and their frequency in a binary
4.  //          search tree, number of characters read
5.  Algorithm getFrequency(P, bst)
6.  {
7.      n_ch := 0; // number of chars read
8.      ch := next character in P;
9.      // read the paragraph
10.       while (ch is not EndOfFile) do {
11.           n_ch := n_ch + 1; // # of chars read plus 1
12.           tmp := bst_find(ch); // find the char in bst
13.           if (tmp = NULL) {
14.               // add the char if not in bst
15.               bst_insert(ch);
16.           }
17.           else {
18.               // increase the frequency of the char by 1
19.               tmp→n_ch := tmp→n_ch + 1;
20.           }
21.           ch := next character in P;
22.       }
23.
24.       return n_ch;
25.  }
```

The while loop will execute as many times as the number of characters in the paragraph, denoted $n$. In the while loop, we always perform a search in the binary search tree, which takes $\mathcal{O}(\lg n)$ on average and $\mathcal{O}(n)$ for the worst case. Besides, we sometimes need to perform insertion to a binary search tree, which has the same average and worst-case time complexity as searching. Therefore, one iteration of the while loop has an average and worst-case time complexity $\mathcal{O}(\lg n)$ and $\mathcal{O}(n)$ respectively, making `getFrequency` an algorithm with time complexity $\mathcal{O}(n \lg n)$ and $\mathcal{O}(n^2)$ for average and worst case.

## 2.3   Binary Merge Tree Algorithm

In this phase, we already have the character frequencies. Thus, we can perform the Binary Merge Tree algorithm to get the optimal merge order.

```
1. // Generate binary merge tree from list of n nodes
2. // which contains the character and its frequency
3. // Input: int n, list of nodes
4. // Output: optimal merge order
5. Algorithm Tree(n, list)
6. {
7.      for i := 1 to (n - 1) do {
8.          pt := new node;
9.          // find and remove min from list
10.         pt→lchild := Least(list);
11.         pt→rchild := Least(list);
12.         pt→n_ch := (pt→lchild)→n_ch + (pt→rchild)→n_ch;
13.         Insert(list, pt);
14.     }
15.     return Least(list);
16. }
```

As we can see, the algorithm above involves finding minimum in a list.

An intuitive way of implementing `Least` is to enforce the min heap property to the list. Therefore, retrieve the minimum will simply be an $\mathcal{O}(1)$ operation.

However, our nodes are stored in a binary search tree right now. We need to first transfer all the nodes to an array then enforce the min heap property.

```
1. // Store the binary search tree in a list.
2. // Input: current node node, list arr, index i
3. // Output: index i
4. // Initiate the call with
5. //         bst_to_array(bst, arr, 1)
6. Algorithm bst_to_array(node, arr, i)
7. {
8.     if (node != NULL) {
9.         arr[i] := node; // add the node to array
10.         i := i + 1;
11.          // travel to left child
12.          i := bst_to_array(node→l, arr, i);
13.          // travel to right child
14.          i := bst_to_array(node→r, arr, i);
15.         // become a leaf node for merging the tree later
16.          node→l := NULL;
17.          node→r := NULL;
18.     }
19.     return i;
20. }
```

```
1. // Make the array a min heap.
2. // Input: array A with n elements
3. // Output: array A with min heap property
4. Algorithm array_to_minHeap(A, n)
5. {
6.     // Enforce min heap property to non leaf node
7.     // in bottom-up order.
8.     for i := ⌊n / 2⌋ to 1 step -1 do
9.         minHeapify(A, i, n);
10. }
```

In bst_to_array, it is a basic tree traversal operation; hence it has

$\mathcal{O}(N)$ time complexity, where $N$ denotes the number of nodes in the binary

search tree, i.e. the number of unique characters in the paragraph.

In `array_to_minHeap`, it is exactly the same as the first loop of heap sort except that heap sort uses `maxHeapify`. Therefore, this function has $O(N \lg N)$ time complexity.

With our list possessing min heap property, the function `Least` and `Insert` in the algorithm `Tree` should be implemented as follows respectively:

```
1. // Remove minimum from the min heap.
2. // Input: min heap in an array A with n elements
3. // Output: minimum of the min heap
4. Algorithm minHeapRemoveMin(A, n)
5. {
6.     if (A = NULL) return NULL; // empty heap
7.     min := A[1]; // minimum is the root
8.     A[1] := A[n]; // move last node to root
9.     minHeapify(A, 1, n - 1); // recover min heap
10.    return min;
11. }
```

```
1. // Insert a node to min heap.
2. // Input: min heap in an array A with n elements
3. //        new node to be inserted nn
4. // Output: none
5. Algorithm minHeapInsertion(A, n, nn)
6. {
7.      i := n; // start at last node
8.     A[n] := nn; // put new node at last node
9.     while ((i > 1) and (A[⌊i / 2⌋]→n_ch > nn→n_ch)) do {
10.           A[i] := A[⌊i / 2⌋]; // parent should be larger
11.           i := ⌊i / 2⌋; // move up one layer
12.        }
13.     A[i] := nn; // put new node at proper place
14. }
```

In minHeapRemoveMin, although retrieving the minimum is a constant time operation, we still need to restore the min heap property since we will reuse the min heap. Therefore, the time complexity is dominated by minHeapify, which is $\mathcal{O}(\lg N)$.

In minHeapInsertion, the while loop will execute at most $\mathcal{O}(\lg N)$ since the index is divided by 2 every time, making the insertion an $\mathcal{O}(\lg N)$ operation.

With the above analysis, we can go back to obtain the time complexity of Tree $T_{\text{Tree}}$:

$$T_{Tree} = (N - 1) \cdot (2T_{\text{Least}} + T_{\text{Insert}})$$

$$= (N - 1) \cdot \left(2 \cdot \mathcal{O}(lgN) + \mathcal{O}(lgN)\right)$$

$$= \mathcal{O}(N \ lgN)$$

Recall that `bst_to_array` has $\mathcal{O}(N)$ and `array_to_minHeap`

has $\mathcal{O}(N \ lgN)$ time complexity. Therefore, in this phase, the time

complexity is $\mathcal{O}(N \ lgN)$, where $N$ is the number of unique characters and

is pretty much a constant for an English paragraph with moderate length.

## 2.4  Retrieve Huffman Code

Now that we have the binary merge tree that produces the optimal code.

What is left to do is to parse the tree and retrieve the code for every character.

```
1. // Print Huffman code.
2. // Input: node in the merge tree node,
3. //          array for the code code
4. //          i-th bit for the code bit
5. // Output: none
6. // Initiate the call with
7. //    printHuffmanCode(root of merge tree, code, 1, -1)
8. Algorithm printHuffmanCode(node, code, i, bit)
9. {
10.      if (bit != -1) { // not initial call
11.          code[i - 1] := bit; // set the bit
12.      }
13.      if (node→ch != -1) { // leaf nodes
14.          write(node→ch) // print chars
15.          // print the code for the char
16.          for i := 1 to i - 1 do {
17.              write(code[j]);
18.          }
19.      }
20.      else { // non leaf nodes
21.          // left child, next bit is 0
22.          printHuffmanCode(node→l, code, i + 1, 0);
23.          // right child, next bit is 1
24.          printHuffmanCode(node→r, code, i + 1, 1);
25.      }
26. }
```

printHuffmanCode is again a traversal algorithm of the binary merge tree. The fact that the number of nodes in the merge tree is $2N - 1$, where $N$ is the number of tree nodes, can be observed by considering that every time we create a new node using two existing nodes, and will not stop if there is more than one node left. Thus, the process will create $N - 1$ merged nodes in addition to $N$ leaf nodes. This results a $\mathcal{O}(N)$ time

complexity.

Besides traversing the tree, when being at leaf nodes, we will print out the corresponding Huffman code, which can be proved to have a length $log_2 \frac{1}{f}$, where $f$ is the normalized frequency of the character. Therefore, to print out all the codes, it will take $\sum_{i=1}^{N} log_2 \frac{1}{f_i}$. Consider the case that $f$ is uniform, i.e., $f_i = \frac{1}{N}$, then $\sum_{i=1}^{N} log_2 \frac{1}{f_i} = N \, log_2 N$. While in general the frequencies we get from a random paragraph will not have a uniformly distributed characters, I think $\mathcal{O}(N \, lgN)$ can be our compromised approximation since $\sum_{i=1}^{N} lg \frac{1}{f_i}$ is unbounded.

With the analysis above, we can conclude that the time complexity of `printHuffmanCode` is $\mathcal{O}(N) + \mathcal{O}(N \, lgN) = \mathcal{O}(N \, lgN)$

## 2.5  Summary

To summarize, we list the time complexity of all phases.

**Obtain Character Frequencies**: $\mathcal{O}(n \, \lg n)$ (average)

$$\mathcal{O}(n^2) \text{ (worst-case)}$$

**Binary Merge Tree Algorithm**: $\mathcal{O}(N \, lgN)$

**Retrieve Huffman Code**: $\mathcal{O}(N \, lgN)$

Generally speaking, $n \gg N$. As a result, the time complexity of the procedure of constructing the Huffman code including collecting the

character frequencies is dominated by the first phase.

# 3. Result and Observation

To observe the relationship between code length l and the normalized frequency

p, I print out some meaningful numeric.

```
Letter   p      log₂(p) l code
----------------------------
' ': 0.17474 2.51670 3 111
e:   0.09047 3.46645 3 001
t:   0.07616 3.71488 4 1100
o:   0.06469 3.95028 4 1000
a:   0.06218 4.00738 4 0111
n:   0.04938 4.34003 4 0101
i:   0.04410 4.50294 4 0001
r:   0.04168 4.58460 4 0000
≈
K:   0.00017 12.54460 13 1010111110100
6:   0.00017 12.54460 13 1010111111100
4:   0.00008 13.54460 13 0100011100110
8:   0.00008 13.54460 13 0100011100111
$:   0.00008 13.54460 14 10101111101010
U:   0.00008 13.54460 14 10101111101011
Number of Chars read: 11949
  Huffman Coding needs 53830 bits, 6729 bytes
  Ratio = 56.3143 %
```

We can see that the code length is bounded by $\log_2 \frac{1}{p} + 1$, which is the optimal

solution for minimizing the average code length $\sum_{i=1}^{N} p_i l_i$, subject to the Kraft

Inequality $\sum_{i=1}^{N} 2^{-l_i} \leq 1$. Kraft Inequality is the necessary condition for a binary prefix

code e.g. Huffman code. It is good to have the constraint because a prefix code can be

decoded without reference to future codewords.