

EE3980 Algorithms

hw09 Encoding ASCII Texts

106061146 陳兆廷

Introduction:

In this homework, I will be analyzing, implementing, and observing 1 algorithm.

The goal of the algorithm is to encode a list of characters into Huffman codes. The input of them will be an essay, and the output of them will be a list of characters used in the essay and their Huffman code. Furthermore, the ratio between the space used to store ASCII and Huffman will be calculated.

During the analysis process, I will first introduce why Heap Sort and Binary Merge Tree can be used in this task. Then, I will be using counting method to calculate the time complexities of the algorithm. Finally, I will calculate their space complexity for the total spaces used by the algorithm.

The implementation of the algorithm on C code will find the Huffman Code for the provided data, course.dat.

Analysis:

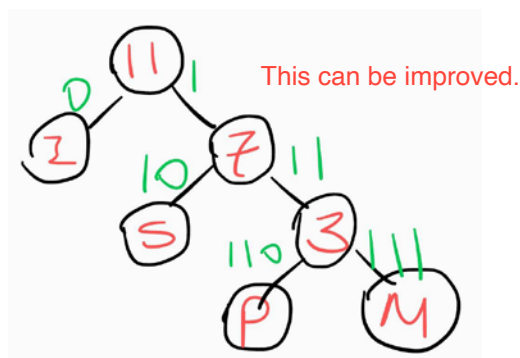
1. Huffman Code:

Huffman encoding is a way of representing a set of characters based on their appearing frequencies. Take the word Mississippi for an example. The data

table will be each existed character and frequency table will store their appearing frequencies.

	[0]	[1]	[2]	[3]
data	M	i	s	p
freq	1	4	4	2

After Sorting characters depending on their frequencies, we implement Binary Merge Tree method to merge the least two elements until there is only one element left, and that will be the root of our Huffman Tree. Take Mississippi as an example again. After sorting it into ascending order, we get M and P and merge frequency $2 + 1 = 3$ into the tree. By doing so until there is only element 11 left.



Finally, we traverse the Huffman Tree to print out our Huffman Encoding Table.

We use Heap sort as the sorting method during Huffman encoding process since there's sorting and it constantly pops the minimum value. Using heap sort ideally has $O(n \lg n)$ and $O(1)$ for the two tasks.

2. Data Structure:

a. Node:

A Node has 4 properties:

properties	definition
Node->(char)data	character
Node->(int)freq	character's frequency
Node->(Node) left	Left node
Node->(Node) right	Right node

b. Heap:

A Heap is to store the sorted character/frequency array and store the merged Binary merge tree.

properties	definition
Heap->(int)size	character
Heap->(Node)*array	Sorted array

3. Heap Sort:

a. Abstract:

Heap sort is a sorting method derives from Heap tree. With heap-representation of the array, heap sort can easily sort input array by *Heapify*, which is an algorithm that maintain heap-property on specific node.

This method of Heapify is slightly different from the original one. It constantly Heapify the child item of the inserted node (i). The process is similar to the old one that it constantly get the heaped first item and

Heapify the rest of them.

b. Algorithm:

```
1. // To enforce min heap property for n-element Heap with root i.
2. // Input: size n max heap array Heap->array, root i
3. // Output: updated Heap.
4. Algorithm Heapify(list, i, size)
5. {
6.     s := i;
7.     left := 2 * i + 1;
8.     right := 2 * i + 2;
9.
10.    if (left < size && list[left]->freq < list[s]->freq) {
11.        s := left;
12.    }
13.    if (right < size && list[right]->freq < list[s]->freq) {
14.        s := right;
15.    }
16.    if (s != i) {
17.        t := list[s];
18.        list[s] := list[i];
19.        list[i] := t;
20.        Heapify(list, s, size);
21.    }
22. }
```

```
1. Algorithm HeapSort(list, n)
2. {
3.     for i := (n - 1) / 2 to 0 step -1 do {
4.         Heapify(list, i, size);
5.     }
6. }
```

c. Time Complexity:

For the Heapify process, it traverse and heapify until the input node is a leaf node, therefore the time complexity will be $O(\lg n)$. For the HeapSort process, it does Heapify for $n/2$ times, therefore the total time complexity will be $n/2 * O(\lg n)$ equals **$O(n \lg n)$** .

d. Space Complexity:

The algorithm uses several integers and Node(t) and a Heap array list[n]. **The space complexity would be $O(N)$.**

4. Binary Merge Tree, BMT():

a. Abstract:

BMT() gets the least two element, which are characters that has least two frequencies, and merge them into one element with frequencies added.

b. Algorithm:

```
1. // Generate binary merge tree from list of n files.
2. // Input: int n, list of files
3. // Output: optimal merge order.
4. Algorithm BMT(n, list)
5. {
6.     for n > 0 do {
7.         pt := new node ;
8.         pt -> lchild := GetMin(list) ; // Find and remove min from list.
9.         pt -> rchild := GetMin(list) ;
10.        pt -> w := (pt -> lchild) -> w + (pt -> rchild) -> w ;
11.        while (n && pt -> freq < list[(n - 1) / 2]->freq) {
12.            list[n] =list[(n - 1) / 2];
13.            n := (n - 1) / 2;
14.        }
```

```

15.         list[n] := pt;
16.     }
17.     return GetMin(list) ;
18. }
19.
20. Algorithm GetMin(list)
21. {
22.     temp := list[0];
23.     list[0] := list[size - 1];
24.     size--;
25.     Heapify(list, 0, size);
26.     return temp;
27. }

```

c. Time complexity:

For GetMin(), there is one Heapify in it, therefore it's time complexity will be $O(\lg n)$.

For BMT(), there is GetMin and a while loop with $O(\lg n)$ frequency (since it iterates with $n, n/2, n/4 \dots$ etc). The time complexity for BMT() will be **$O(n \lg n)$** .

d. Space Complexity:

The algorithm uses several integers and Nodes(temp) and a Heap array list[n]. **The space complexity would be $O(N)$.**

5. Huffman Tree Traversal, PrintCode():

a. Abstract:

PrintCode() prints each input character's Huffman Code by traversing

the Huffman Tree.

b. Algorithm:

```
1. // Generate Huffman Code from Heap Tree.
2. // Input: Node Node, int code[], int top
3. // Output: Huffman Codes
4. Algorithm PrintCode(Node, code, top)
5. {
6.     if Node is a leaf do {
7.         for i := 0 to n do {
8.             print(code[i]);
9.         }
10.    }
11.    if (node->left) {
12.        code[top] = 0;
13.        printCodes(node->left, code, top + 1);
14.    }
15.    if (node->right) {
16.        code[top] = 1;
17.        printCodes(node->right, code, top + 1);
18.    }
19. }
```

c. Time Complexity:

The printing process prints n character's Huffman Code and traverse the Huffman tree. The time complexity would be **$O(n)$** .

d. Space Complexity:

The algorithm uses several integers, code array and a Heap tree. **The space complexity would be $O(N)$.**

6. Overall Algorithm, Huff():

a. Abstract:

Gathering all elements in the above analysis, we can construct an algorithm to implement Huffman Code Encoding.

b. Algorithm:

```
1. // Generate Huffman Code from Heap Tree.
2. // Input: Node Node, int code[], int top
3. // Output: Huffman Codes
4. Algorithm Huff(data, freq, size)
5. {
6.     for i := 0 to size do {
7.         HeapTree->array[i] := newNwode(data[i], freq[i])
8.     }
9.
10.    HeapTree = HeapSort(HeapTree, size);
11.    root = BMT(HeapTree);
12.    PrintCodes(root);
13. }
```

c. Time Complexity:

The time complexity for Huff() will be:

$O(n)$ for initializing.

$O(n \lg n)$ for Heap Sorting.

$O(n \lg n)$ for setting Binary Merge Tree.

$O(n)$ for printing Huffman Codes.

The total time complexity is **$O(n \lg n)$** .

d. Space Complexity:

The algorithm uses the spaces those algorithms I analyzed above

takes. **The space complexity would be $O(N)$.**

7. Time & Space:

	<i>Huff()</i>
Time complexity	$O(N \lg N)$
Space complexity	$O(N)$

Implementation:

1. Result:

	ASCII (Bytes)	Huffman (bits)	Huffman (Bytes)	Ratio (%)
talk1.txt	11949	53830	6729	56.3143
talk2.txt	17654	79601	9951	56.3668
talk3.txt	15944	70812	8852	55.5193
talk4.txt	11014	50144	6268	56.9094
talk5.txt	11802	53813	6727	56.9988

Average Ratio: 56.42172 %

Observation:

1. Ratio between storing ASCII and Huffman codes:

Using Huffman encoding to store a list of characters saves **56.42172 %** of the space needed comparing with storing ASCII.

Conclusions:

1. It takes Heap Sort, Binary Merge Tree to implement a Huffman encoding process.
2. Time and space complexities of *Huff()*:

	<i>Huff()</i>
--	---------------

Time complexity	$O(N \lg N)$
Space complexity	$O(N)$

3. The average ratio between storing ASCIIs and Huffman codes is 56.42172 %.

hw09.c

```
1 // EE3980 HW09 Encoding ASCII Texts
2 // 106061146, Jhao-Ting, Chen
3 // 2020/05/14
4
5 #include <stdio.h>
6 #include <stdlib.h>
7 #include <string.h>
8 #include <sys/time.h>
9
10 typedef struct sNode {                // store new node
11     char data;
12     int freq;
13     struct sNode *left, *right;
14 } Node;
15
16 typedef struct sHeap {                // store new heap
17     int size;
18     Node** array;
19 } Heap;
20
21 void readInput(void);                 // read all inputs
22 void printInput(void);                // print Input
23 Node* newNode(char data, int freq);   // initialize a new node
24 void Heapify(Heap* Heap, int idx);    // Heapify function
25 Heap* HeapSort(Heap* Heap, int n);    // Heap Sort function
26 Node* GetMin(Heap* Heap);             // get min from heap and remove it
27 Node* BMT(Heap* HeapTree);           // Binary Merge Tree
28 void printCodes(Node* root, int *code, int top); // Print Huff code
29 void Huff(char *data, int *freq, int size); // Main Huff function
30
31 char *data;                          // charactor array
32 int *freq;                           // frequency array
33 int N, charnum, bitnum;               // calculate # of characters and bits
34
35 int main()
36 {
37     readInput();
38     //printInput();
39     Huff(data, freq, N);
40
41     return 0;
42 }
43
44 void readInput(void)                  // read all inputs
45 {
46     char c;                          // initialize
47     int i, found = 0;
```

```

48     N = 0;
49     charnum = 0;
50     bitnum = 0;
51     while ((c = getchar()) != EOF) {           // before text ends
52         if (N == 0) {                         // first text
53             data = (char *)malloc(1);
54             freq = (int *)malloc(1);
55             data[0] = c;
56             freq[0] = 1;
57             N++;
58         } else {                               // if existed
59             for (i = 0; (i < N); i++) {
60                 if (data[i] == c) {           // store and add frequency
61                     found = 1;
62                     freq[i]++;
63                 }
64             }
65             if (found == 0) {                 // if new character
66                 data = (char *)realloc((data), (N + 1)* sizeof(char));
67                 freq = (int *)realloc((freq), (N + 1)* sizeof(int));
68                 data[N] = c;                  // initialize
69                 freq[N] = 1;
70                 N++;
71             }
72             found = 0;
73         }
74         charnum++;
75     }
76 }
77
78 void printInput(void)                        // print character and frequency list
79 {
80     int i;
81
82     for (i = 0; i < N; i++) {
83         printf("%c: %d times\n", data[i], freq[i]);
84     }
85 }
86
87 Node* newNode(char data, int freq)           // initialize a new node
88 {
89     Node* tmp = (Node*)malloc(sizeof(Node));
90
91     tmp->left = NULL;
92     tmp->right = NULL;
93     tmp->data = data;
94     tmp->freq = freq;
95
96     return tmp;
97 }

```

```

98
99 void Heapify(Heap* Heap, int i)           // Heapify function
100 {
101     int s = i;
102     int l = 2 * i + 1;                    // left child
103     int r = 2 * i + 2;                    // right child
104     Node* t;
105
106     if (l < Heap->size && Heap->array[l]->freq < Heap->array[s]->freq) {
107         s = l;                            // left smaller
108     }
109     if (r < Heap->size && Heap->array[r]->freq < Heap->array[s]->freq) {
110         s = r;                            // right smaller
111     }
112     if (s != i) {                         // swap with smaller one
113         t = Heap->array[s];
114         Heap->array[s] = Heap->array[i];
115         Heap->array[i] = t;
116         Heapify(Heap, s);                 // heapify children
117     }
118 }
119
120 Heap* HeapSort(Heap* Heap, int n)         // calling Heapify
121 {
122     int i;
123     for (i = (n - 1) / 2; i >= 0; i--) {
124         Heapify(Heap, i);
125     }
126     return Heap;
127 }
128
129 Node* GetMin(Heap* Heap)                  // get the minimum node and remove it
130 {
131     Node* temp = Heap->array[0];           // get minimum
132
133     Heap->array[0] = Heap->array[Heap->size - 1]; // get the biggest
134     Heap->size--;                          // remove one
135     Heapify(Heap, 0);                     // Heapify again
136
137     return temp;                          // return min
138 }
139
140 Node* BMT(Heap* HeapTree)                 // Binary Merge Tree function
141 {
142     Node *left, *right, *mid;             // Initialize
143     int m;
144
145     while (HeapTree->size != 1) {
146
147         right = GetMin(HeapTree);         // get left child

```

```

148     left = GetMin(HeapTree);           // get right child
149
150     mid = newNode('\0', left->freq + right->freq); // create merged node
151
152     mid->left = left;                    // assign child to merged node
153     mid->right = right;
154
155     HeapTree->size++;
156     m = HeapTree->size - 1;              // remove 2 element and add 1
157
158     while (m && mid->freq < HeapTree->array[(m - 1) / 2]->freq) {
159         HeapTree->array[m] = HeapTree->array[(m - 1) / 2];
160         m = (m - 1) / 2;                // find position and insert merged node
161     }
162     HeapTree->array[m] = mid;
163 }
164
165 return GetMin(HeapTree);                // return the root of HeapTree
166 }
167
168 void printCodes(Node* node, int *code, int top) // Print Huff Codes
169 {
170     int i, b = 0;
171     if (!(node->left) && !(node->right)) {      // if it is a leaf
172         if (node->data == '\n') {              // start printing
173             printf(" '\n': ");
174         } else if (node->data == ' ') {
175             printf(" ' ': ");
176         } else {
177             printf(" %c: ", node->data);
178         }
179
180         for (i = 0; i < top; i++) {
181             printf("%d", code[i]);
182             b++;
183         }
184         printf("\n");
185
186         for (i = 0; i < N; i++) {              // calculate # of bits
187             if (data[i] == node->data) {
188                 bitnum = bitnum + b * freq[i];
189             }
190         }
191     }
192
193     if (node->left) {                          // if traverse to left child
194         code[top] = 0;                        // code add 0
195         printCodes(node->left, code, top + 1);
196     }
197

```

```

198     if (node->right) {                // if traverse to right child
199         code[top] = 1;                // code add 1
200         printCodes(node->right, code, top + 1);
201     }
202 }
203
204 void Huff(char *data, int *freq, int size)        // main Huff function
205 {
206     Heap* HeapTree;                    // Initialize a Heap tree
207     Node *root;                        // initialize a root node
208     int i;
209     int code[100], top = 0, mod;        // code to store huff codes
210
211     HeapTree = (Heap*)malloc(sizeof(Heap));      // open space
212     HeapTree->size = size;
213     HeapTree->array = (Node**)malloc(size * sizeof(Node*));
214
215     for (i = 0; i < size; ++i) {            // assigned original array
216         HeapTree->array[i] = newNode(data[i], freq[i]);
217     }
218
219     HeapTree = HeapSort(HeapTree, size - 1);    // sort heap array
220
221     root = BMT(HeapTree);                    // create binary merge tree
222
223     printf("Huffman coding:\n");
224
225     printCodes(root, code, top);              // print huffman codes
226
227     printf("Number of Chars read: %d\n", charnum);
228     mod = bitnum % 8;
229     if (mod) mod = 1;
230     printf(" Huffman Coding needs %d bits, %d bytes\n",
231           bitnum, bitnum / 8 + mod);
232     printf(" Ratio = %g%%\n",
233           (float)(bitnum / 8 + mod) / (float)charnum * 100);
234
235 }

```

[Program Format] can be improved.

[Coding] hw09.c spelling errors: agian(1), charactor(1)

[CPU time] 0.218094 sec

[Report] should explain what is n and how are A and F constructed.

Score: 90