

# 浙江大学



## Huffman Codes

Advanced Data Structures and Algorithm Analysis  
Research Project 5

Wang Zhongwei  
王中伟

Lin Jiafeng  
林家丰

Jin Yuruo  
金雨若

May 13, 2019

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Problem Description . . . . .	2
1.2	Background of Project . . . . .	2
1.2.1	Min Heap . . . . .	2
1.2.2	Greedy Algorithm . . . . .	2
<b>2</b>	<b>Algorithm Specification</b>	<b>3</b>
2.1	Definition . . . . .	3
2.1.1	function . . . . .	3
2.1.2	Node . . . . .	3
2.1.3	CF . . . . .	3
2.1.4	CC . . . . .	4
2.2	Huffman coding . . . . .	4
2.2.1	Huffman tree . . . . .	4
2.2.2	Greedy algorithm . . . . .	5
2.2.3	Correctness of the algorithm . . . . .	5
<b>3</b>	<b>Testing Results</b>	<b>7</b>
3.1	Correctness Test . . . . .	7
3.1.1	Test data 1 . . . . .	7
3.1.2	Test data 2 . . . . .	7
3.1.3	Test data 3 . . . . .	8
3.1.4	Test data 4 . . . . .	9
3.1.5	Test data 5 . . . . .	9
<b>4</b>	<b>Analysis and Comments</b>	<b>11</b>
4.1	Min-Heap Analysis . . . . .	11
4.2	Build Tree . . . . .	12
4.3	Huffman Tree Analysis . . . . .	12
	<b>Appendices</b>	<b>13</b>
<b>A</b>	<b>Source Code (in C/C++)</b>	<b>14</b>
<b>B</b>	<b>Declaration and Signatures</b>	<b>23</b>

# Chapter 1

## Introduction

### 1.1 Problem Description

In this project, We are supposed to write a program to judge if the answers the students hand in is right when solving problem on Huffman codes. Because the Huffman codes sometimes can be not unique, we have to think out another way to judge the correctness of the answers rather than just compare its output with the correct answer produced by our program.

The input is a sequence of character and frequency pairs, which is the problem needed to solve by building Huffman codes. And then follow M answers that the students' code produce, for each character, it gives the character and its code. The output we are supposed to give is whether the answer is right, if right, output "Yes" and output "No" otherwise.

### 1.2 Background of Project

#### 1.2.1 Min Heap

Min heap is a specialized tree-based data structure that the parent node is less than its child. Although there is no particular relationship among the siblings, it is a useful data structure when we need to find and remove the highest priority node. And in this project we use min heap to get the lowest frequency characters.

#### 1.2.2 Greedy Algorithm

Greedy algorithm is a useful algorithm which works well in optimizing problems, it choose the choice that can get most benefits for the optimizing target in each stage. Sometimes it can not get the global optimal solution, because there is a great possibility the local optimal solution.

But in this project we solve the problem on Huffman codes, and it is guaranteed to be a global optimal solution.

# Chapter 2

## Algorithm Specification

### 2.1 Definition

#### 2.1.1 function

```
1 void percup(Ptr *list, int i);
2 void ptrcdown(Ptr *list, int i);
3 Ptr deletemin(Ptr *list);
4 void insert(Ptr *list, Ptr node);
5 int compute(Ptr tree);
6 Ptr buildprotree(int N, cf* pro);
7 Ptr buildanstree(int N, cc *ans);
8 void deletetree(Ptr tree);
```

#### 2.1.2 Node

```
1 typedef struct node *Ptr;
2 struct node
3 {
4     bool leaf;
5     int weight;
6     int depth;
7     Ptr left;
8     Ptr right;
9 };
```

There are five variables, the leaf that mark if the node is a leaf node, the weight record the weight of the sub-tree with this node to be the root, and depth is the depth of the node, and the two pointers left and right is used to record the two child node of the node.

#### 2.1.3 CF

```
1 struct cf
2 {
```

```

3   char character;
4   int freq;
5 };

```

CF is used to record the character and frequency pairs.

### 2.1.4 CC

```

1 struct cc
2 {
3     char anscharacter;
4     int freq;
5     string code;
6 };

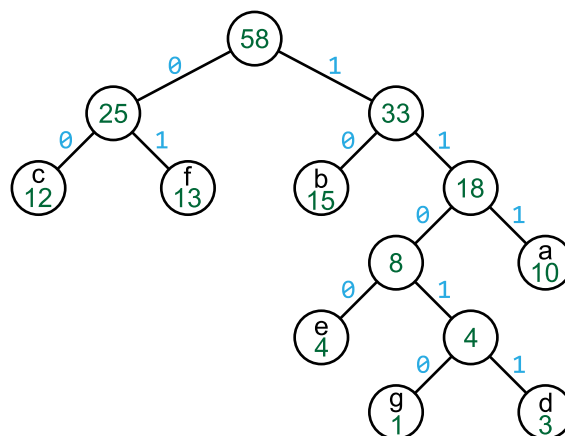
```

To record the character, frequency and code submitted by the students

## 2.2 Huffman coding

### 2.2.1 Huffman tree

Huffman tree is a binary prefix tree, and all strings in Huffman tree are leaf nodes. The Huffman code of a character is corresponding string on the Huffman tree of the data.



character	a	b	c	d	e	f	g
frequency	10	15	12	3	4	13	1

Figure 2.1: Huffman tree

Frequencies of every character in the data are different. The length of compressed data is  $\sum_{i=1}^n L_i F_i$  ( $L_i$  is the length of corresponding string of  $i$ -th character,

$F_i$  is the frequency of occurrence of  $i$ -th character in the data.) Our target is to minimize  $\sum_{i=1}^n L_i F_i$  by arranging the structure of the tree.

### 2.2.2 Greedy algorithm

As we can see, the length of a character with low frequency is better to be long. By this idea, we use a greedy strategy to build the Huffman tree. Assume all characters are a node of tree. Every time we select two nodes of lowest frequency until the size of queue is 1, then merge them (the frequency of new node is sum of two) and push into queue. (the queue here is a priority queue implemented with binary heap)

---

#### Algorithm 1 Build Huffman tree

---

```

1: function HUFFMAN-TREE( $C$ )
2:    $Q := C$ 
3:   for  $i = 1$  to  $|C| - 1$  do
4:     allocate a new node  $z$ 
5:      $x := \text{EXTRACT-MIN}(Q)$ 
6:      $y := \text{EXTRACT-MIN}(Q)$ 
7:      $z.\text{left} := x$ 
8:      $z.\text{right} := y$ 
9:      $z.\text{freq} := x.\text{freq} + y.\text{freq}$ 
10:     $\text{INSERT}(Q, z)$ 
11:   end for
12:   return  $\text{EXTRACT-MIN}(Q)$ 
13: end function

```

---

### 2.2.3 Correctness of the algorithm

#### The greedy-choice property

*Lemma:* Let  $C$  be an alphabet in which each character  $c \in C$  has frequency  $c.\text{freq}$ . Let  $x$  and  $y$  be two characters in  $C$  having the lowest frequencies. Then there exists an optimal prefix code for  $C$  in which the codewords for  $x$  and  $y$  have the same length and differ only in the last bit.

#### The optimal substructure property

*Lemma:* Let  $C$  be a given alphabet with frequency  $c.\text{freq}$  defined for each character  $c \in C$ . Let  $x$  and  $y$  be two characters in  $C$  with minimum frequency. Let  $C'$  be the alphabet  $C$  with a new character  $z$  replacing  $x$  and  $y$ , and  $z.\text{freq} = x.\text{freq} + y.\text{freq}$ . Let  $T'$  be any tree representing an optimal prefix code for the alphabet  $C'$ . Then the tree  $T$ , obtained from  $T'$  by replacing the leaf node for  $z$  with an internal node having  $x$  and  $y$  as children, represents an optimal prefix code for the alphabet  $C$ .

These two properties are proved in the book "Introduction to Algorithms"<sup>1</sup>. The first property implies that choosing first and second minimum nodes is optimum. The second property implies that the weight of a node form from two nodes is sum of two nodes. Thus, the properties provide the correctness of the algorithm.

---

<sup>1</sup>Introduction to Algorithms, 3rd Edition: Ch.16.3; Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest and Clifford Stein. The MIT Press. 2009

# Chapter 3

## Testing Results

### 3.1 Correctness Test

#### 3.1.1 Test data 1

##### Purpose

Different characters with same frequencies.

##### Input

```
1 4
2 A 4 B 2 C 1 D 1
3 2
4 A 0
5 B 10
6 C 110
7 D 111
8 A 0
9 B 10
10 C 111
11 D 110
```

##### Output

```
1 Yes
2 Yes
```

#### 3.1.2 Test data 2

##### Purpose

Three kinds of characters



**Input**

```
1 7
2 9 1 B 1 C 1 D 3 e 3 f 6 _ 6
3 4
4 9 00000
5 B 00001
6 C 0001
7 D 001
8 e 01
9 f 10
10 _ 11
11 9 01010
12 B 01011
13 C 0100
14 D 011
15 e 10
16 f 11
17 _ 00
18 9 000
19 B 001
20 C 010
21 D 011
22 e 100
23 f 101
24 _ 110
25 9 00000
26 B 00001
27 C 0001
28 D 001
29 e 00
30 f 10
31 _ 11
```

**Output**

```
1 Yes
2 Yes
3 No
4 No
```

**3.1.3 Test data 3****Purpose**

Not optimum

**Input**

```
1 4
2 a 4 b 2 c 1 d 1
3 2
4 a 0
5 b 10
6 c 110
7 d 111
8 a 0
9 b 10
10 c 1100
11 d 1101
```

**Output**

```
1 Yes
2 No
```

**3.1.4 Test data 4****Purpose**

One code is a prefix of another code

**Input**

```
1 4
2 a 4 b 2 c 1 d 1
3 1
4 a 0
5 b 10
6 c 1110
7 d 1101
```

**Output**

```
1 No
```

**3.1.5 Test data 5****Purpose**

The input is exactly the best peak shape.

**Input**

```
1 2
2 a 2 b 1
3 1
4 a 0
5 b 1
```

**Output**

```
1 Yes
```

# Chapter 4

## Analysis and Comments

### 4.1 Min-Heap Analysis

The time complexity of heap sorting is mainly in the process of initializing the heap process and rebuilding the heap after each maximum number is selected; Initialize the heap : time is  $O(n)$  Process:

Assume that the height is  $k$ , then from the node on the right of the penultimate layer, the nodes of this layer must perform the comparison of the child nodes and then exchange (if the order is correct, no exchange is needed); if the countdown layer is the third layer, it will select Sub-nodes are compared and exchanged, and if they are not exchanged, they can no longer be implemented. If it is exchanged, then a sub-tree should be selected for comparison and exchange. Then the total time is calculated as:  $s = 2^{(i-1)*(k-i)}$ ; where  $i$  represents the layer, and  $2^{i-1}$  represents how many elements are on the layer,  $(k - i)$  Indicates the number of times to be compared in the subtree. If it is in the worst condition, it will be exchanged after the number of comparisons; because this is a constant, it can be ignored after it is proposed;

$$S = 2^{k-2} * 1 + 2^{k-3} * 2..... + 2 * (k - 2) + 2^0 * (k - 1)$$

Because the leaf layer is not exchanged,  $i$  starts from  $k-1$  to 1;

Solve this equation: The equation is multiplied by 2 and then subtracted from the original equation, it becomes:

$$S = 2^{k-1} + 2^{k-2} + 2^{k-3}..... + 2 - (k - 1)$$

Except for the last item, it is a series of equal numbers, and the sum formula is used directly:  $S = \frac{a[1-(q^n)]}{1-q}$ ;  $S = 2^k - k - 1$ ; and because  $k$  is the depth of the complete binary tree,  $(2^k) \leq n < (2^k - 1)$ , in short, it can be assumed that:  $k = \log n$  (The actual calculation should be  $\log(n+1) < k \leq \log(n)$ ); In summary, we get:  $S = n - \log n - 1$ .

so the time complexity is:  $O(n)$  Rebuild heap time after changing heap elements:  $O(n \log n)$  Process: Loop  $n - 1$  times, each time from the root node to find the next cycle, so each time is  $\log n$ , the total time:  $\log n (n-1) = n \log n - \log n$ ; To sum up: the time complexity of heap sorting is:  $O(n \log n)$  Spatial complexity: The space complexity is  $O(N)$  because it need to store the  $N$  nodes.

## 4.2 Build Tree

The worst situation to build a Huffman Tree is  $O(n)$ , the best situation to build tree is  $O(\log n)$

## 4.3 Huffman Tree Analysis

There are many ways to implement the Huffman tree. We Min-Heap to simply achieve this process. The algorithm is as follows:

1. Add  $n$  terminal nodes to the priority queue, then  $n$  nodes have a priority  $P_i$ ,  $1 \leq i \leq n$
2. If the number of nodes in the queue is greater than 1, then:
  - (1) Remove the two smallest  $P_i$  nodes from the queue, ie make two consecutive removals ( $\min(P_i)$ , Min-Heap)
  - (2) Generate a new node, this node is the parent node of the removed node of (1), and the weight value of this node is (1) the weight of the two nodes
  - (3) Adding the node generated by (2) to the priority queue
3. The last point in the Min-Heap is the root node of the tree (root) The time complexity of this algorithm is  $O(n \log n)$ ; because there are  $n$  terminal nodes, the tree has  $2n - 1$  nodes in total, and the priority queue must use  $O(\log n)$  for each cycle. The space complexity is  $O(N)$  because it need to store the  $N$  nodes.

# Appendices

# Appendix A

## Source Code (in C/C++)

```
1 #include <iostream>
2 #include <string>
3 #include <time.h>
4
5 #define KTIMES 1000
6 //Number of cycles, time of recording small N
7 #define CLK_TCK CLOCKS_PER_SEC
8
9 using namespace std;
10
11 typedef struct node *Ptr;
12 struct node//define huffman tree
13 {
14     bool leaf;// whether leaf or not;
15     int weight;// used to sort;
16     int depth;//depth;
17     Ptr left;
18     Ptr right;
19 };
20
21 struct cf//use to establish problem huffman tree
22 {
23     char character;
24     int freq;
25 };
26
27 struct cc//use to establish answer huffman tree
28 {
29     char anscharacter;
30     int freq;
31     string code;
32 };
33 //heap operation:
34 void percup(Ptr *list, int i);//used in insert heap
35 void ptrcdwn(Ptr *list, int i);//used in deletemin
```

```

36 Ptr deletemin(Ptr *list);
37 void insert(Ptr *list, Ptr node);
38
39 int compute(Ptr tree);//compute PWL
40 Ptr buildprotree(int N, cf* pro);
41 //build problem huffman tree
42 Ptr buildanstree(int N, cc *ans);
43 //build answer huffman tree
44 void deletetree(Ptr tree);
45 // delete huffman tree
46
47 clock_t start, stop;//used to record operation time
48 double duration[KTIMES];
49
50 int main()
51 {
52     srand((unsigned)time(NULL));//reset clock time
53     long long ticks = 0;
54     double totaltime = 0;
55
56     int N;
57     int M, count = 1;
58     int *result;
59     cin >> N;//get N
60     cf *pro;
61     Ptr protree;
62     Ptr anstree;
63     int proPWL;
64     cc *ans;
65     int ansPWL;
66     //used to One-time processing of multi-clock data
67     while(N != 0)
68     {
69         N = N <= 63 ? N : 63;
70         pro = new cf[N];
71         for (int i = 0; i < N; i++)//Initialization
72         {
73             cin >> pro[i].character;
74             cin >> pro[i].freq;
75         }
76         //Cyclic time measurement
77         for (int e = 0; e < KTIMES; e++)
78         {
79             start = clock();
80             protree = buildprotree(N, pro);
81             //build problem tree
82             proPWL = compute(protree);
83             //compute problem PWL

```



```

84         stop = clock();
85         //compute build tree and compute PWL time
86         ticks = ticks + stop - start;
87         duration[e] = ((double)(stop-start))/CLK_TCK;
88         totaltime = totaltime + duration[e];
89     }
90     cout << "Build huffman codes and compute PWL: " << \
91     totaltime / KTIMES << endl;
92     cout << "Ticks: " << ticks << endl;
93     ans = new cc[N];
94     cin >> M;
95     result = new int[M];
96     for (int i = 0; i < M; i++)
97     {
98         for (int k = 0; k < N; k++)//Initialization
99         {
100             cin >> ans[k].anscharacter;
101             for (int j = 0; j < N; j++)
102             {
103                 // get the frequency from
104                 // problem Initialization
105                 if (pro[j].character == \
106                     ans[k].anscharacter)
107                 {
108                     ans[k].freq = pro[j].freq;
109                     break;
110                 }
111             }
112             cin >> ans[k].code;
113         }
114         anstree = buildanstree(N, ans);//build tree
115         ansPWL = compute(anstree);//compute PWL
116         if (ansPWL == proPWL) result[i] = 1;
117         // if problem PWL==answer PWL,
118         // the answer is right
119         else result[i] = 0;
120     }
121     cout << "Case " << count << ":" << endl;
122     for (int i = 0; i < M; i++)
123     {
124         if (result[i] == 1) cout << "Yes" << endl;
125         else cout << "No" << endl;
126     }
127     count++;
128     cin >> N;
129 }
130 system("PAUSE");
131 }

```

```

132
133 void percup(Ptr *list, int i)//heap up
134 {
135     Ptr nownode = list[i];
136     int p;
137     //from up to down, if bigger
138     for (p=i; list[p/2]->weight>nownode->weight; p=p/2)
139     {
140         list[p] = list[p / 2];//up the smaller
141     }
142     list[p] = nownode;//set the node
143 }
144
145 void ptrcdown(Ptr *list, int i)//heap down
146 {
147     Ptr nownode = list[i];
148     int p = i;
149     int child = p * 2;
150     //from down to up, if not out of range
151     for (p = i; p * 2 <= list[0]->depth; p = child)
152     {
153         child = p * 2;
154         if (child + 1 <= list[0]->\
155             depth&&list[child + 1]->weight < \
156             list[child]->weight)//get the smaller child
157         {
158             child++;
159         }
160         //if bigger bobble down
161         if (list[child]->weight < nownode->weight)
162         {
163             list[p] = list[child];
164         }
165         else break;
166     }
167     list[p] = nownode;
168 }
169
170 Ptr deletemin(Ptr *list)//delete min
171 {
172     Ptr min = list[1];
173     int N = list[0]->depth;
174     list[0]->depth--;
175     list[1] = list[N];
176     //get the last to the first , bobble down
177     list[N] = NULL;
178     ptrcdown(list, 1);
179     return min;

```

```

180 }
181
182 void insert(Ptr *list, Ptr node)//insrtion
183 {
184     list[0]→depth++;
185     list[list[0]→depth] = node;
186     percup(list, list[0]→depth);//bobble up
187 }
188
189 int compute(Ptr tree)//compute PWL
190 {
191     if (tree == NULL || (tree→left == NULL &&\
192         tree→right == NULL))
193     {
194         return 0;
195     }
196     else
197     {
198         Ptr queue[63], p;
199         //use tail queue front to get the depth*weight
200         int front = 0;
201         int tail = 0;
202         int PWL = 0;
203         p = tree;
204         p→depth = 0;
205         queue[tail++] = p;
206         while (tail != front)
207         {
208             p = queue[front++];
209             if (front == 63) front = 0;
210             if (p→leaf)//if leaf, compute
211             {
212                 PWL = PWL + p→weight*p→depth;
213             }
214             else
215             {
216                 // It ends before the end, indicating
217                 // that the establishment is incorrect.
218                 if (p→right == NULL || p→left == NULL)
219                 {
220                     return -1;
221                 }
222                 else
223                 {
224                     p→right→depth = p→depth + 1;
225                     //set the right/left child depth
226                     p→left→depth = p→right→depth;
227                     queue[tail++] = p→left;

```

```

228         //set into the queue
229         if (tail == 63) tail = 0; //loop to use
230         else tail = tail;
231         queue[tail++] = p->right;
232         //set into the queue
233         if (tail == 63) tail = 0;
234         else tail = tail;
235     }
236 }
237 }
238 return PWL;
239 }
240 }
241
242 Ptr buildprotree(int N, cf* pro) //build heap
243 {
244     Ptr *heap, ptr, ptr1, ptr2;
245     heap = new Ptr[N + 1];
246     for (int i = 0; i < N + 1; i++) //Initialization
247     {
248         ptr = new node;
249         if (i == 0)
250             ptr->depth = N; //heap[0] used to insert
251         else ptr->depth = 0; //Initialization depth
252         ptr->leaf = true;
253         ptr->left = NULL;
254         ptr->right = NULL;
255         if (i == 0) ptr->weight = -1e8;
256         //Setting an Impossible Value
257         else ptr->weight = pro[i - 1].freq;
258         heap[i] = ptr;
259     }
260
261     for (int i = N / 2; i > 0; i--) ptrcdow(heap, i);
262     //Sort heaps from small to large
263
264     while (heap[0]->depth > 1) //build heap trees by heaps
265     {
266         ptr1 = deletemin(heap); //get the first smallest
267         ptr2 = deletemin(heap); //get the second smallest
268         ptr = new node;
269         ptr->weight = ptr1->weight + ptr2->weight;
270         //get the new node weight
271         ptr->leaf = false;
272         ptr->left = ptr1; //link the old one to the new one
273         ptr->right = ptr2;
274         ptr->depth = 0;
275         insert(heap, ptr); //insert heap into heaps

```

```

276     }
277     ptr = heap[1];
278     delete[] heap;
279     return ptr;
280 }
281
282 Ptr buildanstable(int N, cc *ans)
283 {
284     Ptr tree = new node;
285     Ptr ptr;
286     string nowstring;
287     tree->leaf = false;
288     tree->left = NULL;
289     tree->right = NULL;
290     for (int i = 0; i < N; i++)
291     {
292         ptr = tree;
293         nowstring = ans[i].code;
294         for (int j = 0; j < nowstring.length(); j++)
295         {
296             if (nowstring.at(j) == '0') //0 -> is left
297             {
298                 if (ptr->left == NULL)
299                 {
300                     ptr->left = new node; //initialization
301                     ptr = ptr->left;
302                     ptr->left = NULL;
303                     ptr->right = NULL;
304                     //if loop into the last one character
305                     if (j + 1 == nowstring.length())
306                     {
307                         ptr->leaf = true; // set is leaf
308                         ptr->weight = ans[i].freq;
309                         //set weight
310                         break;
311                     }
312                     else
313                     {
314                         ptr->leaf = false;
315                     }
316                 }
317                 // not null and not the string ends
318                 // but is leaf,
319                 // not the Huffman code answer
320                 else if (ptr->left->leaf == true)
321                 {
322                     deletetree(tree);
323                     return NULL;

```

```

324         }
325         else//loop in
326         {
327             ptr = ptr->left;
328         }
329     }
330     else//1-> is right
331     {
332         if (ptr->right == NULL)
333         {
334             ptr->right = new node;//initialization
335             ptr = ptr->right;
336             ptr->left = NULL;
337             ptr->right = NULL;
338             //if loop into the last one character
339             if (j + 1 == nowstring.length())
340             {
341                 ptr->leaf = true;// set is leaf
342                 ptr->weight = ans[i].freq;
343                 //set weight
344                 break;
345             }
346             else
347             {
348                 ptr->leaf = false;
349             }
350         }
351         // not null and not the string ends
352         // but is leaf,
353         // not the huffman code answer
354         else if (ptr->right->leaf == true)
355         {
356             deletetree(tree);
357         }
358         else
359         {
360             ptr = ptr->right;
361         }
362     }
363 }
364 }
365 return tree;
366 }
367
368 void deletetree(Ptr tree)//delete tree
369 {
370     if (tree->left != NULL) deletetree(tree->left);
371     if (tree->right != NULL) deletetree(tree->right);

```

```
372     delete(tree);  
373 }
```

# Appendix B

## Declaration and Signatures

### Declaration

*We hereby declare that all the work done in this project titled "Huffman Codes" is of our independent effort as a group.*

### Signatures

王中伟

林家丰

金雨若