

浙江大学



# Advanced Data Structures and Algorithm Analysis

Laboratory Projects

Huffman Codes

Group 15

聂俊哲 张琦 张溢弛

Date: 2020-04-20

# Content

**Chapter 1: Introduction**

**Chapter 2: Data Structure / Algorithm Specification**

**Chapter 3: Testing Results**

**Chapter 4: Analysis and Comments**

**Appendix: source code**

**References**

**Declaration**

**Signatures**

# Chapter 1: Introduction

## 1.1 Problem Description

### 1.1.1 Introduce

- In this project, we will write a program to judge whether the Huffman codes designed by the students are proper answers or not because the Huffman codes are not unique at most of the time, we have find a different way to check the correctness of the students' hand-in answers instead of list all the probable results of a given situation of several characters with their frequency.

### 1.1.2 Input format

- The input format in this project is including the following things
  - A integer  $n$ , which means the number of characters.
  - $n$  characters and their frequency in pairs.
  - A integer  $m$ , which means the number of hand-in answers from students.
  - $m$  groups of hand-in answers from students, each contains  $n$  pairs of character and its Huffman code.

### 1.1.3 Output format

- We will check whether the hand-in answer is correct or not and output a "Yes" or "No" to show the result of our program.

### 1.1.4 Something to note:

- There are some traps in this topic
  - The Huffman codes are not unique
  - The hand-in results from students might be unordered.

## 1.2 Background data structures and algorithm

### 1.2.1 Huffman tree

- Huffman tree is a binary prefix tree which store the characters in its leaf nodes. The Huffman code of a character is corresponding string on the Huffman tree of the data
- The total length of compressed characters is  $\sum_{i=1}^n L_i f_i$  where  $L_i$  is the length of the Huffman code of the  $i$ -th character and  $f_i$  is the frequency.
- Such a total length is a **smallest result** of the characters with their frequency, and any possible result of the Huffman code has the same total length. The other result of coding is larger than the smallest result.
- A Huffman tree can be implentmented like this

```

1 typedef struct node
2 {
3     char name;
4     int weight, parent;
5     int lchild, rchild;
6     node():name(), weight(0), parent(0), lchild(0),rchild(0){}
7 }HTNode,*HuffmanTree;

```

we have implement the data structure in the test program, and in the main program, we use a more simple method using the min heap to simulate a Huffman tree rather than directly build a Huffman tree.

### 1.2.2 Min Heap

- Min heap is a special tree-based data structure that the parent node is less than its non-empty children, but the siblings doesn't have any special relations. It can help us find the highest priority node.
- In this project Huffman code, we will use heap to get the lowest frequency characters. And we could use STL in C++ to create a min heap for us directly, which is in head file <queue>
- We will use such an algorithm in the PPT of class08 to build the min heap and calculate the total weight of the tree.
  - The time complexity is  $O(N\log N)$ , we will prove the theory in the Part 4.

```

1 void Huffman ( PriorityQueue heap[ ], int C )
2 { consider the C characters as C single node binary trees,
3   and initialize them into a min heap;
4   for ( i = 1; i < C; i++ ) {
5       create a new node;
6       /* be greedy here */
7       delete root from min heap and attach it to left_child of node;
8       delete root from min heap and attach it to right_child of node;
9       weight of node = sum of weights of its children;
10      /* weight of a tree = sum of the frequencies of its leaves */
11      insert node into min heap;
12  }
13 }

```

### 1.2.3 Greedy Algorithm

- Greedy algorithm works in optimizing problems, it choose the choice that can get most benefits. But this algorithm works if and only if the local optimum is equal to the global optimum
- In this project, we could use greedy algorithm because Huffman code problem is guaranteed to be a optimal solution.

## Chapter 2: Data Structure / Algorithm Specification

### 2.1 Data Structure and Function

### 2.1.1 Two important function

```
1 //The function to check whether some of the code become prefix of others
2 bool check_prefix(vector<string> code);
3 //The function to calculate the smallest sum of weight in the Huffman tree
4 int calculate_weight(priority_queue<int,vector<int>,greater<int>>& huffman_tree);
```

### 2.1.2 The heap structure definition

```
1 priority_queue<int,vector<int>,greater<int>> huffman_tree;
```

- Using the STL priority\_queue to simplify the code

## 2.2 Algorithm Specification

### 2.2.1 pseudo-code of the whole program

- The pseudo-code of our program is as following

```
1 Algorithm:Project 4: Huffman Codes
2 Input: The input format in the 1.1.2
3 Output: m results of the checking result "Yes" or "No"
4
5 read in characters and their frequency;
6 build_min_heap();
7 min_length=calculate_weight(min_heap H);
8 read in m;
9 for(i=0;i<m;i++)
10 {
11     total_length=0;
12     for(j=0;j<n;j++)
13     {
14         read in character and code;
15         total_length+=frequency[character]*code.length();
16     }
17     if(total_length==min_length&&check_prefix)
18     {
19         printf("Yes\n");
20     }
21     else
22     {
23         printf("No\n");
24     }
25     return 0;
26 }
```

### 2.2.2 The important algorithm of the two functions

- Function 1: calculate\_weight

```
1 function calculate_weight(H)
2     result:=0
3     while(true)
4     {
5         weight:=H.top()
6         H.pop()
7         if(H.size()==0)
8             break
9         weight+=H.top()
10        H.pop()
11        H.push(weight)
```

```

12     result+=weight
13 }
14 return result

```

- Function 2: check\_prefix

```

1 function check_prefix(code)
2     sort(code.begin(),code.end())
3     for i=0 to code.size() do
4         len:=code[i].size()
5         for j=i+1 to code.size() do
6             sub_prefix=code[j].substr(0,len)
7             if sub_prefix==code[i]
8                 return true
9             else
10                return false

```

## Chapter 3: Testing Results

- We made up some test data using the test program test.cpp and other methods. You can run the test program to get some new test data. And some large-size test data

### Test 1: Sample Data

- It is a bottom line of test

```

1 input:
2 7
3 A 1 B 1 C 1 D 3 E 3 F 6 G 6
4 4
5 A 00000
6 B 00001
7 C 0001
8 D 001
9 E 01
10 F 10
11 G 11
12 A 01010
13 B 01011
14 C 0100
15 D 011
16 E 10
17 F 11
18 G 00
19 A 000
20 B 001
21 C 010
22 D 011
23 E 100
24 F 101
25 G 110
26 A 00000
27 B 00001
28 C 0001
29 D 001
30 E 00
31 F 10
32 G 11
33 Output:

```

```
34 Yes
35 Yes
36 No
37 No
```

### Test 2:

- have characters of digit, '\_' and letters

```
1 input:
2 7
3 9 1 B 1 C 1 D 3 e 3 f 6 _ 6
4 4
5 9 00000
6 B 00001
7 C 0001
8 D 001
9 e 01
10 f 10
11 _ 11
12 9 01010
13 B 01011
14 C 0100
15 D 011
16 e 10
17 f 11
18 _ 00
19 9 000
20 B 001
21 C 010
22 D 011
23 e 100
24 f 101
25 _ 110
26 9 00000
27 B 00001
28 C 0001
29 D 001
30 e 00
31 f 10
32 _ 11
33 Output:
34 Yes
35 Yes
36 No
37 No
```

### Test 3:

- Have some not optimum condition

```
1 Input:
2 4
3 a 4 c 1 b 2 d 1
4 2
5 a 0
6 b 10
7 c 110
8 d 111
9 a 0
10 c 11000
11 d 10001
12 b 11111
13 Output:
14 Yes
```

15 | No

#### Test 4:

- Have some code be prefix

```
1 | Input:
2 | 7
3 | A 1 B 1 C 1 D 3 E 3 F 6 G 6
4 | 1
5 | A 00000
6 | B 00001
7 | C 0001
8 | D 011
9 | E 01 ---here is a prefix
10 | F 10
11 | G 11
12 | Output:
13 | No
```

#### Test 5:

- A smallest size of test data.

```
1 | Input
2 | 2
3 | a 1 b 2
4 | 2
5 | a 0
6 | b 1
7 | a 1
8 | b 10
9 | Output:
10 | Yes
11 | No
```

#### Test 6:

- A largest size of test data

```
1 | Input
2 | 63 characters with 50 hand-in answers
3 | The test data is too big, you can find the data in the testdata.txt
4 | Output
5 | Yes
6 | Yes
7 | No
8 | No
9 | No
10 | Yes
11 | Yes
12 | Yes
13 | No
14 | No
15 | Yes
16 | Yes
17 | No
18 | No
19 | No
20 | Yes
21 | Yes
22 | Yes
23 | No
```



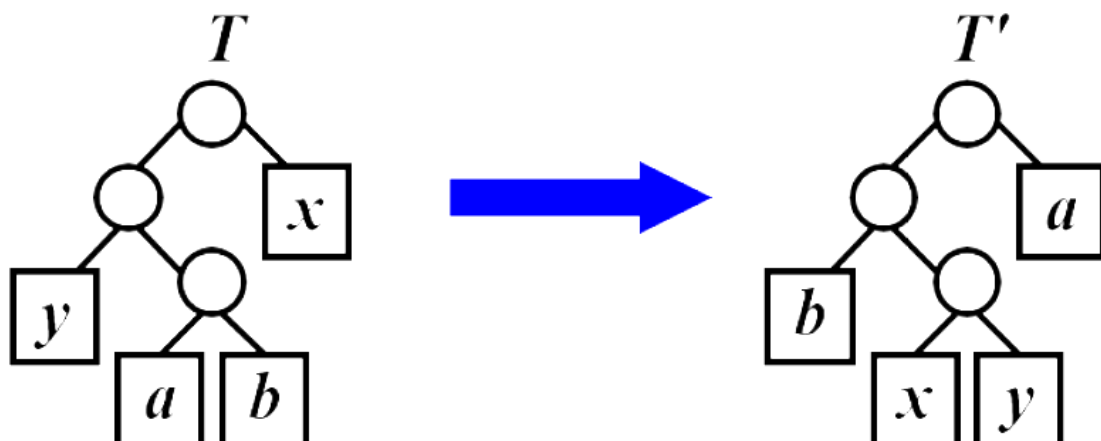
24	No
25	Yes
26	Yes
27	No
28	No
29	No
30	Yes
31	Yes
32	Yes
33	No
34	No
35	Yes
36	Yes
37	No
38	No
39	No
40	Yes
41	Yes
42	Yes
43	No
44	No
45	Yes
46	Yes
47	No
48	No
49	No
50	Yes
51	Yes
52	Yes
53	No
54	No

## Chapter 4: Analysis and Comments

### 4.1 Correctness of the Greedy Algorithm

#### 4.1.1 The Greedy choice property

- **【Lemma】** Let  $C$  be an alphabet in which each character  $c$  in  $C$  has the frequency  $c.\text{freq}$ . Let  $x$  and  $y$  be two characters in  $C$  having the lowest frequencies. Then there must exist 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.



### 4.1.2 The optimal substructure property

- Let  $C$  be a given alphabet with frequency  $c.freq$  defined for each character  $c$  in  $C$ . Let  $x$  and  $y$  be 2 characters in  $C$  with minimum frequency. . Let  $C'$  be the alphabet  $C$  with a new character  $z$  replacing  $x$  and  $y$ , and  $z.freq = x.freq + y.freq$ . 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  $C$ .

## 4.2 Analysis of Min Heap

- The time complexity of the min-heap operation cost most in our program, we have several operations in the program
  - **Build min-heap:** It is  $O(N)$  obviously to build a min heap because each insert cost  $O(\log N)$
  - **Calculate the total weight:** In the while loop, each time we pop 2 elements from the top of the heap(delete min) and insert 1 element into the heap. Delete min from heap and insert into heap both need  $O(\log N)$  each time, so each time we need  $O(3\log N + c)$  time. And we need to do the loop  $N$  times to calculate the total weight of the Huffman tree, so the time complexity is  $O(N\log N)$ .
- The space complexity of the min-heap is obviously  $O(N)$

## 4.3 Analysis of check prefix

- The time complexity is  $O(N^2)$  because we need a double loop to check the prefix
- The space complexity is  $O(1)$  because it needs only constant extra space.

## 4.4 Analysis of algorithm complexity

- Time complexity: The whole program contains several parts:
  - **Read in the data:** It needs  $O(2N + MN) = O(MN)$  time complexity
  - **Heap Operations:** It needs  $O(N\log N)$  according to 4.1
  - **Check each group of answer:** It has  $m$  rounds and each round needs  $O(N^2)$  because of the prefix checking according to 4.2, so it needs  $O(MN^2)$  in total
  - All in all, the time complexity is  $O(MN^2)$
- Space complexity: The whole program contains several parts:
  - Build the heap:  $O(N)$
  - Store the map from character to Huffman code:  $O(N)$
  - Store the result of Yes and No to output them together finally:  $O(M)$
  - So the time complexity is  $O(M + N)$

## Appendix:

## Source Code

```
1  #include<iostream>
2  #include<vector>
3  #include<map>
4  #include<string>
5  #include<queue>
6  #include<algorithm>
7
8  using namespace std;
9
10 //The function to check whether some of the code become prefix of others
11 bool check_prefix(vector<string> code);
12 //The function to calculate the smallest sum of weight in the Huffman tree
13 int calculate_weight(priority_queue<int,vector<int>,greater<int>>&
    huffman_tree);
14
15 int main()
16 {
17     //some basic variables in the program
18     int i,j,n,m,fluquency;
19     char character;
20     //using STL in C++ to build a min-heap for the huffman tree
21     priority_queue<int,vector<int>,greater<int>> huffman_tree;
22     //a map to store the weight of each character
23     map<char,int> huffman_code;
24     //store the result of each group of the input huffman code
25     vector<int> result;
26     //read in the fluquency of each character
27     cin>>n;
28     for(i=0;i<n;i++)
29     {
30         cin>>character>>fluquency;
31         huffman_code[character]=fluquency;
32         //build a min-heap as huffman tree
33         huffman_tree.push(fluquency);
34     }
35     //calculate the total weight of the huffman tree
36     int total_weight=calculate_weight(huffman_tree);
37     // cout<<total_weight<<endl;
38     //check each group of the huffman code
39     cin>>m;
40     for(i=0;i<m;i++)
41     {
42         map<char,string> temp;
43         vector<string> code;
44         char x;
45         string y;
46         int temp_weight=0;
47         //read in the characters and calculate the total weight
48         for(j=0;j<n;j++)
49         {
50             cin>>x>>y;
51             code.push_back(y);
52             temp_weight+=huffman_code[x]*y.length();
53         }
54
55         //a correct group of huffman code should satisfy the 2 conditions
56         //1. the total weight is equal to the smallest weight
57         //2. none of the huffman code is others' prefix
58         if(temp_weight==total_weight&&check_prefix(code))
59         {
60             result.push_back(1);
61         }
62         else{
63             result.push_back(0);
64         }
65     }
```

```

66 //output the result: Yes or No
67 for(i=0;i<result.size();i++)
68 {
69     if(result[i])
70         cout<<"Yes"<<endl;
71     else
72         cout<<"No"<<endl;
73 }
74 return 0;
75 }
76
77 //an algorithm to check if exist prefix
78 bool check_prefix(vector<string> code)
79 {
80     //sort the code in a increasing order
81     sort(code.begin(), code.end());
82     int i,j;
83     //check each huffman code whether is is others' prefix or not
84     //It need O(N^2) time complexity
85     for(i=0;i<code.size();i++)
86     {
87         int len=code[i].size();
88         for(j=i+1;j<code.size();j++)
89         {
90             //get the substring of the first len characters of the huffman code
91             string sub_prefix=code[j].substr(0,len);
92             //check whether they are the same
93             if(sub_prefix==code[i])
94             {
95                 return false;
96             }
97         }
98     }
99     return true;
100 }
101
102 // calculate the total weight of the huffman tree
103 // the result is also a smallest weight
104 int calculate_weight(priority_queue<int,vector<int>,greater<int>>&
huffman_tree)
105 {
106     int result=0;
107     while(true)
108     {
109         //using a greedy method to calculate the total weight of Huffman tree
110         int weight=huffman_tree.top();
111         huffman_tree.pop();
112         if(huffman_tree.size()==0){
113             break;
114         }
115         //simulation the merge of the 2 smallest nodes in the huffman tree
116         weight+=huffman_tree.top();
117         huffman_tree.pop();
118         huffman_tree.push(weight);
119         result+=weight;
120     }
121     return result;
122 }

```

## Test Program

```

1 #include<iostream>
2 #include <cstring>
3 #include<algorithm>
4
5 using namespace std;
6

```

```

7  typedef struct node
8  {
9      char name;
10     int weight, parent;
11     int lchild, rchild;
12     node():name(), weight(0), parent(0), lchild(0), rchild(0){}
13 }HTNode,*HuffmanTree;
14
15 typedef char **HuffmanCode;
16 typedef struct
17 {
18     int weight, locate;
19 }TNode,*Temp;
20
21 void CreatHuffmanTree(HuffmanTree &HT, int n);
22 void CreatHuffmanCode(HuffmanTree HT, HuffmanCode &HC, int n, int m);
23 int main()
24 {
25     HuffmanTree HT;
26     HuffmanCode HC;
27     char p = 'y';
28     int t,n;
29     while (p=='y')
30     {
31         cout << "Please enter the number of characters :";
32         cin >> n;
33         CreatHuffmanTree(HT, n);
34         cout << "Please enter the number of test cases :";
35         cin >> t;
36         for(int m = 0; m < t; m++)
37             CreatHuffmanCode(HT, HC, n , m);
38         cout << "input again: y, quit: n" << endl;
39         cin >> p;
40     }
41     system("pause");
42     return 0;
43 }
44
45 bool compare(const TNode &x, const TNode &y) {
46
47     return x.weight < y.weight;
48 }
49 //Generate Huffman Tree
50 void CreatHuffmanTree(HuffmanTree &HT, int n)
51 {
52     int i, j, k, s1, s2, num;
53     HT = new HTNode[2 * n];
54     cout << "Enter the frequency of characters in format(character
55 frequency):"<<endl;
56     for (i = 1; i <= n; i++)
57     {
58         cin >> HT[i].name >> HT[i].weight;
59     }
60     for (i = n + 1; i < 2 * n; i++)
61     {
62         /*select the subscripts of the two nodes whose parent
63         is 0 and the weight is the smallest from 1 ~ i-1, and
64         count the number of nodes whose parents are 0 in 1 ~ i-1*/
65         num = 0;
66         for (k = 1; k < i; k++)
67         {
68             if (HT[k].parent == 0)
69             {
70                 num++;
71             }
72         }
73         Temp T;
74         T = new TNode[num];

```

```

75         for (j = 0, k = 1; k < i; k++)
76         {
77             if (HT[k].parent == 0)
78             {
79                 T[j].weight = HT[k].weight;
80                 T[j].locate = k;
81                 j++;
82             }
83         }
84         //sort
85         sort(T,T+num,compare);
86         s1 = T[0].locate;
87         s2 = T[1].locate;
88         HT[s1].parent = i;
89         HT[s2].parent = i;
90         HT[i].lchild = s1;
91         HT[i].rchild = s2;
92         HT[i].weight = HT[s1].weight + HT[s2].weight;
93         delete T;
94     }
95 }
96 //Create Huffman Code
97 void CreatHuffmanCode(HuffmanTree HT, HuffmanCode &HC, int n ,int m)
98 {
99     int i, c, f,start;
100    char *cd;
101    HC = new char*[n + 1];
102    cd = new char[n];
103    cd[n - 1] = '\0';
104    for (i = 1; i <= n; i++)
105    {
106        start = n - 1;
107        c = i;
108        f = HT[i].parent;
109        while (f != 0)
110        {
111            start--;
112            if (HT[f].lchild == c)
113            {
114                cd[start] = (m % 2)?'0':'1';
115            }
116            else
117            {
118                cd[start] = (m % 3)?'1':'0';
119            }
120            c = f;
121            f = HT[f].parent;
122        }
123        HC[i] = new char[n - start];
124        strcpy(HC[i], &cd[start]);
125    }
126    delete cd;
127    cout << "Haffman code:";
128    for (i = 1; i <= n; i++)
129    {
130        cout << HT[i].name<<" "<< HC[i] << endl;
131    }
132    cout << endl;
133 }

```

## References

List all the references here in the following format:

[1] 《算法导论》第三版

[2] PTA website <https://pintia.cn/>

[3] Course Slides of ADS

## Author List

- Zhang Yichi 3180103772
- Zhang Qi 3180103162
- Nie Junzhe 3180103501

## Declaration

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

## Signatures

夏俊哲

张琦

张溢邦