



Huffman Codes

Group 15
2020/05/13

THE MAIN CONTENTS

01

Introduction

02

Data Structure and Algorithm Specification

03

Analysis and Comments

04

Testing Results

第

1

部分

Introduction

Problem Description

- 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.
- 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.

Input format

Output format

第

2

部分

Data Structure and Algorithm Specification

Huffman tree

```
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;
```

- 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
- 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

Greedy Algorithm

```
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 }
```

- 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.

02 Data Structure and Algorithm

pseudo-code of the whole program

```
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(code))
18     {
19         printf("Yes\n");
20     }
21     else
22     {
23         printf("No\n");
24     }
25     return 0;
26 }
```


The most important two functions 1

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
```

The most important two functions 2

```
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 false
9     return true
```

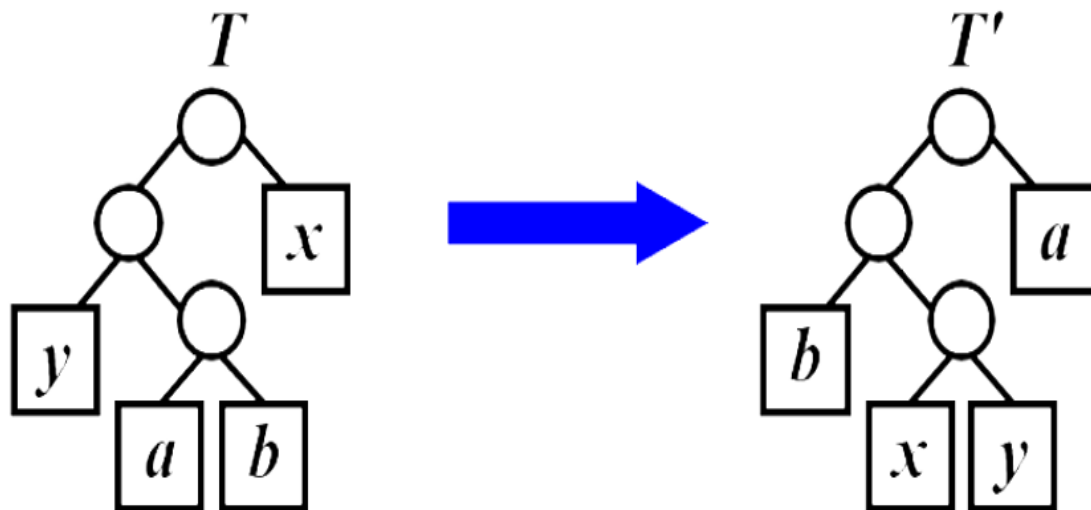
Function 2:
check prefix

第

3

部分

Analysis and Comments



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.

The optimal substructure property

- Let C be a given alphabet with frequency $c.\text{freq}$ 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 $z.\text{freq} = x.\text{freq} + y.\text{freq}$. 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 .

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 $O(N)$

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 need $O(N^2)$ because if 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 : $O(M+N)$

第

4

部分

Testing Results

04 Testing Results

```
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
```

Test 1: Sample Data

```
33 Output:
34 Yes
35 Yes
36 No
37 No
```


04 Testing Results

```
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
```

Test 2: Have characters of digit, '_' and letters

```
33 Output:
34 Yes
35 Yes
36 No
37 No
```

04 Testing Results

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
```

04 Testing Results

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
```

.....



感谢您的观看与聆听



Group 15
2020/05/13