# Project

## Huffman Codes

## Group 17

## Date:2022-05-08

# Chapter 1: Introduction

In 1953, David A. Huffman published his paper "A Method for the Construction of Minimum-Redundancy Codes", and hence printed his name in the history of computer science. Huffman coding is a coding method, and Huffman coding is a kind of variable word length coding (VLC). This method constructs codewords with the shortest average length of different prefixes based on the probability of occurrence of characters. It is sometimes called the best coding. However, there is a big problem that the Huffman codes are not unique. For example, given a string "aaaxuaxz", we can observe that the frequencies of the characters 'a', 'x', 'u' and 'z' are 4, 2, 1 and 1, respectively. We may either encode the symbols as {'a'=0, 'x'=10, 'u'=110, 'z'=111}, or in another way as {'a'=1, 'x'=01, 'u'=001, 'z'=000}, both compress the string into 14 bits. Another set of code can be given as {'a'=0, 'x'=11, 'u'=100, 'z'=101}, but {'a'=0, 'x'=01, 'u'=011, 'z'=001} is NOT correct since "aaaxuaxz" and "aazuaxax" can both be decoded from the code 00001011001001. In this project, we will make a computer program to help determine which ones are correct and which ones are not.

This project is about verifying whether one coding method can be equal to Huffman code in efficiency, or it is already another form of Huffman code.

# Chapter 2: Algorithm Specification

### Main Program

The main program uses function preinput and readsol to read in data, huffman to process the data and, verify to get the answer(Yes/No).

```
1  signed main(signed, char**, char**) {
2      preinput();
3      huffman();
4      For i = 0 To m - 1 Step 1
5          readsol();
6          Print verify();
```

```
7  }
```

**preinput**

This part is to read in the number of characters, characters, their frequencies and, finally, the number of student submissions.

```
1  void preinput() {
2      Input n;
3      For i = 0 To n - 1 Step 1
4          Input idor[i];
5          id[idor[i]] = i;    // Map characters into integer [0,n)
6          Input freq[i];        // Frequency
7      Input m;
8  }
```

**huffman**

This part is a simplified huffman-code process, only to calculate the coded length of the original string under Huffman code.

```
1  void huffman() {
2      priority_queue<int, vector<int>, greater<int>> heap;    //Using a min-hea
3      For i = 0 To n - 1 Step 1
4          heap.push(freq[i]);
5      while (heap.size() != 1) {
6          int a = heap.top();
7          heap.pop();
8          int b = heap.top();
9          heap.pop();
10         len = len + a + b;
11         heap.push(a + b);
12     }
13 }
```

**readsol**

This part is to read in student submissions.

```
1  void readsol() {
2      For i = 0 To n - 1 Step 1
3          Input cdi[i];
4          Input code[id[cdi[i]]];
5  }
```

**prefix**

This part is to determine prefix to verify the correctness of a huffman-code.

```
1  bool prefix(int a, int b) {
2      if (code[a].length() > code[b].length()) then
3          return 0;                    // Length exceeded, apparently not a prefix
4      For i = 0 To code[a].length() - 1 Step 1
5          if (code[a][i] != code[b][i]) then
6              return 0;                // If different, then not a prefix
7      return 1;
8  }
```

**verify**

This part is to verify the correctness of a huffman-code using function prefix and to return Yes/No.

If the coded length is unequal to the Huffman-coded length, then it is surely not a Huffman coding; Another prespect, if the code contain one or more prefix, then it's not a prefix code, and surely not a Huffman code.

```
1  string verify() {
2      int slen = 0;
3      For i = 0 To n - 1 Step 1
4          slen = slen + freq[i] * code[i].length();
5      if (slen != len) then
6          return "No";         // The coded length is different
7      sort(code, code + n);
8      For i = 1 To n - 1 Step 1
9          if (prefix(i - 1, i)) then
10             return "No";    // Contain a prefix
11     return "Yes";            // If both not detected, then it's correct.
```

```
12  }
```

## Chapter 3: Testing Results

| input | output | note |
|---|---|---|
| in1.txt | out1.txt | same as sample |
| in2.txt | out2.txt | 0-1 reversed; Two groups of data overlapping. Only lowercase letters |
| in3.txt | out3.txt | Check whether Huffman code has the same length and prefix error |
| in4.txt | out4.txt | Maximum N & M |
| in5.txt | out5.txt | Minimum N & M |

Besides these cases, you can also test the correctness in PTA *Data Structures and Algorithms (English)* 7-9.

## Chapter 4: Analysis and Comments

N: number of distinct characters

M: number of student submissions

**Space Complexity**

- The program use O(N) space to store the data and results.

**Time Complexity**

- preinput and readsol: both loop N times to input, so it's $O(N)$.

- huffman: $O(N)$, for each character shall get a code.

- prefix: loop code[a].length times, so the average time complexity is $O(N)$.

- verify: two loop each loop n times, the second loop contains one prefix, with a quick-sort on $N$ objects. So the time complexity is $O(N^2 + NlogN) = O(N^2)$.

- main program: preinput + huffman + M * (readsol + verify), so the time complexity is $O(MN^2)$.

## Appendix(readme.txt)

```
We suppose read the source file to get more comments.

We have generated some data (and their coresponding output) in data
folder.
You could use these data to check the program.

We recommend using TDM-GCC9.2.0 to compile,
the C++ standard shall be C++11 or above.
```