哈尔滨工业大学计算学部

# 实验报告

课程名称：数据结构与算法

课程类型：专业基础（必修）

实验项目：树形结构及其应用

实验题目：哈夫曼编码与译码方法

实验日期：2022 年 10 月 20 日

班级：2103601

学号：2021112845

姓名：张智雄

| 设计成绩 | 报告成绩 | 指导老师 |
|---|---|---|
|  |  | 李秀坤 |

## 一、实验目的

哈夫曼编码是一种以哈夫曼树（最优二叉树，带权路径长度最小的二叉树）为基础变长编码方法。其基本思想是：将使用次数多的代码转换成长度较短的编码，而使用次数少的采用较长的编码，并且保持编码的唯一可解性。在计算机信息处理中，经常应用于数据压缩。是一种一致性编码法（又称"熵编码法"），用于数据的无损压缩。要求实现一个完整的哈夫曼编码与译码系统。

## 二、实验要求及实验环境

实验要求：

1. 从文件中读入任意一篇英文文本文件，分别统计英文文本文件中各字符（包括标点符号和空格）的使用频率；

2. 根据已统计的字符使用频率构造哈夫曼编码树，并给出每个字符的哈夫曼编码（字符集的哈夫曼编码表）；

3. 将文本文件利用哈夫曼树进行编码，存储成压缩文件（哈夫曼编码文件）；

4. 计算哈夫曼编码文件的压缩率；

5. 将哈夫曼编码文件译码为文本文件，并与原文件进行比较。

6. 利用堆结构，优化的哈夫曼编码算法
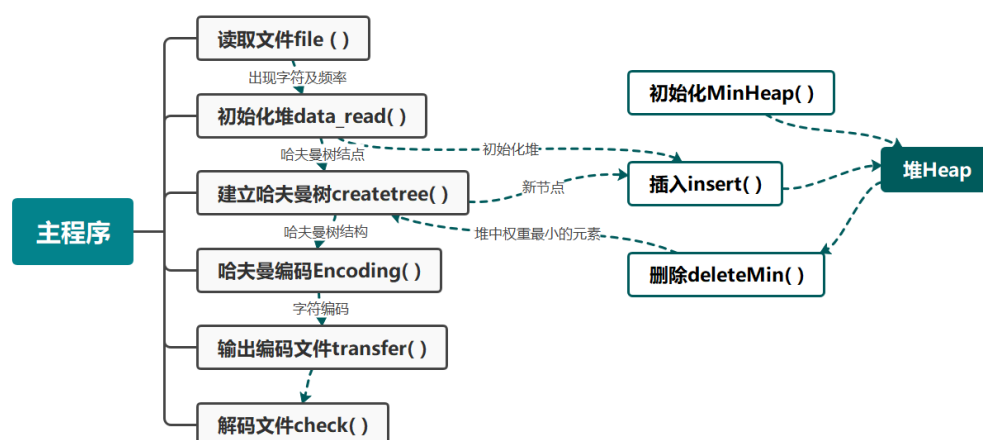
实验环境：Windows11 操作系统+ VS Code 编译器

## 三、设计思想

### 🞿 数据结构：

本实验中的堆结构以及哈夫曼树均采用顺序结构，即逻辑与物理结构相统一的结构，通过**结构体数组**的形式储存节点信息。具体使用数据类型及功能如下表：

| 序号 | 名称 | 功能 | 内含参数 |
|---|---|---|---|
| **1** | **struct data** | 读取数据临时储存 | 文本出现字符**char alphabet**，该字符频率**double frequency** |
| **2** | **struct node(HTNODE)** | 哈夫曼树节点 | 自身编号**int self**，左儿子**int lchild**，右儿子**int rchild**，父节点**int parent**，储存字符**char alphabet**，字符权重**double weight** |
| **3** | **struct code** | 编码结构 | 字符**char ch**，字符对应编码**char bits[MAX+1];** |
| **4** | **struct heap(HEAP)** | 堆结构 | 顺序储存堆结构**HTNODE heap[2*MAX];** 堆内数据个数**int n;** |

## 🔸 函数以及程序主要流程：

函数间调用关系以及程序主要流程图如下图：



各函数具体参数、功能以及返回值如下表：

| 序号 | 函数名 | 函数功能 | 函数参数 | 函数返回值 |
|---|---|---|---|---|
| 1 | **MinHeap( )** | 最小堆结构初始化 | **HEAP \*minheap** | 无 |
| 2 | **insert( )** | 向堆结构中插入新节点 | **HEAP \*minheap, HTNODE x** | 无 |
| 3 | **deleteMin( )** | 删除堆结构最小节点 | **HEAP \*minheap** | 堆最小节点 **HTNODE** |
| 4 | **file ( )** | 读取文件信息 | 字符计数数组**int x[], **信息储存数组**struct data Youth[]** | 等长编码文件长度**pre** |
| 5 | **data_read( )** | 文件信息转移到树节点，并初始化堆 | 信息储存数组 **struct data Youth[], HTNODE T[], HEAP \*minheap** | 无 |
| 6 | **createtree( )** | 建立哈夫曼树 | **HTNODE T[], HEAP \*minheap** | 无 |
| 7 | **Encoding( )** | 哈夫曼编码 | **(HTNODE T[], **编码信息数组**struct code HAff[]** | 无 |
| 8 | **transfer( )** | 输出编码文件，返回压缩后的长度 | 编码信息数组 **struct code HAff[], **输出路径**char path[]** | 压缩后文件长度**now** |
| 9 | **check( )** | 解码文件 | **HTNODE T[], **输出路径**char path[]** | 无 |

## 🔸 核心算法的主要步骤：

## A. 堆结构的插入以及删除

## 插入：

1. 每次插入都是将先将新数据放在当前树的结尾位置，

2. 随后将其与其父节点数据相比较，若其权值小于其父节点，则交换两节点位置，重复此过程，直至小于当前父节点的权值停止。

3. 在插入过程中，最大交换次数为树的高度（log n）

**删除：**

堆中每次都只能删除堆顶元素。删除后需要重建堆，实际的操作是从根节点开始遍历（无需回溯），依次将左右儿子中的较小值赋值给其父节点，最后将未删除前的结尾节点赋给当前遍历到的尾节点。

## B. 哈夫曼树的建立

每次从堆中弹出两个最小的哈夫曼树节点，作为新哈夫曼树节点的左右儿子（默认左儿子权值较小），新哈夫曼树节点的权值等于其左右儿子的权值之和，随后将新节点加入堆中，直至新节点权值为 1 时停止，此时哈夫曼树建立完成。

## C. 哈夫曼编码以及解码

**编码：** 编码过程，自哈夫曼树的叶节点由下而上遍历，通过 start 指针记录当前编码在数组中位置（初始为数组结尾），当为其父节点左儿子时，在编码数组里输入 0；反之为 1。随后 start 前移一位，当遍历至根节点时，将编码数组从 start 所指位置开始复制到编码节点中。

**解码：** 自根节点由上而下开始遍历，读入'0'则移动至其左儿子；'1'为其右儿子，直至遍历到叶节点时输出字符，并将直至重置到根节点，直至文件读取结束。

## D. 压缩率的计算

设文本内使用的字符种类为$n(n \leq 256)$，则使用等长编码长度$e = [log_2 n] + 1$（**取对后向上取整**）；设文本包含字符个数为$X$，则压缩率：

$$\eta = \frac{H(实际哈夫曼译码文件长度)}{eX}$$

# 四、测试结果

测试用例：

Youth is not a time of life; it is a state of mind; it is not a matter of rosy cheeks, red lips and supple knees; it is a matter of the will, a quality of the imagination, a vigor of the emotions; it is the freshness of the deep springs of life. Youth means a temperamental predominance of courage over timidity, of the appetite for adventure over the love of ease. This often exists in a man of 60 more than a boy of 20. Nobody grows old merely by a number of years. We grow old by deserting our ideals.
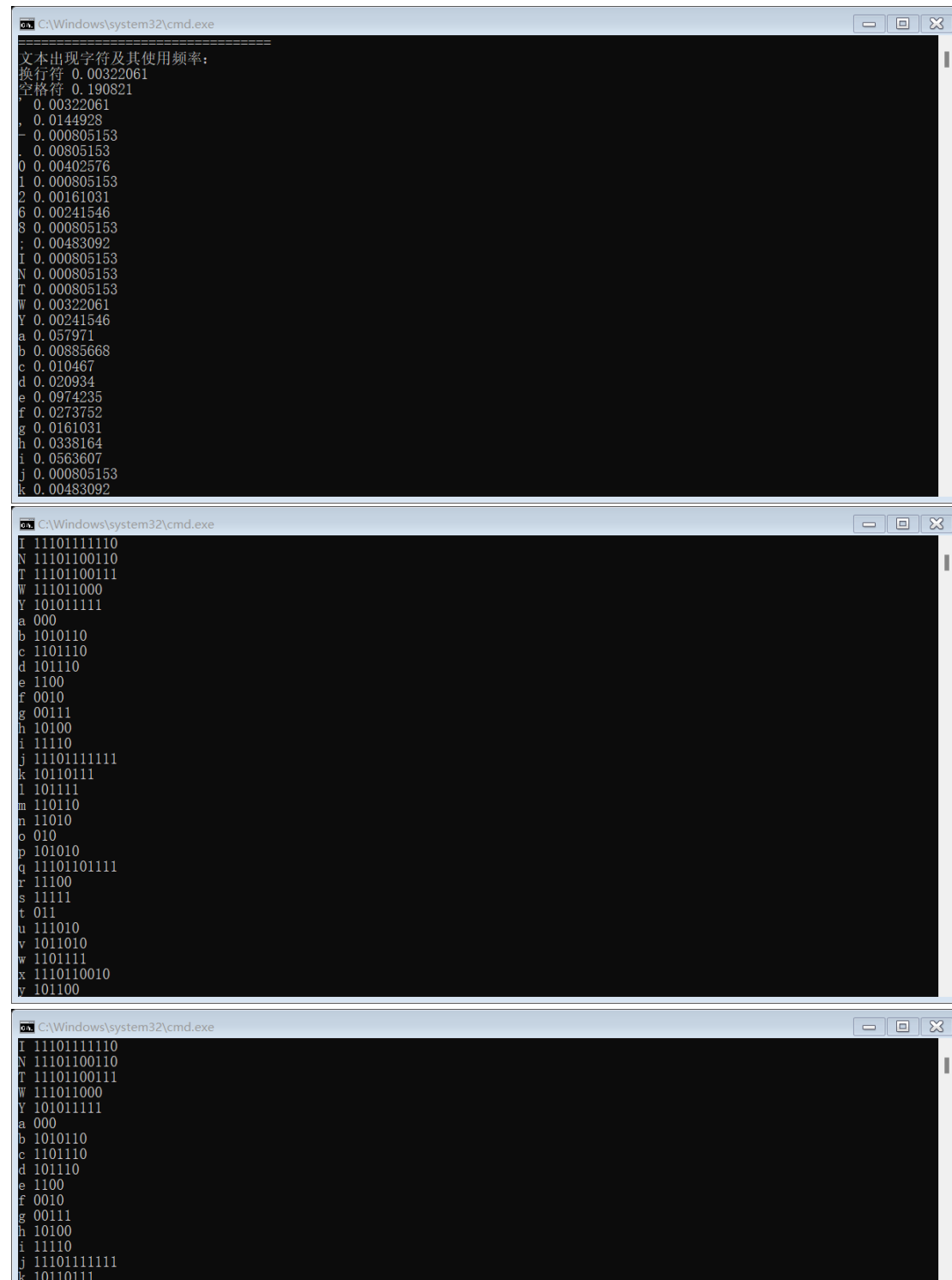
Years may wrinkle the skin, but to give up enthusiasm wrinkles the soul. Worry, fear, self-distrust bows the heart and turns the spirit back to dust.

Whether 60 or 16, there is in every human being's heart the lure of wonders, the unfailing appetite for what's next and the joy of the game of living. In the center of your heart and my heart, there is a wireless station; so long as it receives messages of beauty, hope, courage and power from man and from the infinite, so long as you are young. When your aerials are down, and your spirit is covered with snows of cynicism and the ice of pessimism, then you've grown old, even at 20; but as long as your aerials are up, to catch waves of optimism, there's hope you may die young at 80.

编码结果：（均为同次程序运行结果）

1．运行结果：



```
===============================
文本出现字符及其使用频率：
换行符 0.00322061
空格符 0.190821
'  0.00322061
,  0.0144928
-  0.000805153
.  0.00805153
0  0.00402576
1  0.000805153
2  0.00161031
6  0.00241546
8  0.000805153
;  0.00483092
I  0.000805153
N  0.000805153
T  0.000805153
W  0.00322061
Y  0.00241546
a  0.057971
b  0.00885668
c  0.010467
d  0.020934
e  0.0974235
f  0.0273752
g  0.0161031
h  0.0338164
i  0.0563607
j  0.000805153
k  0.00483092
```



```
I  11101111110
N  11011100110
T  11011100111
W  111011000
Y  101011111
a  000
b  1010110
c  1101110
d  101110
e  1100
f  0010
g  00111
h  10100
i  11110
j  11101111111
k  10110111
l  101111
m  110110
n  11010
o  010
p  101010
q  11101101111
r  11100
s  11111
t  011
u  111010
v  1011010
w  1101111
x  1110110010
y  101100
```



```
I  11101111110
N  11011100110
T  11011100111
W  111011000
Y  101011111
a  000
b  1010110
c  1101110
d  101110
e  1100
f  0010
g  00111
h  10100
i  11110
j  11101111111
k  10110111
```

```
k 1011011
l 101111
m 110110
n 11010
o 010
p 101010
q 11101101111
r 11100
s 11111
t 011
u 111010
v 1011010
w 1101111
x 1110110010
y 101100
```

2. 整理结果：

| 字符 | 概率 | 哈夫曼编码 |
|---|---|---|
| 换行符 | 0.00322061 | 111011010 |
| 空格符 | 0.190821 | 100 |
| ' | 0.00322061 | 111011110 |
| , | 0.0144928 | 00110 |
| − | 0.000805153 | 11101101100 |
| . | 0.00805153 | 11101110 |
| 0 | 0.00402576 | 10101110 |
| 1 | 0.000805153 | 11101101101 |
| 2 | 0.00161031 | 1110111110 |
| 6 | 0.00241546 | 101011110 |
| 8 | 0.000805153 | 11101101110 |
| ; | 0.00483092 | 10110110 |
| I | 0.000805153 | 11101111110 |
| N | 0.000805153 | 11101100110 |
| T | 0.000805153 | 11101100111 |
| W | 0.00322061 | 111011000 |
| Y | 0.00241546 | 101011111 |
| a | 0.057971 | 000 |
| b | 0.00885668 | 1010110 |
| c | 0.010467 | 1101110 |
| d | 0.020934 | 101110 |

| 字符 | 概率 | 哈夫曼编码 |
|---|---|---|
| e | 0.0974235 | 1100 |
| f | 0.0273752 | 0010 |
| g | 0.0161031 | 00111 |
| h | 0.0338164 | 10100 |
| i | 0.0563607 | 11110 |
| j | 0.000805153 | 11101111111 |
| k | 0.00483092 | 10110111 |
| l | 0.0217391 | 101111 |
| m | 0.0241546 | 110110 |
| n | 0.047504 | 11010 |
| o | 0.0676329 | 010 |
| p | 0.0161031 | 101010 |
| q | 0.000805153 | 11101101111 |
| r | 0.0507246 | 11100 |
| s | 0.0563607 | 11111 |
| t | 0.0668277 | 011 |
| u | 0.0249597 | 111010 |
| v | 0.010467 | 1011010 |
| w | 0.0120773 | 1101111 |
| x | 0.00161031 | 1110110010 |
| y | 0.020934 | 101100 |

3. 编码文件（节选）：

10101111101011101001101001001111011111100110100100111000001000111111011011011001000100010100101011111111000
10110010110110100111001100111101111110000010011111011000011100100010001010011011011110110101011101011101
10100111100111001110101111110011010010011100000100110110000011011110011100100010001010011100010111111011001
00110111010100110011001011011111111001101001100110010110100101111111101010101011111100000110101011101001111
11111010101010101011111100100101101111101011001100111111101101101001110011100111101111110000010011011010
00011011110011100100010001010001110100110010011011111111010111110111100110100000100110110111111101000010
1111111100110110010001000101000111010011001001111011011000000111111101101000001111100101101000110100001
00101101011100011010110010001000101000111010011001001100110110010011111100101101011111101101101001011100
11100111101111110001110100110010000010110011001111110100110101100111111111111000100010100011101001100100101
11011001100101010100111111010101011001110110100011111111100010001010010111111111000101100111011101101101101
010111111010111010011010010011011011000001101011111110000010001111001101101010101011001110000010110110011011
00110001011111001010101011001100101100101101101111011010000110101101110110100010001010011011100101110101
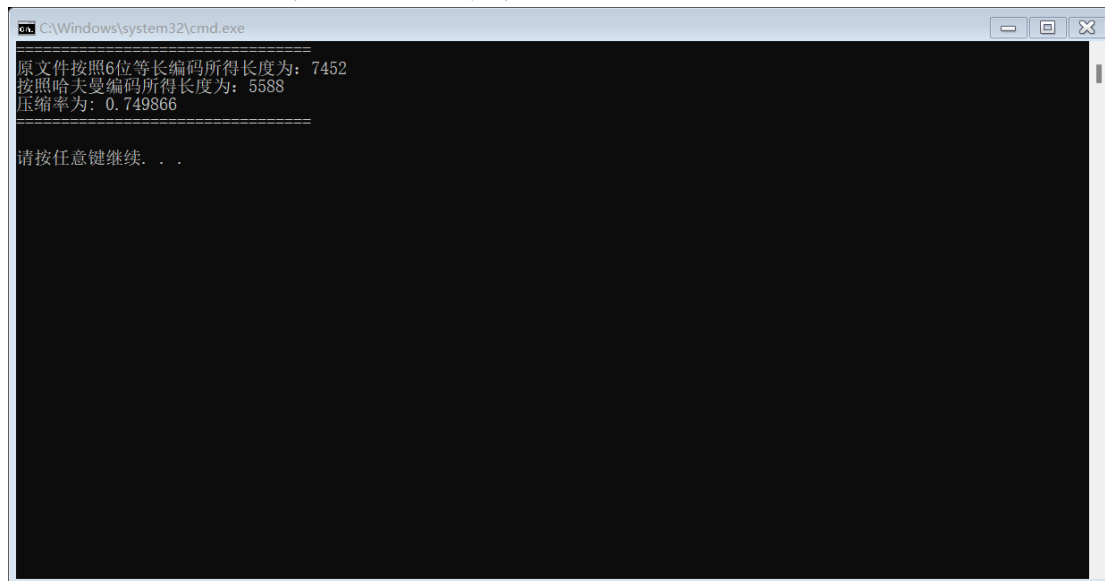
4. 译码结果：

Youth is not a time of life; it is a state of mind; it is not a matter of rosy cheeks, red lips and supple knees; it is a matter of the will, a quality of the imagination, a vigor of the emotions; it is the freshness of the deep springs of life. Youth means a temperamental predominance of courage over timidity, of the appetite for adventure over the love of ease. This often exists in a man of 60 more than a boy of 20. Nobody grows old merely by a number of years. We grow old by deserting our ideals.

Years may wrinkle the skin, but to give up enthusiasm wrinkles the soul. Worry, fear, self-distrust bows the heart and turns the spirit back to dust.

Whether 60 or 16, there is in every human being's heart the lure of wonders, the unfailing appetite for what's next and the joy of the game of living. In the center of your heart and my heart, there is a wireless station; so long as it receives messages of beauty, hope, courage and power from man and from the infinite, so long as you are young.

When your aerials are down, and your spirit is covered with snows of cynicism and the ice of pessimism, then you've grown old, even at 20; but as long as your aerials are up, to catch waves of optimism, there's hope you may die young at 80.

5. 比对结果：正确率 100%，压缩率 74.9866%



## 五、经验体会与不足

### 经验体会：

通过本实验，加深了对哈夫曼树以及堆结构的理解，在实践中体会了哈夫曼编码的前缀性，以及编码的冗余较小的特点，对其中蕴含的贪心思想有了一定了解；而对于堆，使用堆排序可以大大减小时间复杂度，由原来的 O（n）减小到 O（log n），对于本问题的选择优化效果较为明显，但在建立堆以及维护堆的过程中也具有较大的时间复杂度，而且结构本身也增加了空间复杂度。在元素比较多的情况下，还是不错的一个选择。尤其是在解决诸如"前 n 大的数"一类问题时，几乎是首选算法。

### 存在不足：

1. 可以考虑基于单词的压缩，提高压缩效果
2. 可以采用 K 叉的哈夫曼树完成上述工作，实现"K 进制"的编码和译码
3. 在文件读入，信息统计过程中程序还较为冗杂，值得优化

## 六、附录：源代码（带注释）

```cpp
1.  #include <iostream>
2.  #include <fstream>
3.  #include <cstring>
4.  #include <algorithm>
5.  #define SIZE 200
6.  #define MAX 42 //文本所用字符个数
7.
8.  using namespace std;
9.  struct data
10. {
11.     char alphabet;
12.     double frequency = 0;
13. };//读取数据临时储存
14.
15. typedef struct node
16. {
17.     int self;
18.     int lchild = -1;
19.     int rchild = -1;
20.     int parent = -1;
21.     char alphabet;
22.     double weight = 0;
23. }HTNODE;//哈夫曼树节点
24.
25. struct code
26. {
27.     char ch;
28.     char bits[MAX+1];
29. };//编码结构
30.
31. typedef struct heap
32. {
33.     HTNODE heap[2*MAX]; //顺序储存堆结构
34.     int n; //堆内数据个数
35. }HEAP; //堆结构
36.
37. //最小堆结构初始化
38. void MinHeap(HEAP *minheap)
39. {
40.     minheap->n = 0;
41. }
42.
```

```
43.//插入新节点
44.void insert(HEAP *minheap, HTNODE x)
45.{
46.    int i;
47.    if((minheap->n < MAX) && (minheap->n >= 0))
48.    {
49.        i = minheap->n + 1;
50.        while(((i != 1) && (x.weight < minheap->heap[i/2].weight)) || (
   (i != 1) && (x.weight == minheap->heap[i/2].weight) && (x.self < minhea
   p->heap[i/2].self)))
51.        {
52.            minheap->heap[i] = minheap->heap[i/2];
53.            i = i / 2;
54.        }
55.        minheap->heap[i] = x;
56.        minheap->n++;
57.    }
58.}
59.
60.//删除最小节点
61.HTNODE deleteMin(HEAP *minheap)
62.{
63.    int par = 1, child = 2;
64.    HTNODE ele, tmp;
65.    if(minheap->n > 0)
66.    {
67.        ele = minheap->heap[1];
68.        tmp = minheap->heap[minheap->n--];
69.        while(child < minheap->n)
70.        {
71.            if((child < minheap->n) && (minheap->heap[child].weight > m
   inheap->heap[child + 1].weight))
72.                child++;
73.            if(tmp.weight <= minheap->heap[child].weight) break;
74.            minheap->heap[par] = minheap->heap[child];
75.            par = child;
76.            child *= 2;
77.        }
78.        minheap->heap[par] = tmp;
79.        return ele;
80.    }
81.}
82.
83.//读取文件信息，返回等长编码文件长度
```

```
84. int file(int x[], struct data Youth[]);
85. //文件信息转移到树节点，并初始化堆
86. void data_read(struct data Youth[], HTNODE T[], HEAP *minheap);
87. //寻找最小的两个节点（不使用堆）
88. void Select(HTNODE T[], int i, int &p1, int &p2);
89. //建立哈夫曼树
90. void createtree(HTNODE T[], HEAP *minheap);
91. //哈夫曼编码
92. void Encoding(HTNODE T[], struct code HAff[]);
93. //输出编码文件，返回压缩后的长度
94. int transfer(struct code HAff[], char path[]);
95. //解码文件
96. void check(HTNODE T[], char path[]);
97.
98. int main()
99. {
100.         struct data Youth[MAX];
101.         struct code HAFF[MAX];
102.         HTNODE T[2*MAX - 1];
103.         HEAP min;
104.         MinHeap(&min);
105.         int x[SIZE] = {0};
106.         int pre = file(x,Youth);
107.         cout<<"=============================="<<endl;
108.         cout<<"文本出现字符及其使用频率：  "<<endl;
109.         data_read(Youth, T, &min);
110.         cout<<"=============================="<<endl;
111.         createtree(T, &min);
112.         cout<<'\n';
113.         cout<<"=============================="<<endl;
114.         cout<<"各字符哈夫曼编码：   "<<endl;
115.         Encoding(T, HAFF);
116.         cout<<"=============================="<<endl;
117.         int now = transfer(HAFF,"result.txt");
118.         cout<<"原文件按照6位等长编码所得长度为："<<pre<<endl;
119.         cout<<"按照哈夫曼编码所得长度为："<<now<<endl;
120.         cout<<"压缩率为: "<<double(now)/double(pre)<<endl;
121.         cout<<"=============================="<<endl;
122.         check(T, "check.txt");
123.     }
124.
125.     int file(int x[], struct data Youth[])
126.     {
127.         char data;
```

```cpp
128.        ifstream infile;
129.        infile.open("Youth.txt");
130.        int sum = 0;
131.        while(infile.get(data))
132.        {
133.            x[data]++;
134.            sum++;
135.        }
136.        int i = 0;
137.        for(int j = 0; j < SIZE; j++)
138.        {
139.            if(x[j] != 0)
140.            {
141.                Youth[i].alphabet = char(j);
142.                Youth[i].frequency = double(x[j]) / double(sum);
143.                i++;
144.            }
145.        }
146.        infile.close();
147.        return sum*6;
148.    }
149.
150.    void data_read(struct data Youth[], HTNODE T[], HEAP *minheap)
151.    {
152.        //节点初始化
153.        for(int i = 0; i < MAX; i++)
154.        {
155.            if(Youth[i].alphabet == '\n') cout <<"换行符
    "<<" "<< Youth[i].frequency<<endl;
156.            else if(Youth[i].alphabet == ' ')  cout <<"空格符
    "<<" "<< Youth[i].frequency<<endl;
157.            else  cout<<Youth[i].alphabet <<" "<< Youth[i].frequency<
    <endl;
158.            T[i].alphabet = Youth[i].alphabet;
159.            T[i].weight = Youth[i].frequency;
160.            T[i].self = i;
161.            insert(minheap, T[i]);
162.        }
163.    }
164.
165.    void Select(HTNODE T[], int n, int &p1, int &p2)
166.    {
167.        int i, j;
168.        for(i = 0; i < n; i++)
```

```cpp
169.            if(T[i].parent == -1) {p1 = i; break;}
170.        for(j = i + 1; j < n; j++)
171.            if(T[j].parent == -1) {p2 = j; break;}
172.        for(i = 0; i < n; i++)
173.            if((T[p1].weight > T[i].weight) && (T[i].parent == -
    1) && ((i-p2) != 0))  p1 = i;
174.        for(j = 0; j < n; j++)
175.            if((T[p2].weight > T[j].weight) && (T[j].parent == -
    1) && (j != p1))  p2 = j;
176.        //cout<<p1<<p2;
177.    }

179.    void createtree(HTNODE T[], HEAP *minheap)
180.    {
181.        int i, p1, p2,j,k;
182.        for(i = MAX; i < (2*MAX-1); i++)
183.        {
184.            j = deleteMin(minheap).self;
185.            k = deleteMin(minheap).self;
186.            p1 = min(j,k);
187.            p2 = max(j,k);
188.            T[p1].parent = T[p2].parent = i;
189.            T[i].lchild = p1;
190.            T[i].rchild = p2;
191.            T[i].weight = T[p1].weight + T[p2].weight;
192.            T[i].self   = i;
193.            insert(minheap, T[i]);

195.        }
196.        cout << "树建立完毕！";
197.    }

199.    void Encoding(HTNODE T[], struct code HAff[])
200.    {
201.        for(int i = 0; i < MAX; i++)
202.        {
203.            HAff[i].ch = T[i].alphabet;
204.            int start = MAX;
205.            char x[MAX + 1];
206.            x[MAX] = '\0';
207.            int p, q = i;
208.            while((p = T[q].parent) != -1)
209.            {
210.                if(T[p].lchild == q) x[--start] = '0';
```

```cpp
                else if(T[p].rchild == q) x[--start] = '1';
                q = p;
            }
            strcpy(HAff[i].bits, &x[start]);
            if(HAff[i].ch == '\n') cout <<"换行符
    "<<" "<<HAff[i].bits<<endl;
            else if(HAff[i].ch == ' ') cout<<"空格符
    "<<" "<<HAff[i].bits<<endl;
            else cout << HAff[i].ch <<" " <<HAff[i].bits<<endl;
        }
    }

    int transfer(struct code HAff[], char path[])
    {
        ofstream fout(path);
        char x;
        int number = 0;
        if(fout)
        {
            ifstream infile;
            infile.open("Youth.txt");
            while(infile.get(x))
                for(int i = 0; i < MAX; i++)
                    if(HAff[i].ch == x)
                    {
                        fout<<HAff[i].bits;
                        number += strlen(HAff[i].bits);
                    }
        }
        fout.close();
        return number;
    }

    void check(HTNODE T[], char path[])
    {
        char x;
        ifstream infile;
        ofstream fout(path);
        if(fout)
        {
            int p,q;
            infile.open("result.txt");
            for(int i = 0; i < 2*MAX-1;i++)
                if(T[i].weight == 1)  {p = i;q = p;}
```

```cpp
253.          while(infile.get(x))
254.          {
255.              if(x == '0') p = T[p].lchild;
256.              else if(x == '1') p = T[p].rchild;
257.
258.              if((T[p].lchild == -1) && (T[p].rchild == -1))
259.              {
260.                  fout<<T[p].alphabet;
261.                  p = q;
262.              }
263.          }
264.      }
265.      fout.close();
266.      infile.close();
267.  }
```