

Huffman 编码 tsD14627

计 72 陈嘉杰

April 12, 2018

Contents

1	题目说明	1
2	实现思路	1
3	程序编译环境	1
4	实现步骤	2
4.1	数据读入	2
4.2	建 Huffman 树	2
4.3	统计结果	3
5	遇到的问题和得到的收获	4

1 题目说明

文件 `input.txt` 中含有一个字符串，然后要把 Huffman 编码后文本的长度写入 `output.txt` 中。

2 实现思路

先统计词频，然后构造出哈夫曼树。由于这里只需要输出最后文本的长度，所以只需要对这个树进行遍历，对每个叶子结点进行统计即可得到答案。

3 程序编译环境

1. 操作系统：macOS
2. 编译器：LLVM/Clang 6.0.0

4 实现步骤

4.1 数据读入

```
1 freopen("input.txt", "r", stdin);
2 freopen("output.txt", "w", stdout);
3 scanf("%s", buffer);
4 len = strlen(buffer);
5 if (len == 1) {
6     // only one char?
7     // zero entropy here
8     // no reasonable answer
9     printf("1\n");
10    return 0;
11 }
```

首先把文本都保存到一个足够大的数组中。然后这里有一个有争议的地方：题目数据中有一个点，input.txt 中仅有一个字母 a。在这种情况下，Huffman 树仅有一个叶结点。此时信息熵为 0。我不是很确定这里应该写什么答案。还有就是，如果输入的文本是 aaaa，此时信息熵仍然为 0。那此时又如何编码呢。

4.2 建 Huffman 树

首先是数据结构：

```
1 struct Node {
2     Node *left = NULL;
3     Node *right = NULL;
4     int weight = 0;
5     int index = 0;
6
7     bool operator()(const Node *a, const Node *b) {
8         return a->weight > b->weight;
9     }
10 };
```

很直观的一个二叉树，weight 保存子树的词频，index 表示叶结点对应的字母，下面的 operator () 是提供给 priority_queue 的比较器。

接下来，统计词频，建树：

```

1  for (int i = 0; i < len; i++) {
2      freq[buffer[i] - 'a']++;
3  }
4  for (int i = 0; i < 26; i++) {
5      if (freq[i]) {
6          Node *node = new Node;
7          node->weight = freq[i];
8          node->index = i;
9          pq.push(node);
10     }
11 }
12 while (pq.size() > 1) {
13     Node *first = pq.top();
14     pq.pop();
15
16     Node *second = pq.top();
17     pq.pop();
18
19     Node *new_node = new Node;
20     new_node->left = first;
21     new_node->right = second;
22     new_node->weight = first->weight + second->weight;
23     pq.push(new_node);
24 }

```

这就是很常规的建立 Huffman 树的方法。

4.3 统计结果

接下来，遍历这棵树，然后统计结果：

```

1  int calc_length(Node *current, int cur_depth) {
2      if (current->left == NULL && current->right == NULL) {
3          return freq[current->index] * cur_depth;
4      }
5      return calc_length(current->left, cur_depth + 1) +
6             calc_length(current->right, cur_depth + 1);
7  }
8
9  int main() {

```

```
10  // ..snip..  
11  root = pq.top();  
12  printf("%d\n", calc_length(root, 0));  
13  return 0;  
14 }
```

对于字母表中的每个字母，它的高度就是 Huffman 编码中的长度，乘以词频求和即是答案。最后输出即可。

5 遇到的问题和得到的收获

遇到的问题是，第一，priority_queue 的默认比较，由于我这里用的是 Node * 类型，所以会变成指针比较。所以需要自己写一个 comparator，但是全局的 operator < 又不支持指针。所以复用了 struct Node 作为比较器，这样就可以了。第二就是，输入的文本只有一个字母，这个我觉得没有一个很好的答案。为了过 OJ，只写了一个小小的判断，留下我的疑问。