

HW5

1. 假定一个数据文件由 8 位字符组成，其中所有 256 个字符出现的频率大致相同：最高的频率也低于最低频率的 2 倍。证明：在此情况下，赫夫曼编码并不比 8 位固定长度编码更有效。

在Huffman编码中，最开始将两个频率最小的节点合并，由于最高的频率低于最低频率的2倍，故合并后的节点大于最高频率的节点，故只要还有字符存在，就会合并字符中频率最下的两个。因此一开始会形成128个子树，每个子树的孩子都是字符。由于 $2f_{max} \leq 4f_{min}$ ，设128个子树的频率为 f' ，则有 $f'_{max} \leq 2f'_{min}$ ，和之前一样，会形成64个子树，每个子树的孩子节点都为上一轮合并的节点。以此类推，最终得到高度为 $\log(256)=8$ 的Huffman树，且所有的256个字符的深度均为8，故Huffman编码并不比 8 位固定长度编码更有效。

2. 令 S 是一个有限集， S_1, S_2, \dots, S_k 是 S 的一个划分，这些集合都是非空且不相交的。定义结构 (S, I) 满足条件 $I = \{A : |A \cap S_i| \leq 1, i = 1, \dots, k\}$ 。证明： (S, I) 是一个拟阵。也就是说，与划分中所有子集都最多有一个共同元素的集合 A 组成的集合构成了拟阵的独立集。

(1) S 是一个有限集

(2) 若集合 $A \subseteq I$ ，且 $B \subseteq A$ ，则有 $|A \cap S_i| \leq 1, i = 1, \dots, k$ ，由于 B 是 A 的子集，故 $|B \cap S_i| \leq 1, i = 1, \dots, k$ ，进而 $B \subseteq I$ ，故是遗传的

(3) 假设不满足交换性质，也就是说， $A \subseteq I, B \subseteq I$ 且 $|A| < |B|$ ，有 $\forall x \in B - A, A \cup \{x\} \notin I$ ，即 $|\{x\} \cup A \cap S_i| > 1$ ，而 $|A \cap S_i| \leq 1$ ，故有 $x \in S_i$ 。即当 (S, I) 不满足交换性质时， $\forall x \in B - A$ ，有 $x \in S_i$ ，故为使 $|B \cap S_i| \leq 1$ ， $|B - A|$ 应小于等于 1。

当 $|B| > |A| + 1$ 时，当 A 是 B 的真子集时， $|B - A|$ 最小，而此时 $|B - A| > 1$ ，此时不满足 $|B - A|$ 小于等于 1（由上述结论）。

当 $|B| = |A| + 1$ 时，为使 $|B - A|$ 小于等于 1， A 应为 B 的真子集，设 $b \in B - A$ ，则 $b \in S_i$ 。因为 $A \subseteq I, B \subseteq I$ ， B 最多和 S_i 有一个重复元素，而 $b \in S_i$ ，故 B 中除去 b 以外的元素组成的集合与 S_i 没有交集，而 $A = B - \{b\}$ ，即 A 与 S_i 没有交集，此时将 b 加入集合 A 中，新的集合 A' 与 S_i 只有一个重复元素，即 $|\{b\} \cup A| \subseteq I$ ，这与假设不满足交换性质矛盾。

由上可得， (S, I) 满足交换性质。

综上， (S, I) 是一个拟阵。

3. $A = a_1, \dots, a_n$ 表示一个正整数集合。A 中的元素之和为 N 。设计一个 $O(n \cdot N)$ 的算法来确定是否存在一个 A 的子集 B ，使得 $\sum_{a_i \in B} a_i = \sum_{a_i \in A-B} a_i$

即找到子集 B ，使其和为 $N/2$

定义 $S[i, k]$ 表示在 $\{a_1, \dots, a_i\}$ 中若能找到子集，使得其和为 k ，则为 true，否则为 false

$$S[i, k] = \begin{cases} false & \text{其他} \\ true & S[i-1, k-a_i] \text{ 为 true} \end{cases}$$

采用自底向上计算，先将 $S[1, a_1]$ 、 $S[1, 0]$ 设为 true，其余的 $S[1, k]$ 设为 false，再依次计算接下来的行数的 S 的值。

若存在一个 i ，使得 $S[i, N/2] = true$ ，则存在一个 A 的子集 B ，使得 $\sum_{a_i \in B} a_i = \sum_{a_i \in A-B} a_i$

此算法对矩阵S中的每个元素遍历了一遍，故算法的时间复杂度为 $O(n \cdot N/2) = O(n \cdot N)$