



数据结构与算法(五)

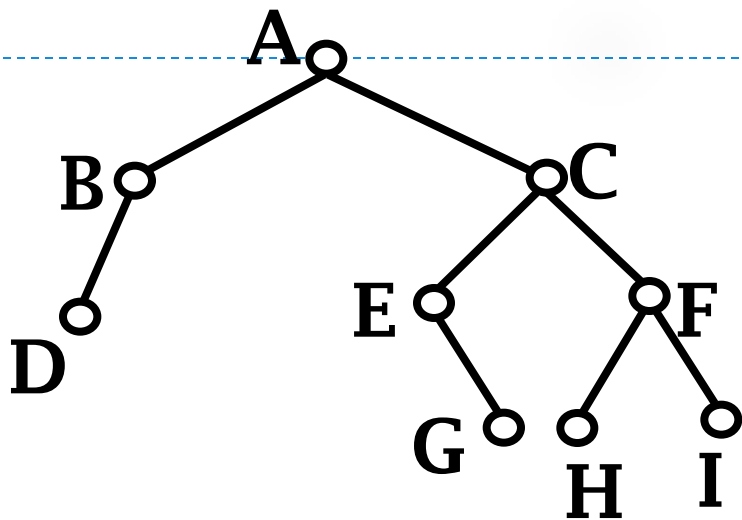
张铭 主讲

采用教材：张铭，王腾蛟，赵海燕 编写
高等教育出版社，2008. 6（“十一五”国家级规划教材）



第五章 二叉树

- 二叉树的概念
- 二叉树的抽象数据类型
 - 深度优先搜索
 - 宽度优先搜索
- 二叉树的存储结构
- 二叉搜索树
- 堆与优先队列
- Huffman树及其应用





等长编码

- 计算机二进制编码
 - ASCII 码
 - 中文编码
- 等长编码
 - 假设所有编码都等长
表示 n 个不同的字符需要 $\log_2 n$ 位
 - 字符的使用频率相等
- 空间效率

数据压缩和不等长编码

- 频率不等的字符

Z	K	F	C	U	D	L	E
2	7	24	32	37	42	42	120

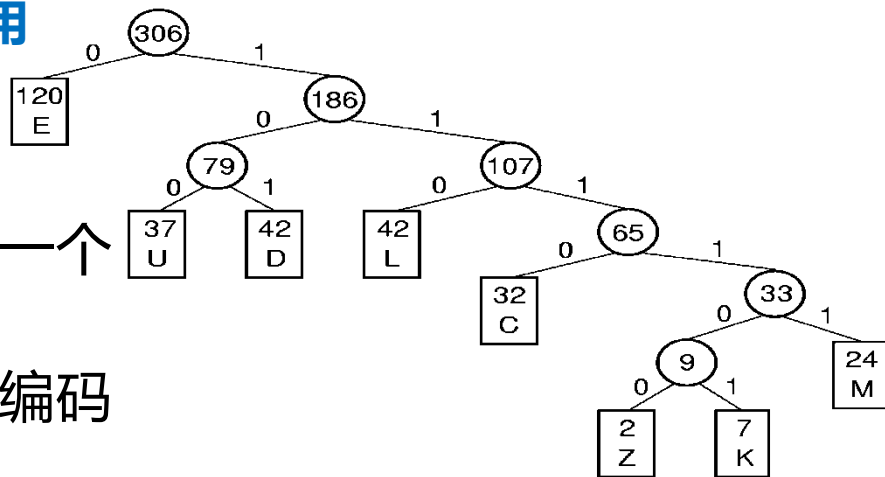
- 可以利用字符的出现频率来编码
 - 经常出现的字符的编码较短，不常出现的字符编码较长
- 数据压缩既能节省磁盘空间，又能提高运算速度
(**外存时空权衡的规则**)



5.6 Huffman树及其应用

前缀编码

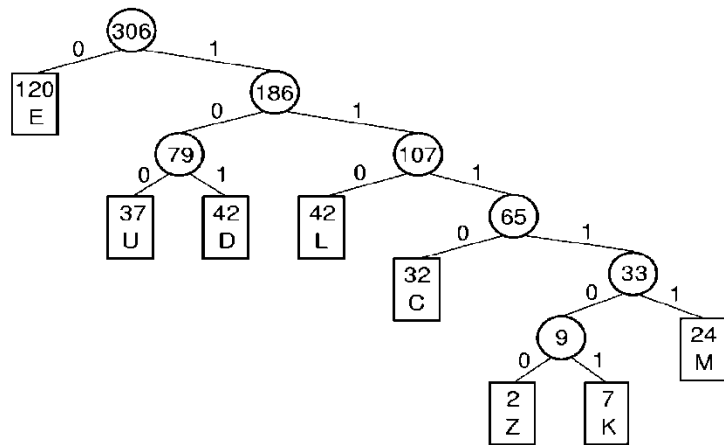
- 任何一个字符的编码都不是另外一个字符编码的前缀
- 这种前缀特性保证了代码串被反编码时，不会有多种可能。例如
- 右图是一种前缀编码，对于 “000110”
，可以翻译出唯一的字符串 “EEEL”。
- 若编码为Z(00), K(01), F(11), C(0), U(1),
D(10), L(110), E(010)。则对应：
“ ZKD” , “ CCCUUC” 等多种可能



- 编码 Z(111100),
K(111101),
F(11111),
C(1110), U(100),
D(101), L(110),
E(0)

Huffman树与前缀编码

- Huffman编码将代码与字符相联系
- 不等长编码
- 代码长度取决于对应字符的相对使用频率或“权”



5.6 Huffman树及其应用

建立Huffman编码树

- 对于n个字符 K_0, K_1, \dots, K_{n-1} ，它们的使用频率分别为 w_0, w_1, \dots, w_{n-1} ，给出它们的前缀编码，使得总编码效率最高
- 给出一个具有n个外部结点的扩充二叉树
 - 该二叉树每个外部结点 K_i 有一个权 w_i 外部路径长度为 l_i
 - 这个扩充二叉树的叶结点带权外部路径长度总和

$$\sum_{i=0}^{n-1} w_i \cdot l_i$$

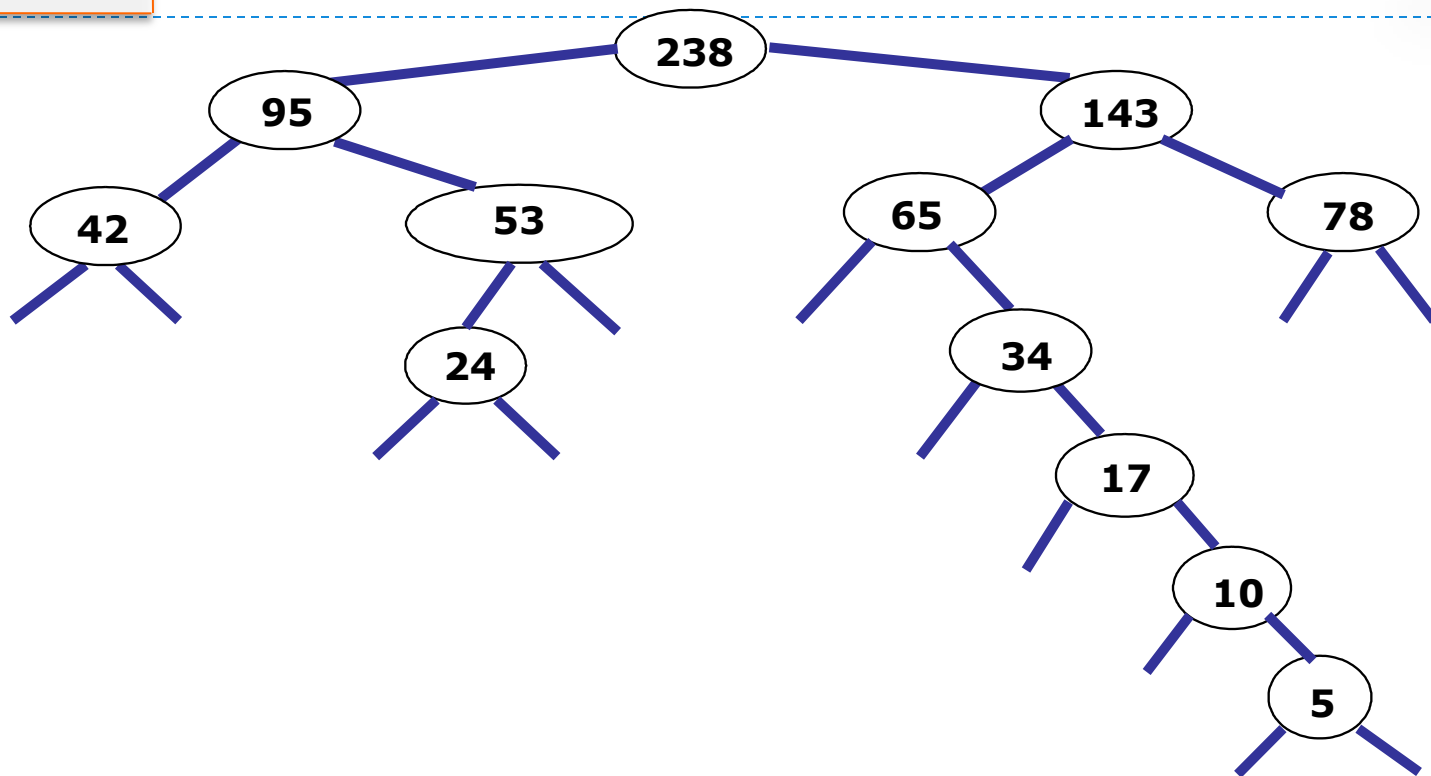
□ 权越大的叶结点离根越近



建立Huffman编码树

- 首先，按照“权”（例如频率）将字符排为一列
 - 接着，拿走**前两个字符**（“权”最小的两个字符）
 - 再将它们标记为Huffman树的树叶，将这两个树叶标为一个分支结点的两个孩子，而该结点的权即为两树叶的权之和
- 将所得“权” **放回序列**，使“权”的顺序保持
- 重复上述步骤**直至序列处理完毕**

5.6 Huffman树及其应用



2

3

5

7

11

13

17

19

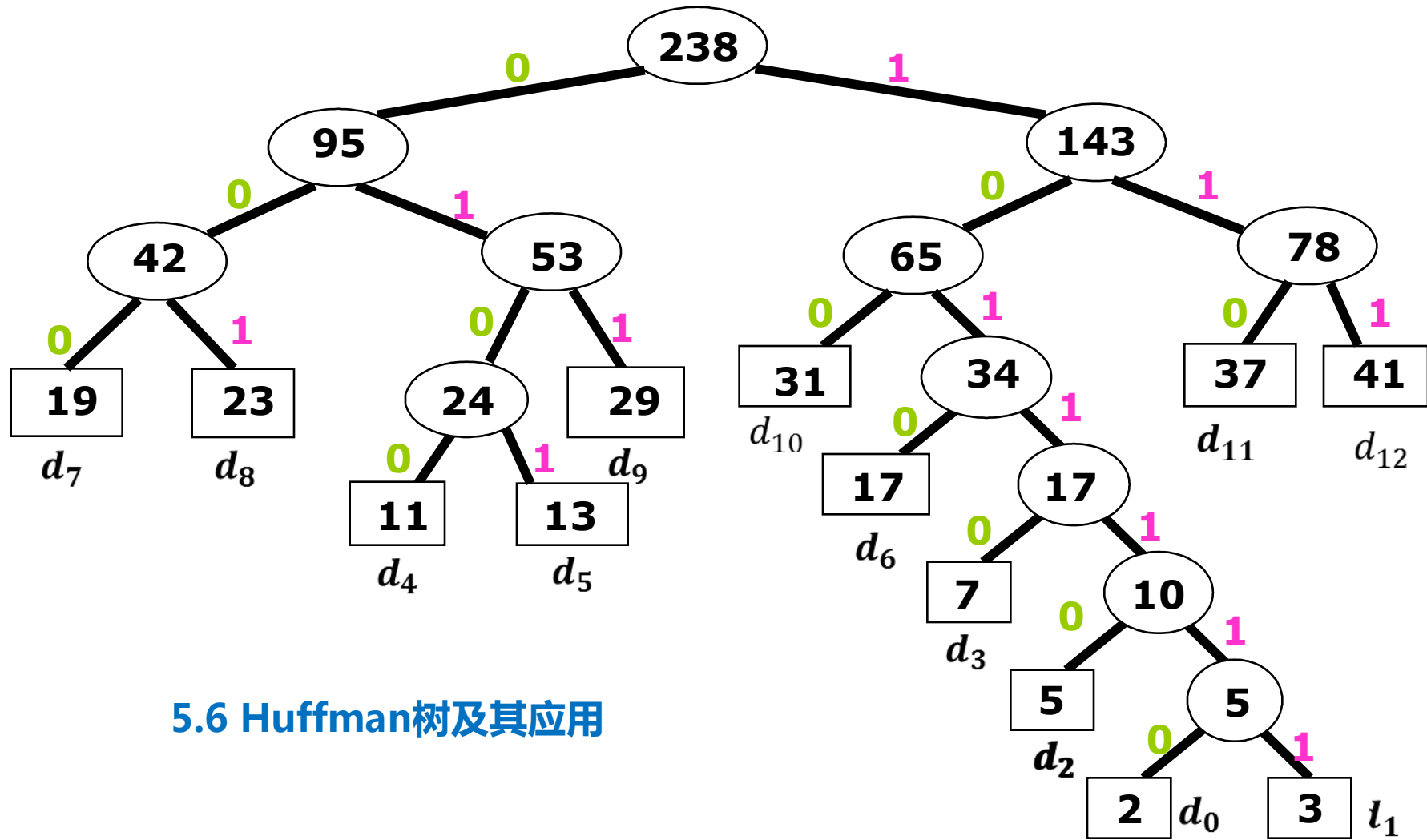
23

29

31

37

41



5.6 Huffman树及其应用

5.6 Huffman树及其应用

频率越大其编码越短

各字符的二进制编码为：

d_0 : 1011110 d_1 : 1011111

d_2 : 101110 d_3 : 10110

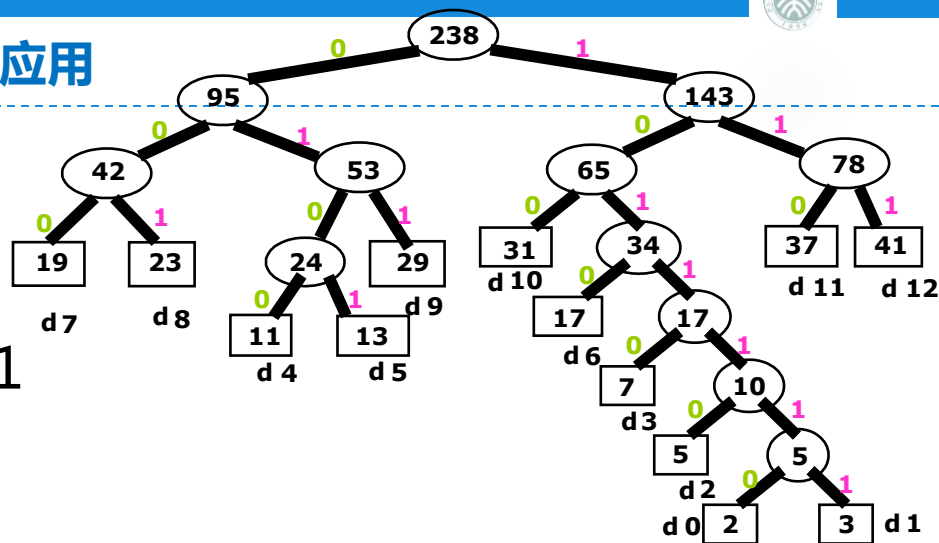
d_4 : 0100 d_5 : 0101

d_6 : 1010 d_7 : 000

d_8 : 001 d_9 : 011

d_{10} : 100 d_{11} : 110

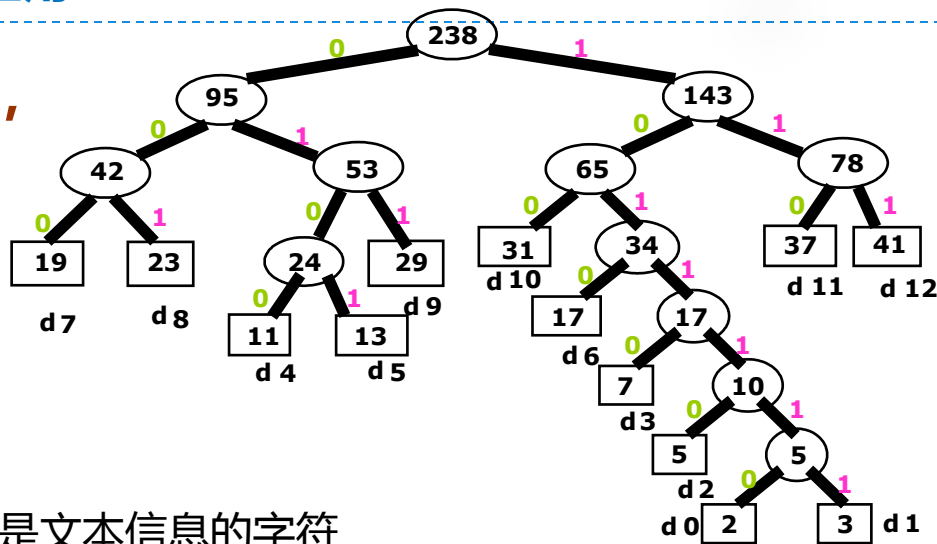
d_{12} : 111

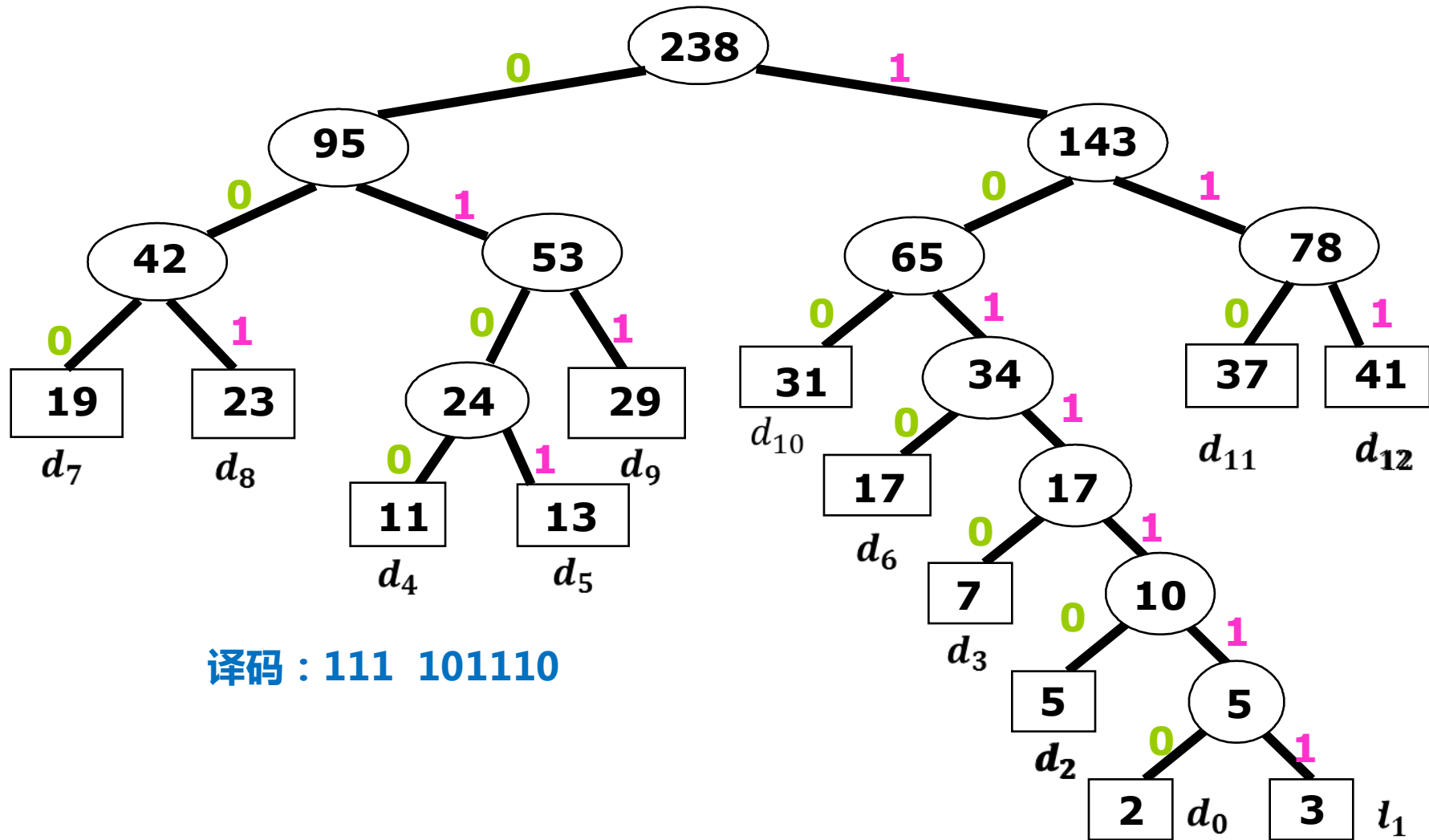


5.6 Huffman树及其应用

译码： 从左至右逐位判别代码串，
直至确定一个字符

- 与编码过程相逆
 - 从树的根结点开始
 - “0” 下降到左分支
 - “1” 下降到右分支
 - 到达一个树叶结点，对应的字符就是文本信息的字符
- 连续译码
 - 译出了一个字符，再回到树根，从二进制位串中的下一位开始继续译码







Huffman树类

```
template <class T> class HuffmanTree {  
private:  
    HuffmanTreeNode<T>* root; // Huffman树的树根  
    // 把ht1和ht2为根的合并成一棵以parent为根的Huffman子树  
    void MergeTree(HuffmanTreeNode<T> &ht1,  
        HuffmanTreeNode<T> &ht2, HuffmanTreeNode<T>* parent) ;  
public:  
    // 构造Huffman树, weight是存储权值的数组, n是数组长度  
    HuffmanTree(T weight[], int n) ;  
    virtual ~HuffmanTree(){DeleteTree(root);} ; // 析构函数  
}
```



5.6 Huffman树及其应用

Huffman树的构造

```
template<class T>
HuffmanTree<T>::HuffmanTree(T weight[], int n) {
    MinHeap<HuffmanTreeNode<T>> heap; //定义最小值堆
    HuffmanTreeNode<T> *parent,&leftchild,&rightchild;
    HuffmanTreeNode<T>* NodeList =
        new HuffmanTreeNode<T>[n];
    for(int i=0; i<n; i++) {
        NodeList[i].element =weight[i];
        NodeList[i].parent = NodeList[i].left
                               = NodeList[i].right = NULL;
        heap.Insert(NodeList[i]); //向堆中添加元素
    } //end for
```

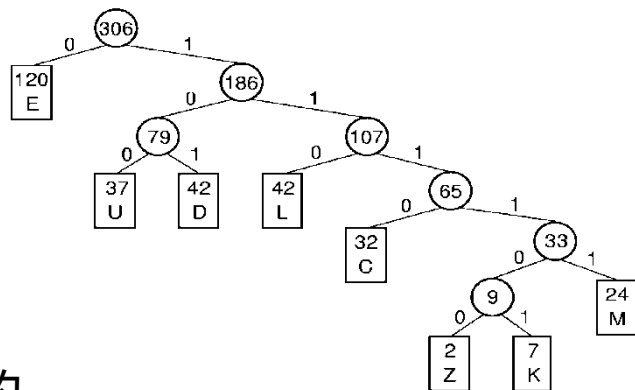


Huffman树的构造

```
for(i=0;i<n-1;i++) {                                //通过n-1次合并建立Huffman树
    parent=new HuffmanTreeNode<T>;
    firstchild=heap.RemoveMin();                     //选值最小的结点
    secondchild=heap.RemoveMin();                   //选值次小的结点
    MergeTree(firstchild,secondchild,parent);        //合并权值最小的两棵树
    heap.Insert(*parent);                            //把parent插入到堆中去
    root=parent;                                    //建立根结点
} //end for
delete []NodeList;
}
```


Huffman方法的正确性证明

- ❑ 是否前缀编码？
- ❑ 贪心法的一个例子
 - ❑ Huffman树建立的每一步，“权”最小的两个子树被结合为一新子树
- ❑ 是否最优解？





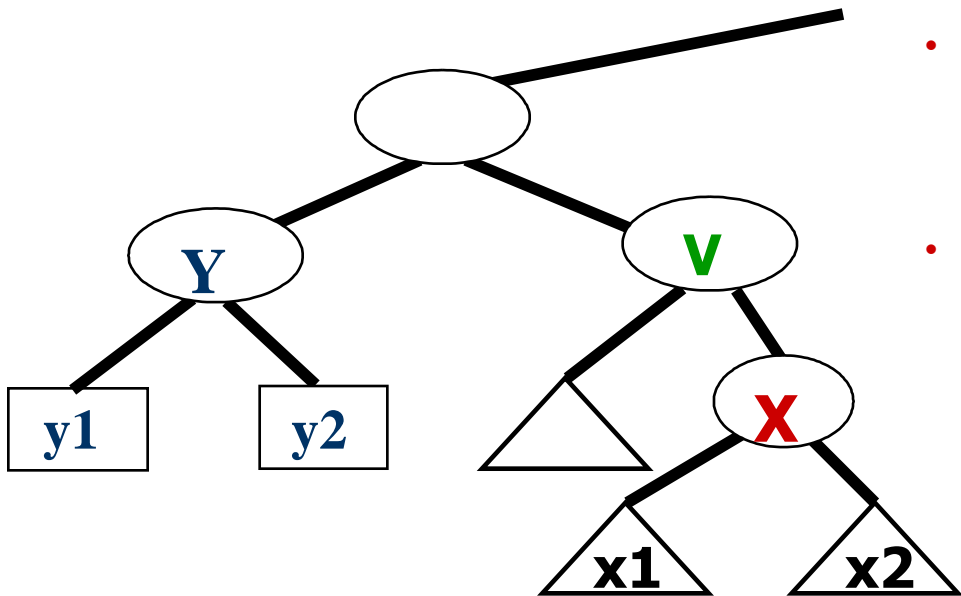
Huffman性质

□ 引理

含有两个以上结点的一棵 Huffman 树中，字符使用频率最小的两个字符是兄弟结点，而且其深度不比树中其他任何叶结点浅

5.6 Huffman树及其应用

证明



- 记使用频率最低的两个字符为 $y1$ 和 $y2$
- 假设 $x1, x2$ 是最深的结点
 - $y1$ 和 $y2$ 的父结点 Y 一定会有比 X 更大的“权”
 - 否则，会选择 Y 而不是 X 作为结点 V 的子结点
- 然而，由于 $y1$ 和 $y2$ 是频率最小的字符，这种情况不可能发生

5.6 Huffman树及其应用

- 定理：对于给定的一组字符，函数HuffmanTree实现了“最小外部路径权重”
- 证明：对字符个数 n 作归纳进行证明
- 初始情况：令 $n = 2$ ，Huffman树一定有最小外部路径权重
 - 只可能有成镜面对称的两种树
 - 两种树的叶结点加权路径长度相等
- 归纳假设：

假设有 $n-1$ 个叶结点的由函数HuffmanTree产生的Huffman树有最小外部路径权重



归纳步骤：

- 设一棵由函数HuffmanTree产生的树 T 有 n 个叶结点， $n > 2$ ，并假设字符的“权” $w_0 \leq w_1 \leq \dots \leq w_{n-1}$
 - 记 V 是频率为 w_0 和 w_1 的两个字符的父结点。根据引理，它们已经是树 T 中最深的结点
 - T 中结点 V 换为一个叶结点 V' （权等于 $w_0 + w_1$ ），得到另一棵树 T'
- 根据归纳假设， T' 具有最小的外部路径长度
- 把 V' 展开为 V （ $w_0 + w_1$ ）， T' 还原为 T ，则 T 也应该有最小的外部路径长度
- 因此，根据归纳原理，定理成立

Huffman树编码效率

- 估计Huffman编码所节省的空间
 - 平均每个字符的代码长度等于每个代码的长度 c_i 乘以其出现的概率 p_i , 即:
 - $c_0p_0 + c_1p_1 + \dots + c_{n-1}p_{n-1}$
或 $(c_0f_0 + c_1f_1 + \dots + c_{n-1}f_{n-1}) / f_T$

这里 f_i 为第 i 个字符的出现频率, 而 f_T 为所有字符出现的总次数

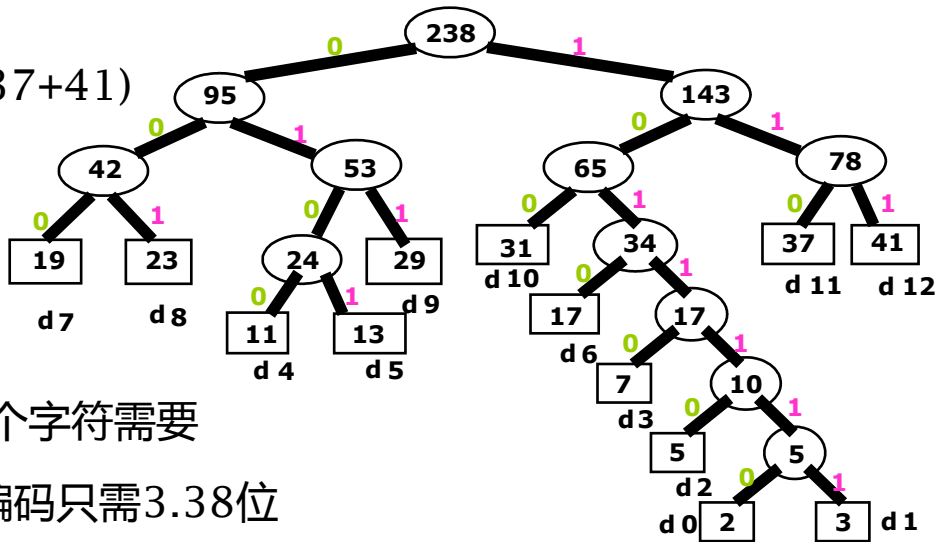
5.6 Huffman树及其应用

Huffman树编码效率 (续)

- 图中，平均代码长度为：

$$\begin{aligned}
 & (3 \times (19 + 23 + 24 + 29 + 31 + 34 + 37 + 41) \\
 & + 4 \times (11 + 13 + 17) \\
 & + 5 \times 7 \\
 & + 6 \times 5 \\
 & + 7 \times (2 + 3)) / 238 \\
 & = 804 / 238 \approx 3.38
 \end{aligned}$$

- 对于这13个字符，等长编码每个字符需要 $\lceil \log 13 \rceil = 4$ 位，而Huffman编码只需3.38位
- Huffman编码预计只需要等长编码 $3.38/4 \approx 84\%$ 的空间





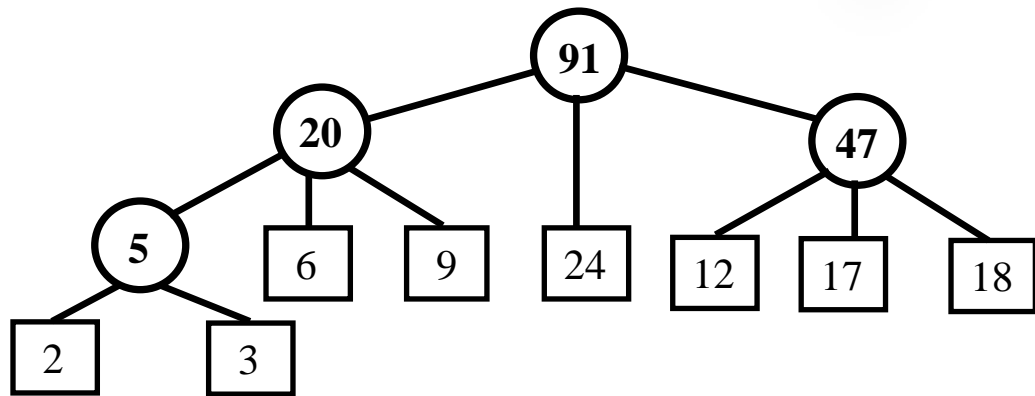
Huffman树的应用

- Huffman编码适合于
字符 **频率不等，差别较大的情况**
- 数据通信的二进制编码
 - 不同的频率分布，会有不同的压缩比率
 - 大多数的商业压缩程序都是采用几种编码方式以应付各种类型的文件
 - Zip 压缩就是 LZ77 与 Huffman 结合
- 归并法外排序，合并顺串

5.6 Huffman树及其应用

思考

- 当外部的数目不能构成满 b 叉 Huffman 树时，需附加多少个权为 0 的“虚”结点？请推导
- R 个外部结点， b 叉树
 - 若 $(r-1) \% (b-1) = 0$ ，则不需要加“虚”结点
 - 否则需要附加 $b - (r-1) \% (b-1) - 1$ 个“虚”结点
 - 即第一次选取 $(r-1) \% (b-1) + 1$ 个非 0 权值
- 试调研常见压缩软件所使用的编码方式





张铭《数据结构与算法》



数据结构与算法

谢谢聆听

国家精品课 “数据结构与算法”

<http://www.jpk.pku.edu.cn/pkujpk/course/sjig/>

张铭，王腾蛟，赵海燕

高等教育出版社，2008. 6。 “十一五” 国家级规划教材