

算法设计与分析

第4章 贪心算法 (2)

谢晓芹

哈尔滨工程大学计算机科学与技术学院



学习要点

- 通过应用范例学习贪心设计策略
 - 哈夫曼编码
 - 最优装载问题
 - 单源最短路径
 - 最小生成树

哈夫曼编码

■ 文件压缩问题

- 输入：给定符号表和一段数据
- 输出：给每个符号赋一个编码，使给定数据在编码下的长度最短

■ 二进制字符编码：每个字符用一个二进制0、1 串来表示

- 固定长编码
每个字符都用相同长的0、1 串表示
- 可变长编码
经常出现的字符用短码，不经常出现的用长码

■ 哈夫曼编码

- 用于数据文件压缩的十分有效的编码方法。压缩率通常在20%~90%之间
- 用字符在文件中出现的频率表来建立一个用0，1串表示各字符的最优表示方式

哈夫曼编码

	a	b	c	d	e	f
频率(千次)	45	13	12	16	9	5
8位ASCII码	01100001	01100010	01100011	01100100	01100101	01100110
定长码	000	001	010	011	100	101
变长码	0	101	100	111	1101	1100

■ 定长码

$$-45 \times 3 + 13 \times 3 + 12 \times 3 + 16 \times 3 + 9 \times 3 + 5 \times 3 = 300$$

■ 变长码

$$-45 \times 1 + 13 \times 3 + 12 \times 3 + 16 \times 3 + 9 \times 4 + 5 \times 4 = 224$$

➤ 频率越高，其编码越短，频率越低，编码越长

减少约
25%

哈夫曼编码

- 可变长码的问题？
 - 译码歧义：编码后的消息可能有多种译码方式
- 前缀码
 - 对每一个字符规定一个0,1串作为其编码，并且任一字符的编码都不是其它字符的编码的前缀。这种编码称为具有前缀性质（或者前缀码）
- 例如
 - 001011101分解为0,0,101,1101, aabe

	a	b	c	d	e	f
频率(千次)	45	13	12	16	9	5
定长码	000	001	010	011	100	101
变长码	0	101	100	111	1101	1100

哈夫曼编码

- 通常用**二叉树**来作为前缀编码的数据结构

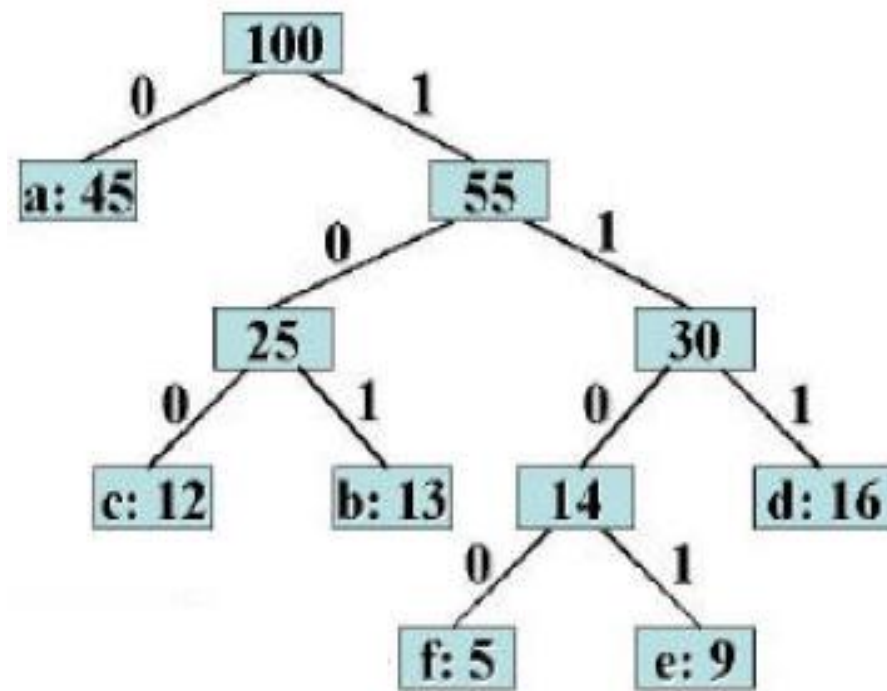
- 编码树

- 叶结点：用**字符**及其出现频率标记。
- 内结点(包括根结点)：其子树的叶子的频率和。
- 边标记：
一个结点的**左侧**边标记0，**右侧**边标记1。

- 前缀码

- 从树根到树叶的一条路径

- 这种树唯一吗？ 怎么构建？



哈夫曼编码 ✨

- C:编码字符集, C中任一字符以频率 $f(c)$ 在数据文件中出现.前缀编码方案对应于二叉树T, 字符c在T中的深度为 $d_T(c)$
- 最优编码树问题
 - 输入: 字母表 $C=\{c_1, c_2, \dots, c_n\}$, 频率表 $F=\{f(c_1), f(c_2), \dots, f(c_n)\}$
 - 输出: C的具有最小平均码长 $B(T)$ 的**前缀编码树**,即哈夫曼树
- 分析
 - (编码方案)**平均码长**定义为:
$$B(T) = \sum_{c \in C} f(c) d_T(c)$$
 - 最优前缀码
 - 使**平均码长达到最小**的前缀码编码方案称为给定编码字符集C的最优前缀码
 - 如何实现: 哈夫曼树
 - 表示最优前缀码的二叉树
 - 总是一棵完全二叉树, 即树中任一结点都有2个儿子结点

哈夫曼编码

■ 哈夫曼编码

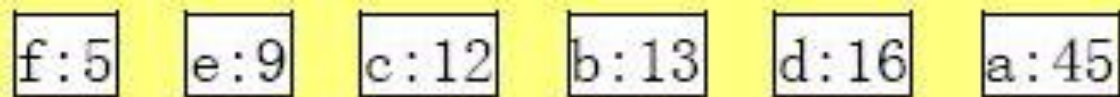
- 哈夫曼提出构造最优前缀码的**贪心算法**，产生的编码方案称为哈夫曼编码
- 基本思想
 - 将字符按照出现频率从高到低排序
 - 给出现**频率高的字符较短**的编码，出现频率较低的字符以**较长**的编码，以大大缩短总码长

■ 频率最高？ 排序？

■ 最短编码？ 哈夫曼树

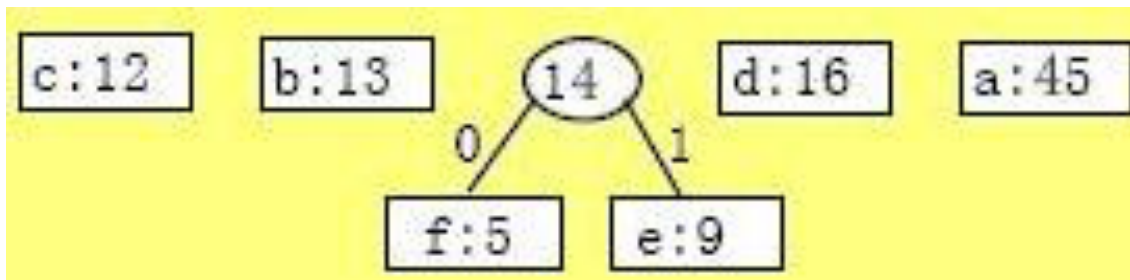
哈夫曼编码

初始：



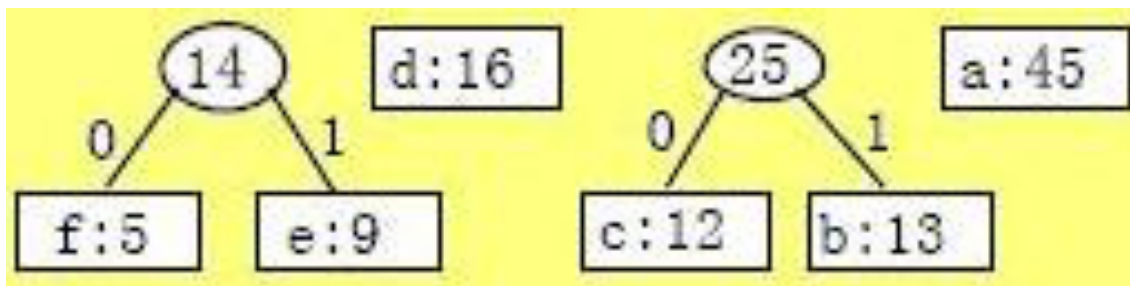
贪什么？

第1步：

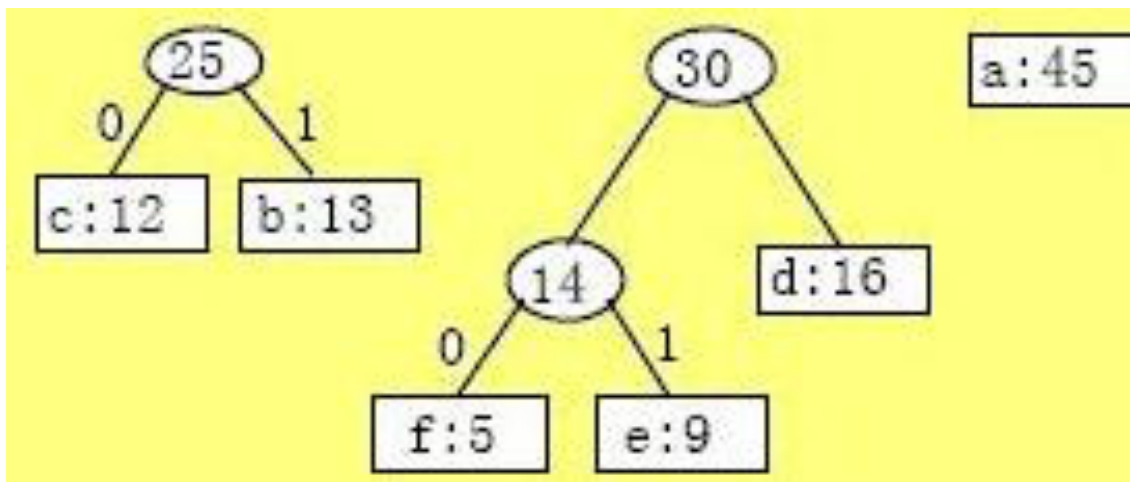


怎么贪？

第2步：

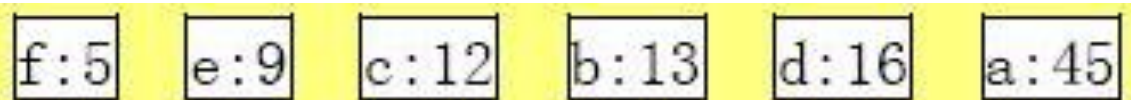


第3步：

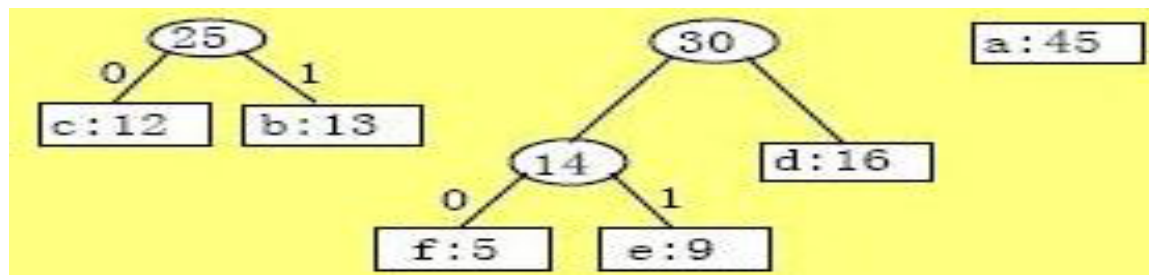


哈夫曼编码

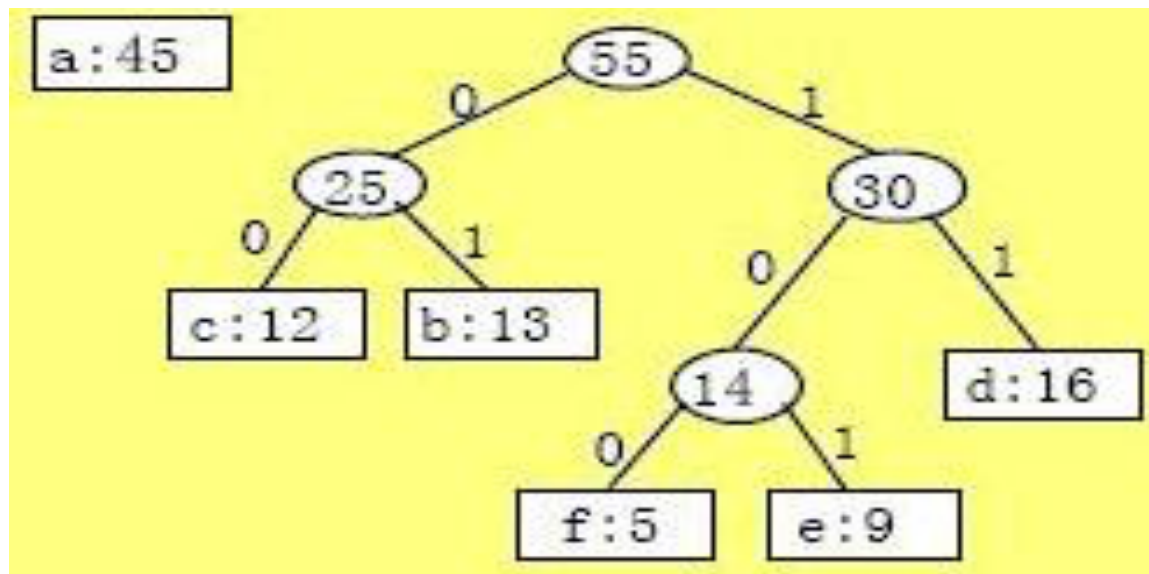
初始：



第3步：



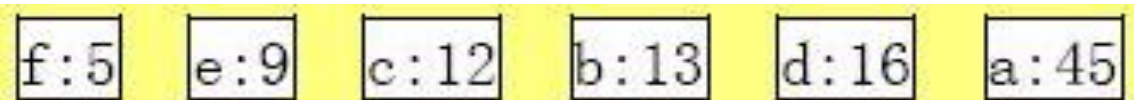
第4步：



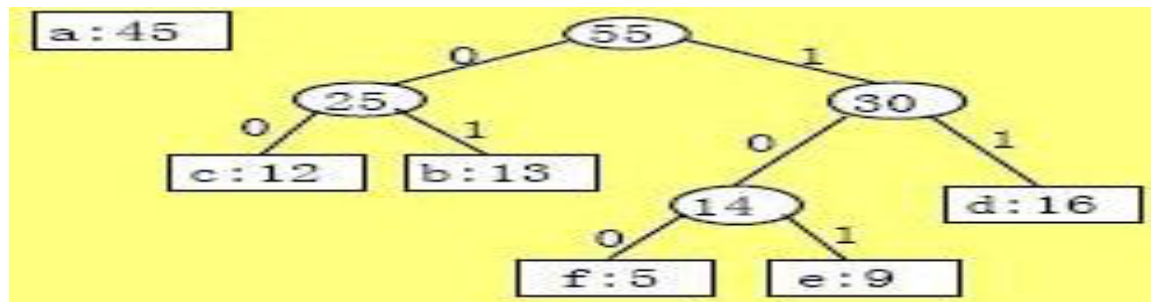
何时结束？

哈夫曼编码

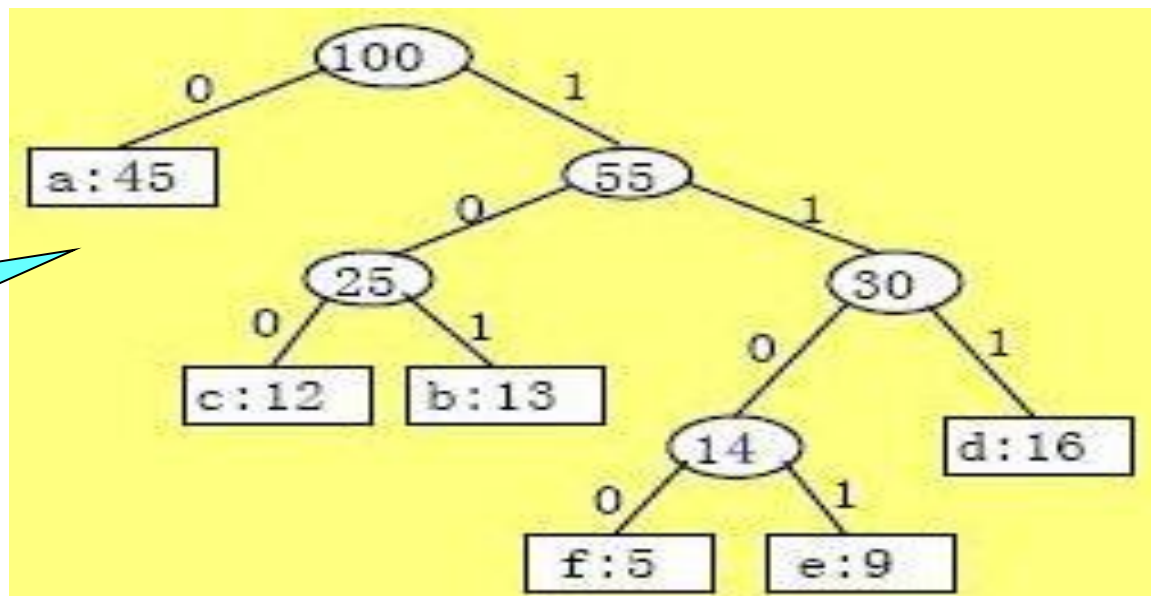
初始：



第4步：



第5步：



编码

选取最小
频率子树

哈夫曼编码

- 哈夫曼算法以**自底向上**的方式构造表示最优前缀码的二叉树T
 - 以 $|C|$ 个叶结点开始, 执行 $|C| - 1$ 次的“合并”运算后产生最终所要求的树T
- 过程
 - 贪什么： 频率最小的两个字符
 - 怎么贪
 - **合并**为一棵树
 - 根：频率和， 叶子：两个字符
 - 何时终止： $n-1$ 个结点全部生成
 - 编码
 - 左边为0， 右边为1
 - 从根到叶结点上所有赋值的连接



哈夫曼树

什么数据结构,方便插入和删除操作?

算法设计技术: 利用堆来管理算法执行中的信息

- 堆: n 个元素的序列 $\{k_1, k_2, \dots, k_n\}$,当且仅当满足以下关系时,称之为堆.

$$\begin{cases} k_i \leq k_{2i} \\ k_i \leq k_{2i+1} \end{cases} \quad \begin{cases} k_i \geq k_{2i} \\ k_i \geq k_{2i+1} \end{cases} \quad (i = 1, 2, \dots, \lfloor n/2 \rfloor)$$

- 含义:

- 完全二叉树中所有非终端结点的值不大于(或不小于)其左右孩子结点的值.
- 堆顶元素必为序列中 n 个元素的最小值(或最大值)

哈夫曼编码

n个字符的集合

频度 $f(c_1), f(c_2), \dots, f(c_n)$

Huffman(C, F) (使用堆操作实现)

$n \leftarrow |C|;$

$Q \leftarrow C;$ /* 用**BUILD-HEAP** 建立堆 */ $O(n)$

for $i \leftarrow 1$ to $n-1$ { //终止条件：n-1个节点全部生成 $n-1$

$z \leftarrow \text{Allocate-Node}();$

$x \leftarrow \text{left}[z] \leftarrow \text{Extract-MIN}(Q);$ /* 堆操作*/ $O(\log n)$

$y \leftarrow \text{right}[z] \leftarrow \text{Extract-MIN}(Q);$ /* 堆操作*/ $O(\log n)$

$f(z) \leftarrow f(x) + f(y);$

$\text{insert}(Q, z);$ } /* 堆操作*/

Return $\text{Extract-MIN}(Q).$

哈夫曼编码

- 哈夫曼算法的计算时间为 $O(n \log n)$
- 分析
 - 初始化优先队列需要 $O(n)$ 计算时间
 - 最小堆的 ExtractMin 和 Insert 运算均需 $O(\log n)$ 时间
 - $n-1$ 次的合并总共需要 $O(n \log n)$ 计算时间。

```
Huffman(C, F) {  
    n ← |C|;  
    Q ← C;      //用BUILD-HEAP 建立堆  
    for i ← 1 to n-1 {    //终止条件  
        z ← Allocate-Node( );  
        x ← left[z] ← Extract-MIN(Q); //堆操作  
        y ← right[z] ← Extract-MIN(Q); //堆操作  
        f(z) ← f(x)+f(y);  
        insert(Q,z); }      //堆操作  
    Return Extract-MIN(Q)}
```

哈夫曼编码

■ 哈夫曼算法的正确性

- 证明哈夫曼算法的正确性，只要证明最优前缀码问题具有如下两个性质：
 - (1)贪心选择性质
 - (2)最优子结构性质。

哈夫曼编码 - 贪心性质证明 ★

- **引理1(贪心选择性)**. 设 C 是字母表, $c \in C$, c 具有频率 $f(c)$, x 、 y 是 C 中具有最少频率的两个字符, 则存在一个 C 的最优前缀树, x 与 y 的编码具有相同长度, 且仅在最末一位不同
- 分析
 - 贪心选择 是贪频率最小的两个字符 → 树上兄弟结点, 左右不同而已
 - 哈夫曼编码得到一个最优解, 是以贪心选择开始

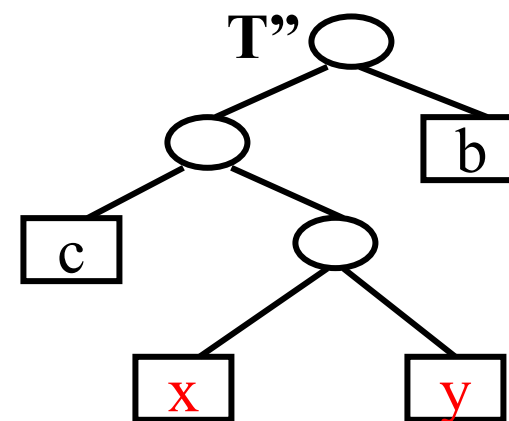
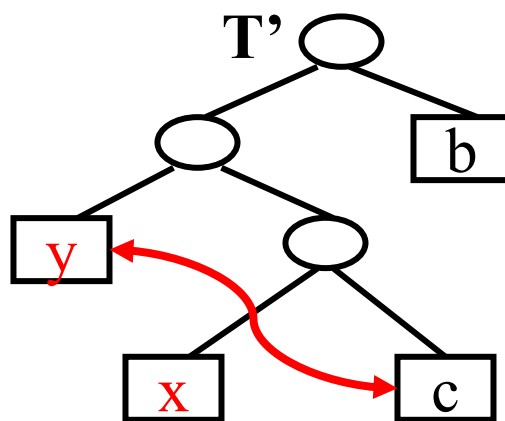
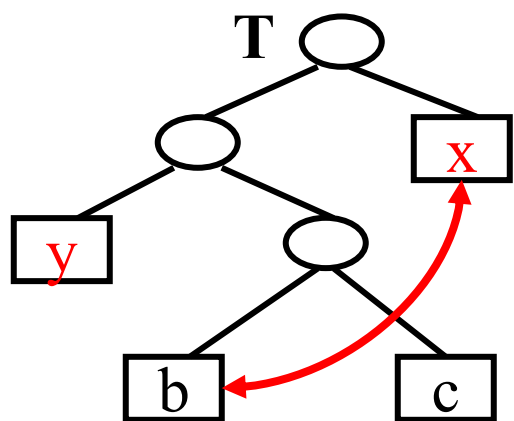
哈夫曼编码 - 贪心性质证明 ★

- **引理1(贪心选择性)**. 设 C 是字母表, $c \in C$, c 具有频率 $f(c)$, x 、 y 是 C 中具有最少频率的两个字符, 则存在一个 C 的最优前缀树, x 与 y 的编码具有相同长度, 且仅在最末一位不同
- **证明**: 设 T 是一个 C 的最优前缀树, 且 b 和 c 是具有最大深度且为兄弟的两个字符 不失一般性, 设 $f(b) \leq f(c)$, $f(x) \leq f(y)$.

因 x 与 y 是具有最低频率的字符, $f(b) \geq f(x)$, $f(c) \geq f(y)$

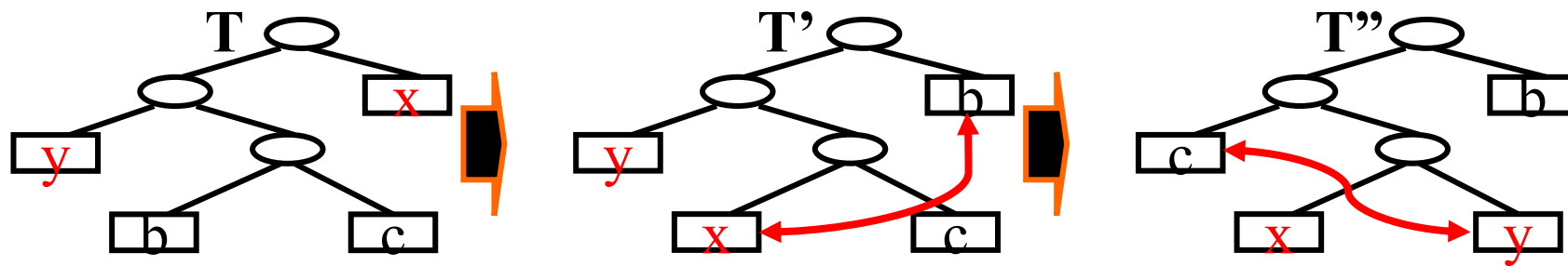
从 T 构造 T' , 交换 T 的 b 和 x ;

从 T' 构造 T'' , 交换 T' 的 y 和 c ;



哈夫曼编码 - 贪心性质证明

T'' 是最优前缀树.



➤ **证:** $B(T) - B(T') = \sum_{c \in C} f(c) d_T(c) - \sum_{c \in C} f(c) d_{T'}(c)$

$$\begin{aligned} &= f(x) d_T(x) + f(b) d_T(b) - f(x) d_{T'}(x) - f(b) d_{T'}(b) \\ &= f(x) d_T(x) + f(b) d_T(b) - f(x) d_{T'}(b) - f(b) d_{T'}(x) \\ &= (f(b) - f(x)) (d_T(b) - d_T(x)) \end{aligned}$$

$\because f(b) \geq f(x), d_T(b) \geq d_T(x)$ (因为b 的深度最大)

$$\therefore B(T) - B(T') \geq 0, B(T) \geq B(T')$$

同理可证 $B(T') \geq B(T'')$. 于是 $B(T) \geq B(T'')$.

由于T 是最优化的, 所以 $B(T) \leq B(T'')$.

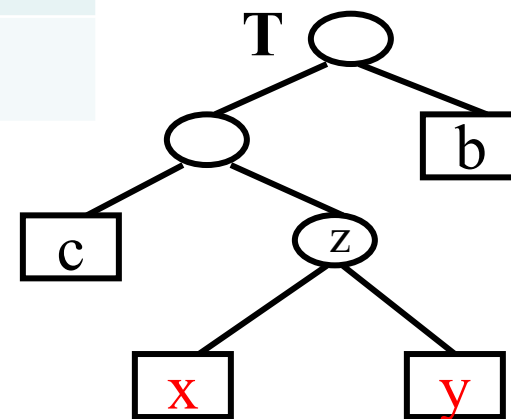
于是, $B(T) = B(T'')$, 所以T''是C 的最优化前缀编码树.

在T''中, x 和y 具有相同长度编码, 而且仅最后一位不同。

哈夫曼编码 - 最优子结构性质证明 ✨

- 引理2(最优子结构). 设 T 是字母表 C 的最优前缀树, $c \in C$, $f(c)$ 是 c 在文件中出现的频率。设 x 、 y 是 T 中任意两个相邻叶结点, z 是它们的父结点, 则 z 看作为频率为 $f(z) = f(x) + f(y)$ 的字符, $T' = T - \{x, y\}$ 表示了字母表 $C' = C - \{x, y\} \cup \{z\}$ 的最优前缀编码树
- 分析

		最优解
原问题	C	T
子问题	$C' = \{C - \{x, y\}\} \cup \{z\}$	$T' = T - \{x, y\}$



哈夫曼编码 - 最优子结构性质证明 ✨

- 引理2(最优子结构). 设 T 是字母表 C 的最优前缀树, $c \in C$, $f(c)$ 是 c 在文件中出现的频率。设 x 、 y 是 T 中任意两个相邻叶结点, z 是它们的父结点, 则 z 看作为频率为 $f(z) = f(x) + f(y)$ 的字符, $T' = T - \{x, y\}$ 表示了字母表 $C' = C - \{x, y\} \cup \{z\}$ 的最优前缀编码树。

➤ 证: 先证 $B(T) = B(T') + f(x) + f(y)$

$\forall v \in C - \{x, y\}, d_T(v) = d_{T'}(v)$ 。于是 $f(v)d_T(v) = f(v)d_{T'}(v)$ 。

$$B(T) = \sum_{v \in C - \{x, y\}} f(v)d_T(v) + f(x)d_T(x) + f(y)d_T(y)$$

$$= \sum_{v \in C - \{x, y\}} f(v)d_T(v) + (f(x) + f(y))(d_{T'}(z) + 1)$$

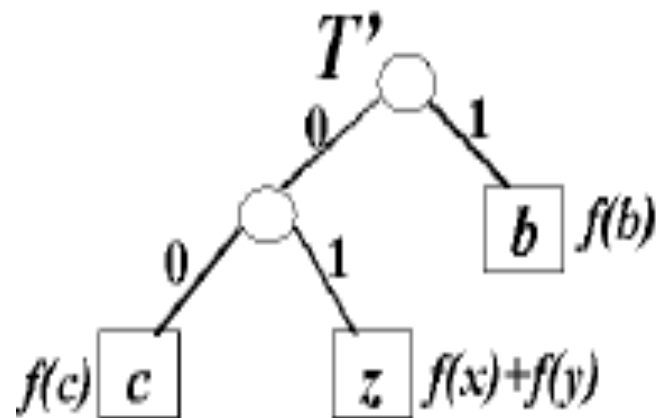
$$B(T') = \sum_{v \in C - \{x, y\}} f(v)d_{T'}(v) + f(z)d_{T'}(z)$$

由于 $f(x) + f(y) = f(z)$, $f(x)d_T(x) + f(y)d_T(y) = f(z)d_{T'}(z) + (f(x) + f(y))$ 。

于是 $B(T) = B(T') + f(x) + f(y)$

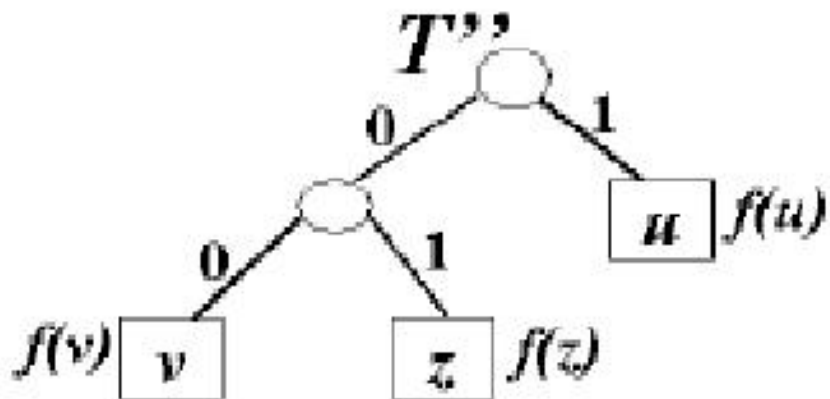
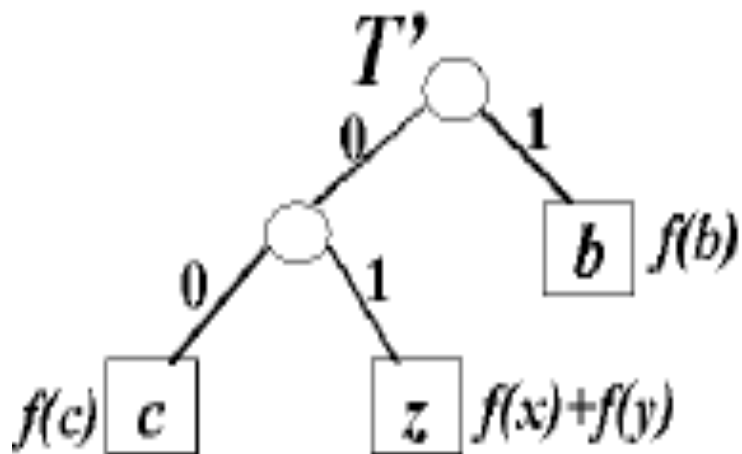
由于 $d_T(x) = d_T(y) = d_{T'}(z) + 1$,
 $f(x)d_T(x) + f(y)d_T(y) = (f(x) + f(y))(d_{T'}(z) + 1)$
 $= (f(x) + f(y))d_{T'}(z) + (f(x) + f(y))$

反证法: 若 T' 是 C' 的非最优前缀编码树, 则必存在 T'' , 使 $B(T'') < B(T')$ 。



哈夫曼编码 - 最优子结构性证明

反证法：若 T' 是 C' 的非最优前缀编码树，则必存在 T'' ，使 $B(T'') < B(T')$ 。

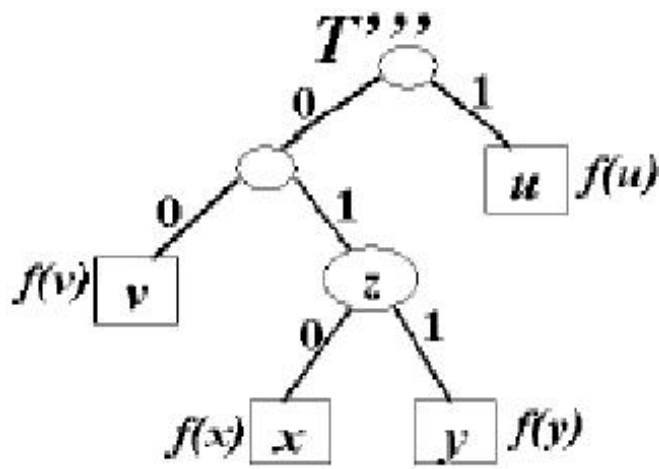


证: (1) 因为 z 在 C' 中是一个字符，它必是 T'' 中的叶子。把结点 x 与 y 加入 T'' ，作为 z 的子结点，则得到 C 的一个如下的前缀编码树 T'''

(2) T''' 代价为：

$$\begin{aligned} B(T''') &= \dots + (f(x)+f(y))(d_{T''}(z)+1) \\ &= \dots + f(z)d_{T''}(z) + (f(x)+f(y)) \\ &= B(T'') + f(x) + f(y) \\ &< B(T') + f(x) + f(y) = B(T) \end{aligned}$$

(3) 与 T 是最优的相矛盾，故 T' 是最优的



哈夫曼编码

- 定理. Huffman 算法产生一个最优前缀编码树。
- 证:
由于引理1、引理2 成立,而且Huffman 算法按照引理1 的贪心选择性确定的规则进行局部优化选择, 所以Huffman 算法产生一个最优前缀编码树。



哈夫曼编码

■ 实例

- 1.周杰伦(Z) , $f(Z)=20$
 - 2.刘德华(L) , $f(L)=7$
 - 3.韩愈(H) , $f(H)=10$
 - 4.李小龙(D) , $f(D)=4$
 - 5.易中天(Y) , $f(Y)=18$
- 构造哈夫曼编码树 , 并计算编码树的平均码长 ?

使用贪心策略的典型步骤

- 步骤1 将原问题表述为一个做出一个选择后剩下唯一的一个子问题的形式
- 步骤2 证明在所有的最优选择里面总有一个是贪心选择
- 步骤3 证明贪心选择加上对剩下子问题的最优解就是原问题的最优解

分治、动态规划和贪心的比较 ✨

	标准分治	动态规划	贪心
适用类型	通用问题	优化问题	优化问题
子问题结构	每个子问题不同	很多子问题重复	只有一个子问题
最优子结构	不要求	必须满足	必须满足
解决子问题数	全部子问题都要解决	全部子问题都要解决	只要解决一个子问题
子问题依赖性	先选择后解决子问题	先解决子问题后选择	先选择后解决子问题
求解顺序	递归与递推都可	从底往上	从顶往下

小结

- 哈夫曼编码
 - 问题描述
 - 求解算法
- 重点
 - 贪心算法的证明过程
 - 方法
 - 构造法，反证法，反证+构造法，数学归纳法
- 难点
 - 怎么表述贪心选择性质
 - 怎么表述最优子结构性性质

实践创新能力训练

- 一辆汽车加满油后可以行使 n 公里。旅途中有若干加油站。请设计一个有效算法，指出应在哪些加油站停靠加油，使得沿途加油次数最少，并证明算法能产生一个最优解。
- 提示：对于给定的 n 和 k 个加油站位置，计算最少加油次数。
- 输入数据： $n, k, d[k]$ 表示第 k 个加油站到第 $k-1$ 个加油站之间的距离。第 0 个加油站表示出发地，汽车已加满油，第 $k+1$ 个加油站表示目的地。
- 输出： (x_1, x_2, \dots, x_k) ， $x_i=1$ 表示停留， $x_i=0$ 表示不停
加油策略 $\text{strat}[]$ ($\text{strat}[i]=j$ 表示第 i 次加油是在第 j 个加油站)