



HIT
CS&E

第五章

贪心算法的设计与分析 原理

程思瑶

计算机科学与技术学院



- 5.1 贪心算法的要素
- 5.2 活动选择(activity-selection)问题
- 5.3 Huffman编码
- 5.4 最小生成树问题
- 5.5 贪心算法的理论基础
- 5.6 任务调度问题



HIT
CS&E

参考资料



Introduction to Algorithms

Chapter 16

Pages 370-405



5.1 贪心算法的要素

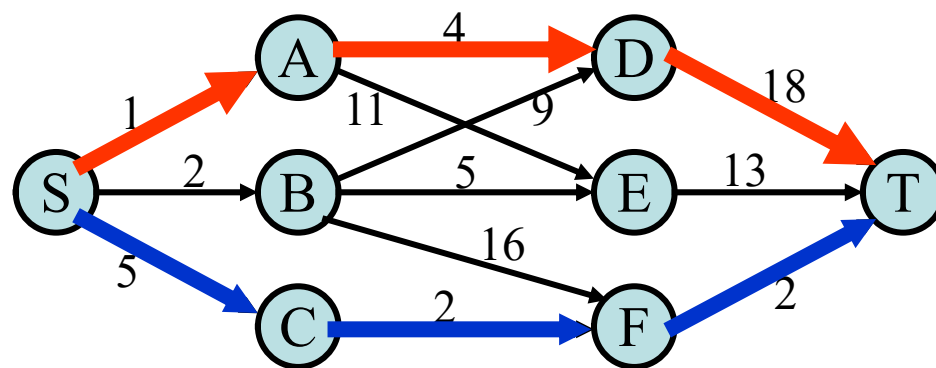
- Greedy算法的基本概念
- Greedy算法与动态规划方法的比较
- Greedy算法的设计步骤



Greedy算法的基本概念

- Greedy算法的实例

- 最短路径问题



- Greedy求解过程

Greedy算法不一定产生正确解



- Greedy算法的基本思想
 - 求解最优化问题的算法包含一系列步骤
 - 每一步都有一组选择
 - 作出在当前看来最好的选择
 - 希望通过作出局部优化选择达到全局优化选择
 - Greedy算法不一定总产生优化解
- Greedy算法产生优化解的条件
 - 优化子结构
 - Greedy选择性(Greedy-choice property)



- Greedy选择性

一个优化问题的全局优化解可以通过局部优化选择得到。

- Greedy选择性需证明

- 归纳法

- 对算法步数归纳或问题规模归纳

- 交换论证法

- 在保证最优性不变前提下，从一个最优解逐步替换，最终得到贪心算法的解



HIT
CS&E

- 优化子结构

若一个优化问题的优化解包含它的(剩余)
子问题的优化解, 则称其具有优化子结构



与动态规划方法的比较

- 动态规划方法
 - 在每一步所做的选择通常依赖于子问题的解
 - 以自底向上方式，先解小子问题，再求解大子问题
- Greedy方法
 - 在每一步先做出当前看起来最好的选择
 - 然后再求解本次选择后产生的剩余子问题
 - 每次选择既不依赖于子问题的解，也不依赖于未来的选择
 - 以自顶向下方式，逐步进行贪心选择，不断减少子问题规模



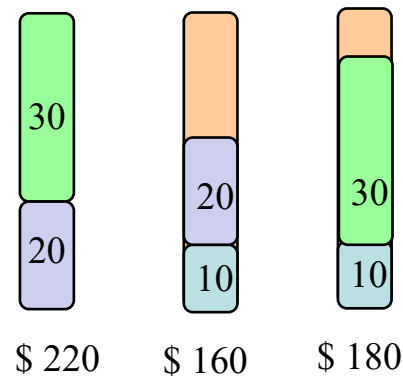
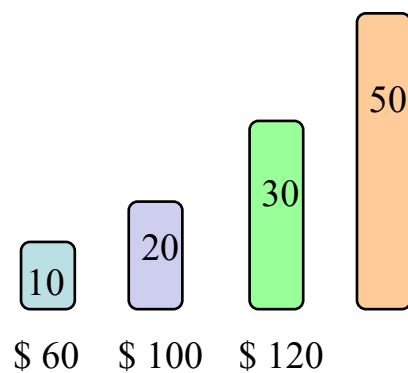
与动态规划方法的比较

- 动态规划方法可用的条件
 - 优化子结构
 - 子问题重叠性
- Greedy方法可用的条件
 - 优化子结构
 - Greedy选择性
- 可用Greedy方法时，动态规划方法可能不适用
- 可用动态规划方法时，Greedy方法可能不适用

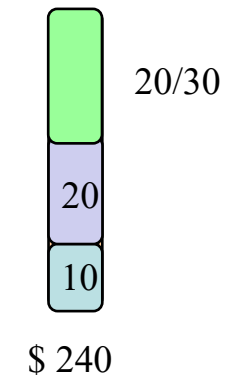


与动态规划方法的比较

- 例如：0-1 背包问题与部分背包问题
 - 都具有优化子结构
 - 但是，部分背包问题可用贪心策略解决，而0-1背包问题却不行！
 - 计算每个物品每磅价值 v_i/w_i ，并按照每磅价值由大到小顺序取物品



0-1背包问题



部分背包问题



准确Greedy算法的设计步骤

1. 设计贪心选择方法：

- 贪心选择方法
- 剩余子问题

很重要！
决定能否得到
全局最优解

2. 证明：对于1中所求解的问题具有优化子结构

3. 证明：对于1中所求解的问题具有贪心选择性

4. 按照1中设计的贪心选择方法设计算法

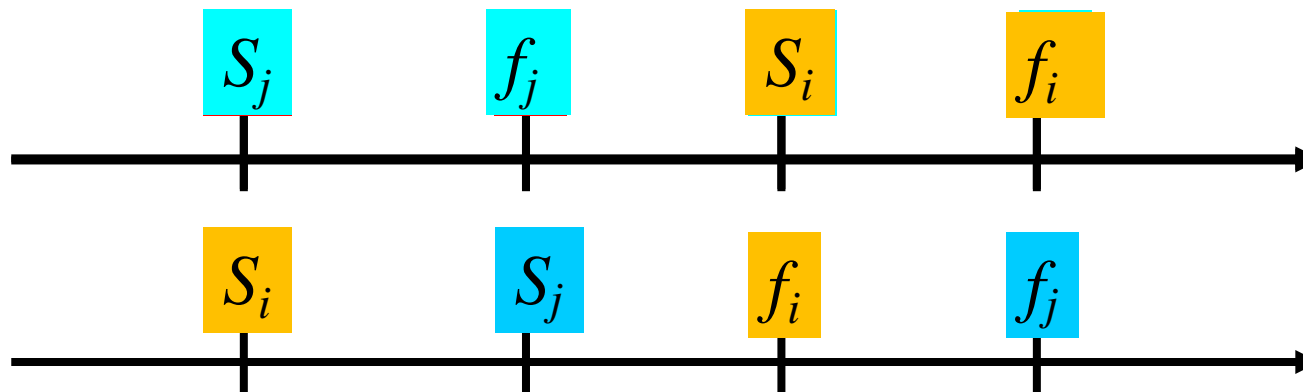


5.2 活动选择 (activity-selection) 问题

- 问题定义
- 问题求解
 - 设计贪心选择方法
 - 优化解的结构分析
 - Greedy选择性证明
 - 算法设计
 - 算法复杂性分析



- 活动
 - 设 $S = \{1, 2, \dots, n\}$ 是 n 个活动的集合，所有活动共享一个资源，该资源同时只能为一个活动使用
 - 每个活动 i 有起始时间 s_i ，终止时间 f_i ， $s_i \leq f_i$
- 相容活动
 - 活动 i 和 j 是相容的，若 $s_j \geq f_i$ 或 $s_i \geq f_j$ ，即





- 活动选择问题定义

- 输入: $S = \{1, 2, \dots, n\}$,

$$F = \{ [s_i, f_i] \}, \quad n \geq i \geq 1$$

- 输出: S 的最大相容活动集合

贪心思想:

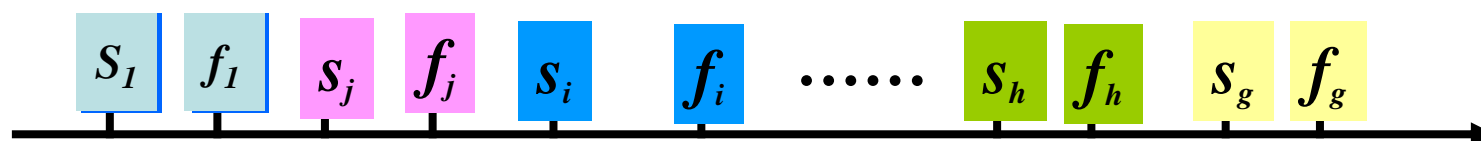
为了选择最多的相容活动, 每次选 f_i 最小的活动, 使我们能够选更多的活动

剩余子问题: $S_i = \{ j \in S \mid s_j \geq f_i \}$



引理1 设 $S=\{1,2,\dots,n\}$ 是 n 个活动集合, $[s_i, f_i]$ 是活动 i 的起始终止时间, 且 $f_1 \leq f_2 \leq \dots \leq f_n$, S 的活动选择问题的某个优化解包括活动1.

证 设 A 是一个优化解, 按结束时间排序 A 中活动, 设其第一个活动为 k , 第二个活动为 j



如果 $k=1$, 引理成立.

如果 $k \neq 1$, 令 $B=A-\{k\} \cup \{1\}$,

由于 A 中活动相容, $f_1 \leq f_k \leq s_j$, B 中活动相容.

因为 $|B|=|A|$, 所以 B 是一个优化解, 且包括活动1.



定理1. 设 $S = \{1, 2, \dots, n\}$ 是 n 个活动集合, $[s_i, f_i]$ 是活动 i 的起始终止时间, 且 $f_1 \leq f_2 \leq \dots \leq f_n$, 设 A 是 S 的调度问题的一个优化解且包括活动 1 , 则 $A' = A - \{1\}$ 是 $S' = \{i \in S / s_i \geq f_1\}$ 的调度问题的优化解.

证. 显然, A' 中的活动是相容的.

我们仅需要证明 A' 是最大的.

定理1说明活动选择问题具有优化子结构

B 是 S 的一个解.

由于 $|A| = |A'| + 1$, $|B| = |B'| + 1 > |A'| + 1 = |A|$, 与 A 最大矛盾.



定理2. 设 $S=\{1, 2, \dots, n\}$ 是 n 个活动集合, $f_1 \leq f_2 \leq \dots \leq f_n$, $f_{l_0}=0$, l_i 是 $S_i = \{j \in S \mid s_j \geq f_{l_{i-1}}\}$ 中具有最小结束时间 f_{l_i} 的活动. 设 A 是 S 的包含活动 1 的优化解, 则 $A = \bigcup_{i=1}^k \{l_i\}$

证. 对 $|A|$ 作归纳法.

$S_1 = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11\}$, $l_0=0$

当 $|A|=1$ 时, 由引理1, 命题成立.

$S_2 = \{4, 6, 7, 8, 9, 11\}$, $l_1=1$

设 $|A| \leq k-1$ 时, 命题成立.

$S_3 = \{8, 9, 11\}$, $l_2=4$

当 $|A|=k$ 时, 由定理1, $A = \{l_1=1\} \cup A_1$.

$S_4 = \{11\}$, $l_3=8$

A_1 是 $S_2 = \{j \in S \mid s_j \geq f_{l_1} = f_1\}$ 的优化解.

$S_1 = \{j \in S \mid s_j \geq f_{l_0} = 0\}$

$S_2 = \{j \in S \mid s_j \geq f_{l_1} = f_1\}$

$S_3 = \{j \in S \mid s_j \geq f_{l_2}\}$

.....

$S_k = \{j \in S \mid s_j \geq f_{l_{k-1}}\}$

$S_{k+1} = \{j \in S \mid s_j \geq f_{l_k}\} = \emptyset$

i	1	2	3	4	5	6	7	8	9	10	11
s_i	1	3	5	7	8	9	10	11	12	13	14
f_i	4	5	6	7	8	9	10	11	12	13	14

由归纳假设, $A_1 = \bigcup_{i=2}^k \{l_i\}$

于是, $A = \bigcup_{i=1}^k \{l_i\}$



- 贪心选择方法

- 选择:

- 每次选择具有最小结束时间的活动 f_i

- 剩余子问题:

- $S_i = \{j \in S \mid s_j \geq f_i\}$



- 算法

(设 $f_1 \leq f_2 \leq \dots \leq f_n$ 已排序)

Greedy-Activity-Selector(S, F)

$n \leftarrow \text{length}(S);$

$A \leftarrow \{1\}$

$j \leftarrow 1$

For $i \leftarrow 2$ To n Do

 If $s_i \geq f_j$

 Then $A \leftarrow A \cup \{i\}; j \leftarrow i;$

Return A



- 算法

(设 $f_1 \leq f_2 \leq \dots \leq f_n$ 已排序)

Greedy-Activity-Selector(S, F)

$n \leftarrow \text{length}(S);$

$A \leftarrow \{1\}$

$j \leftarrow 1$

For $i \leftarrow 2$ To n Do

 If $s_i \geq f_j$

 Then $A \leftarrow A \cup \{i\}; j \leftarrow i;$

Return A

- 如果结束时间已排序

$$T(n) = \theta(n)$$

- 如果结束时间未排序

$$T(n) = \theta(n) + \theta(n \log n) = \theta(n \log n)$$



定理3. Greedy-Activity-Selector 算法能够产生最优解.

证.

- (1) 由定理1可知活动选择问题具有**优化子结构**
- (2) 由定理2知贪心选择方法具有**Greedy选择性**
- (3) Greedy-Activity-Selector 算法**确实按照**定理2的**Greedy选择性**进行局部优化选择.



5.3 Huffman 编码

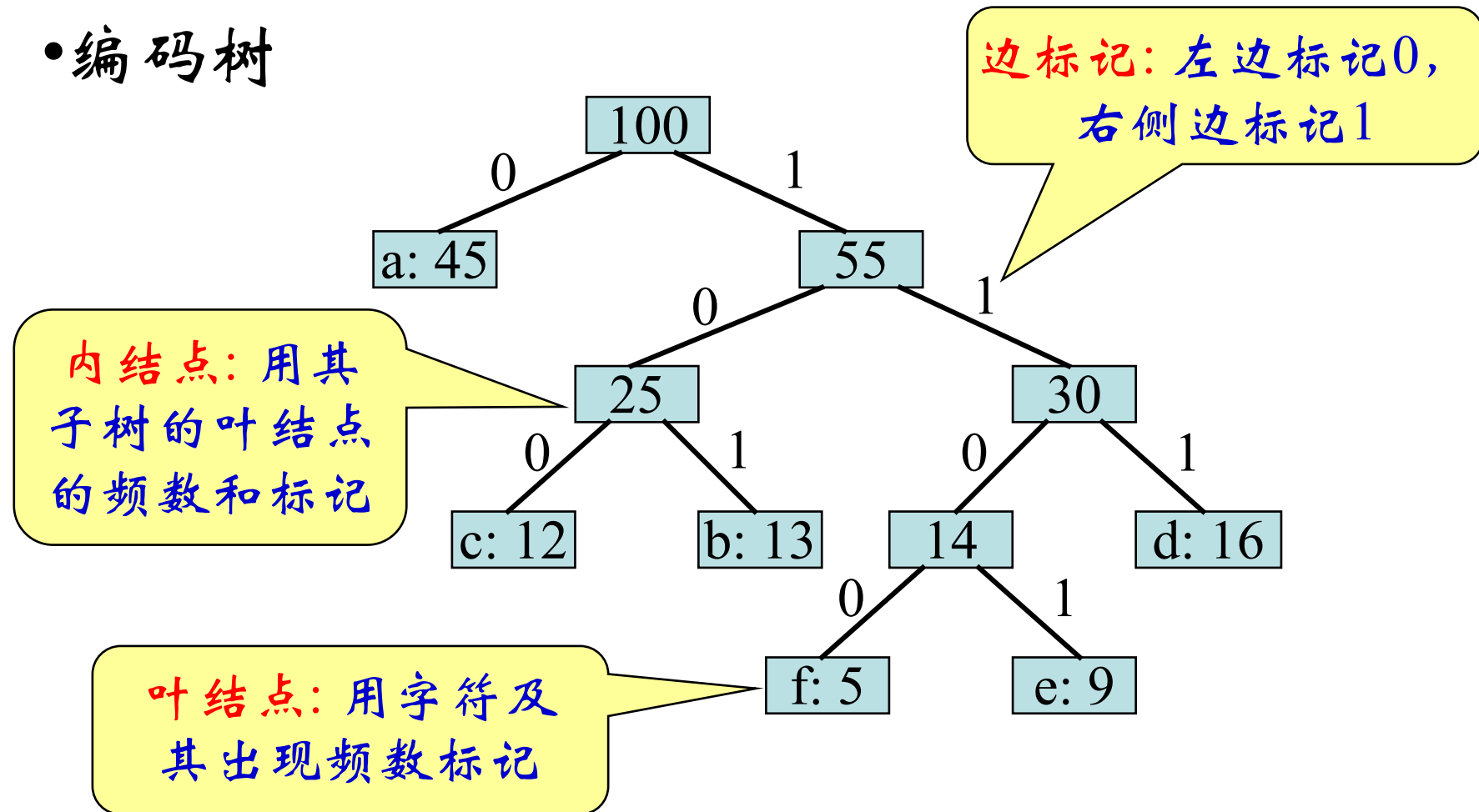
- 问题定义
- 问题求解
 - 设计贪心选择方法
 - 优化解的结构分析
 - Greedy选择性证明
 - 算法设计
 - 算法复杂性分析

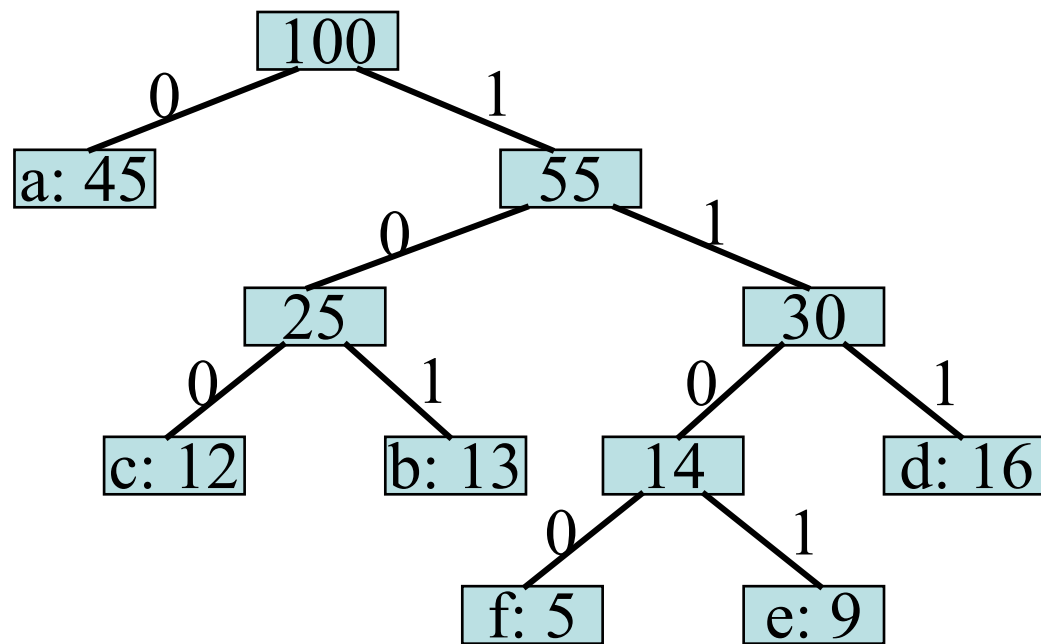


- 二进制字符编码
 - 每个字符用一个二进制0、1串来表示.
- 固定长编码
 - 每个字符都用相同长度的0、1串表示.
- 可变长编码
 - 经常出现的字符用短码, 不经常出现的用长码
- 前缀编码
 - 无任何字符的编码是另一个字符编码的前缀



• 编码树





- 编码树 T 的代价

- 设 C 是字母表(给定文件中的字母集合), $\forall c \in C$
- $f(c)$ 是 c 在文件中出现的频数
- $d_T(c)$ 是叶子 c 在树 T 中的深度, 即 c 的编码长度
- T 的代价是编码一个文件的所有字符的代码位数:

$$B(T) = \sum_{c \in C} f(c) d_T(c)$$



- 优化编码树问题

输入: 字母表 $C = \{c_1, c_2, \dots, c_n\}$,

频数表 $F = \{f(c_1), f(c_2), \dots, f(c_n)\}$

输出: 具有最小 $B(T)$ 的 C 前缀编码树

贪心思想:

循环地选择具有最低频数的两个结点,
生成一棵子树, 直至形成树

剩余子问题: ???

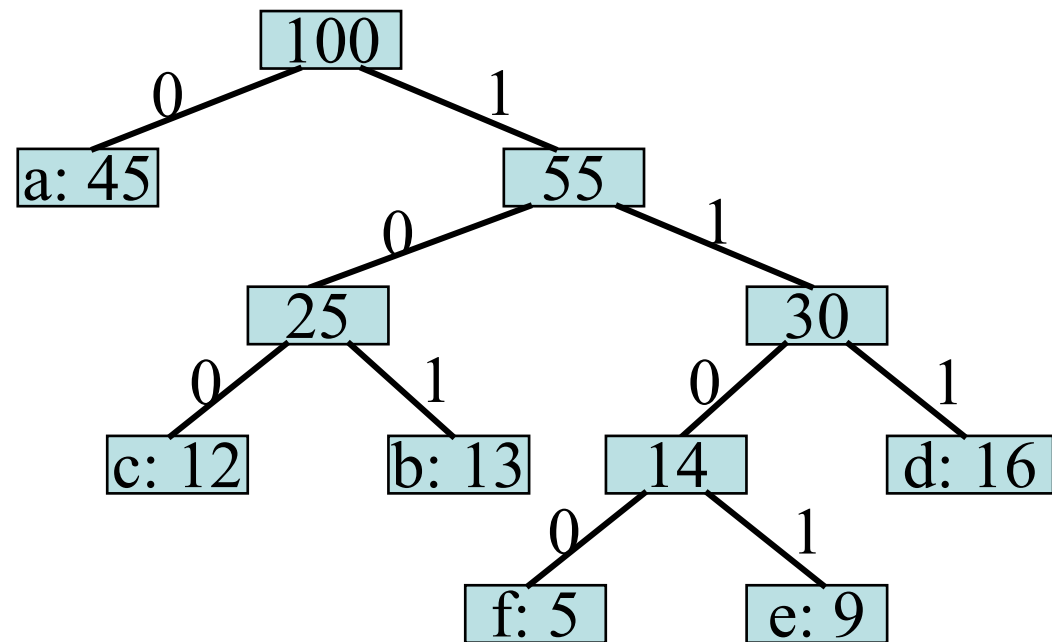


HIT
CS&E

贪心思想:

循环地选择具有最低频数的两个结点,
生成一棵子树, 直至形成树

f: 5 e: 9 c: 12 b: 13 d: 16 a: 45

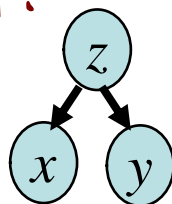




- 贪心选择方法

- 选择方法:

- 每次选择具有最低频数的两个节点 x 和 y ,
构造一个子树:



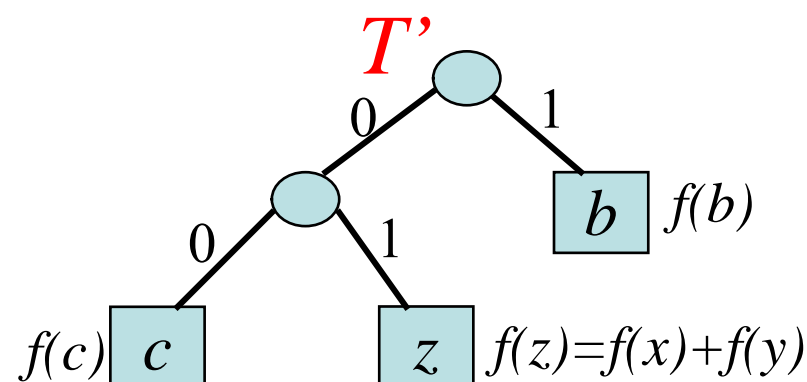
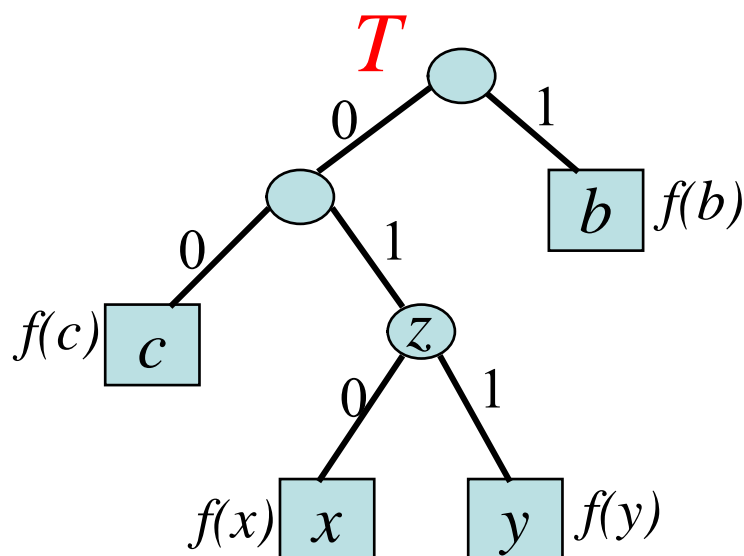
- 剩余子问题的结构:

- $C' = C - \{x, y\} \cup \{z\}$



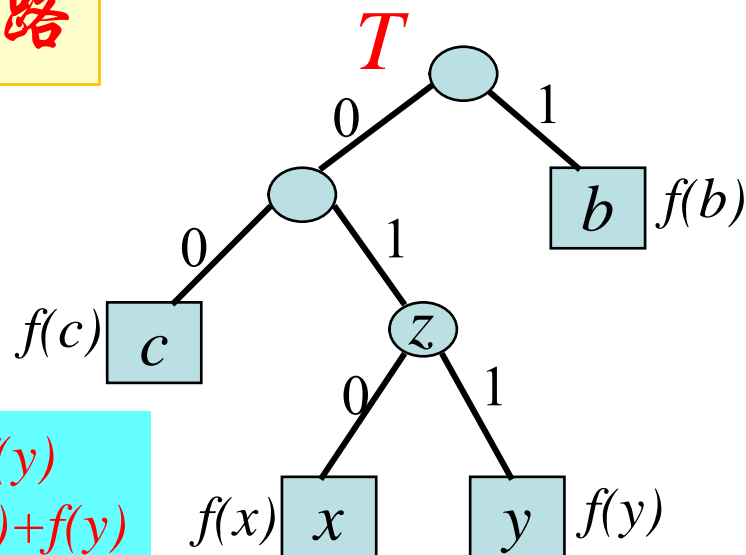
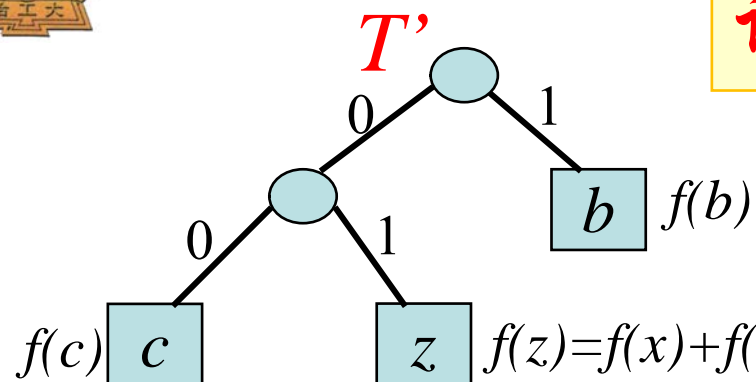
优化解的结构分析

引理1. 设 T 是字母表 C 的优化前缀树, $\forall c \in C$, $f(c)$ 是 c 在文件中出现的频数. 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 作为频数是 $f(z)=f(x)+f(y)$ 的字符, $T'=T-\{x,y\}$ 是字母表 $C'=C-\{x,y\} \cup \{z\}$ 的优化前缀编码树.

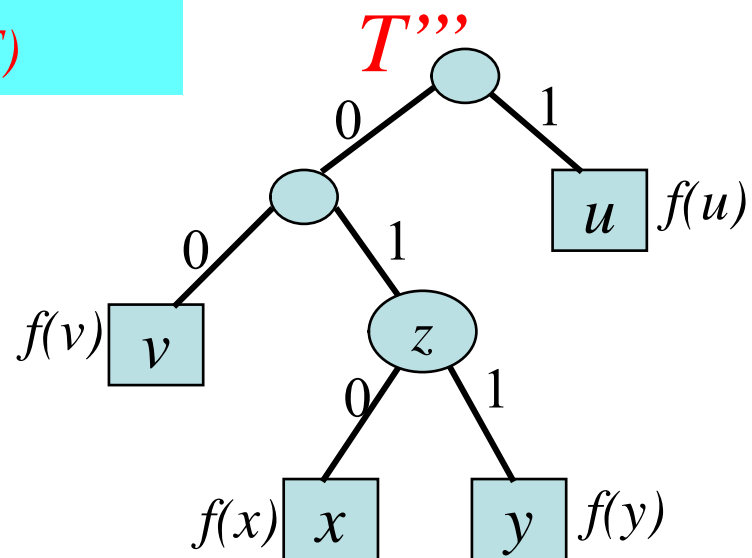
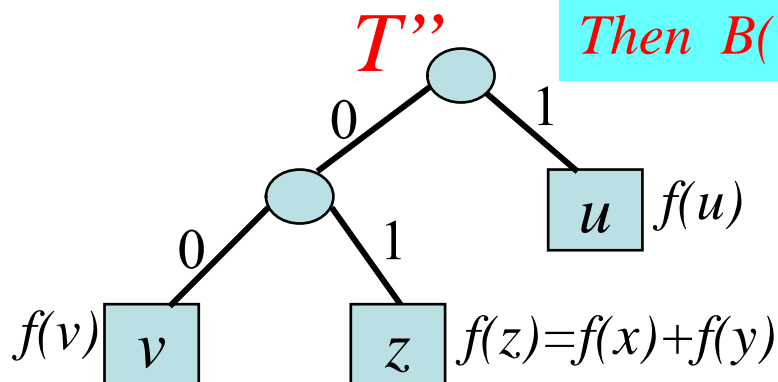




证明思路



$B(T) = B(T') + f(x) + f(y)$
 $B(T'') = B(T') + f(x) + f(y)$
If $B(T'') < B(T')$
Then $B(T'') < B(T)$





证. 往证 $B(T) = B(T') + f(x) + f(y)$.

$$B(T) = \sum_{c \in C} f(c) d_T(c)$$

$$= f(x) d_T(x) + f(y) d_T(y) + \sum_{c \in C' - \{z\}} f(c) d_T(c)$$

由于 $d_T(x) = d_T(y) = d_{T'}(z) + 1$

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

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

由于 $f(x) + f(y) = f(z)$

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

$$= B(T') + f(x) + f(y).$$

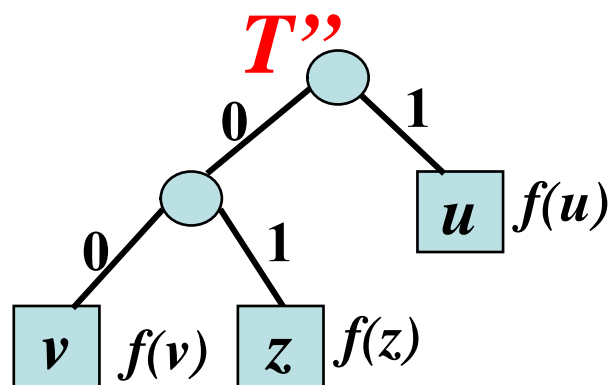
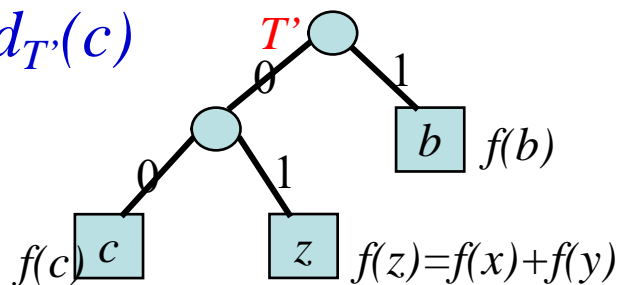
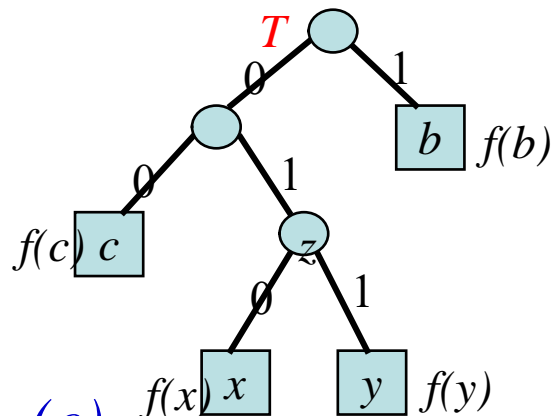
若 T' 不是 C' 的优化前缀编码树,

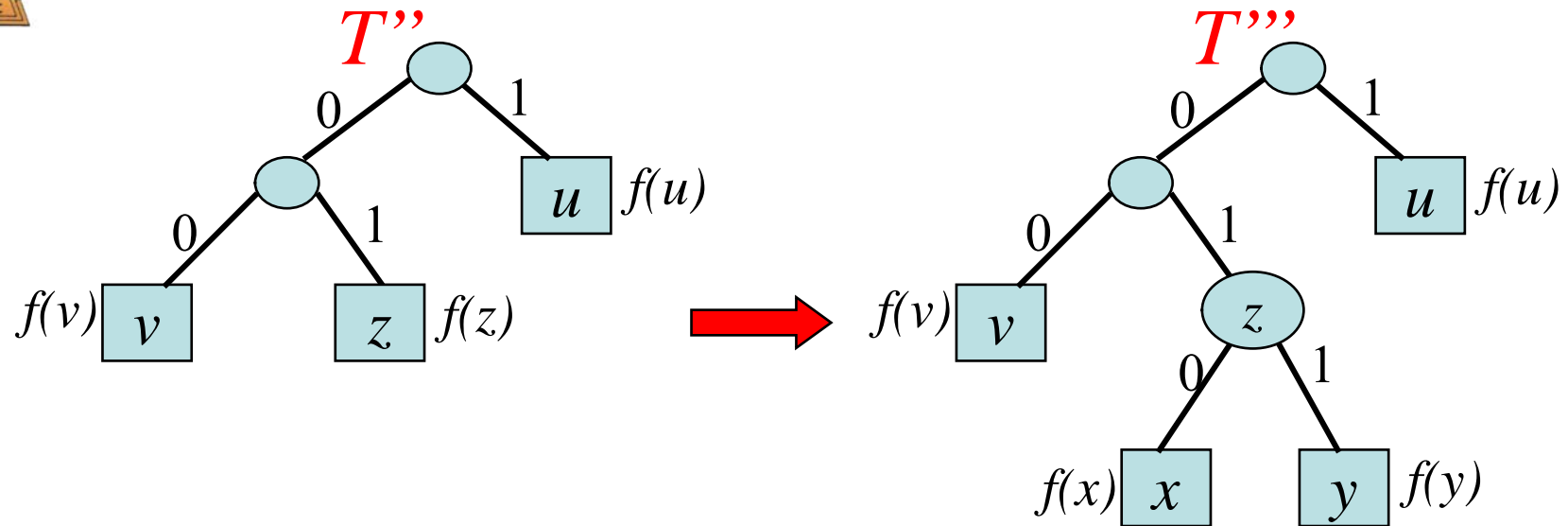
则必存在 T'' , 使 $B(T'') < B(T')$.

因为 z 是 C' 中字符, 它必为 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 的一个如下前缀编码树 T''' :





如上可证：

$$B(T''') = B(T'') + f(x) + f(y)。$$

由于 $B(T'') < B(T')$,

$$B(T''') = B(T'') + f(x) + f(y) < B(T') + f(x) + f(y) = B(T)$$

与 T 是优化的矛盾，故 T' 是 C' 的优化编码树。



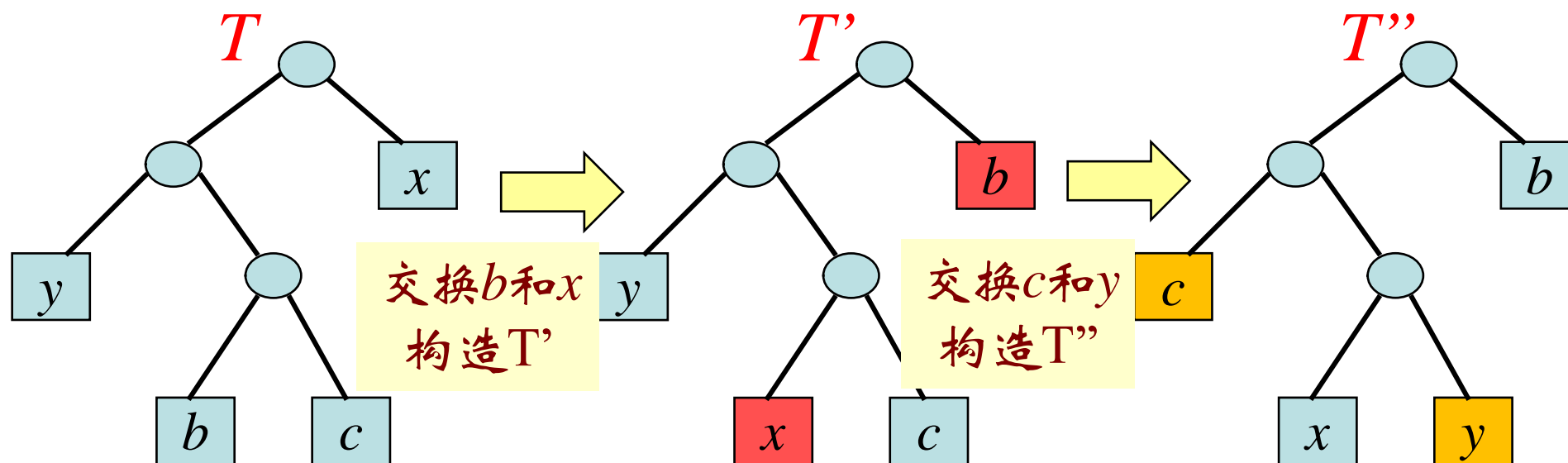
引理2. 设 C 是字母表, $\forall c \in C$, c 具有频数 $f(c)$, x 、 y 是 C 中具有最小频数的两个字符, 则存在一个 C 的优化前缀树, x 与 y 的编码具有相同最大长度, 且仅在最末一位不同.

优化前缀树问题具有 Greedy 选择性.



证：若 T 是 C 的优化前缀树，如果 x 和 y 是具有最大深度的两个兄弟字符，则命题得证。

若不然，设 b 和 c 是具有最大深度的两个兄弟字符：



不失一般性，设 $f(b) \leq f(c)$, $f(x) \leq f(y)$.

因 x 与 y 是具有最低频数的字符， $f(x) \leq f(y) \leq f(b) \leq f(c)$.

交换 T 的 b 和 $\dots x$, 从 T 构造 T' ; 交换 T' 的 c 和 y , 从 T' 构造 T''



往证 T'' 是最优化前缀树.

$$B(T) - B(T')$$

$$= \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C} f(c)d_{T'}(c)$$

$$= f(x)d_T(x) + f(b)d_T(b) - f(x)d_{T'}(x) - f(b)d_{T'}(b)$$

$$= f(x)d_T(x) + f(b)d_T(b) - f(x)d_T(b) - f(b)d_T(x)$$

$$= (f(b) - f(x))(d_T(b) - d_T(x)).$$

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

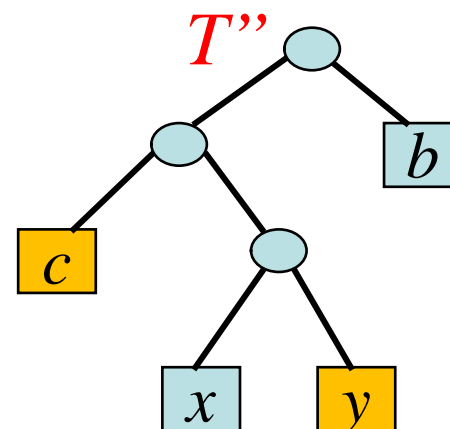
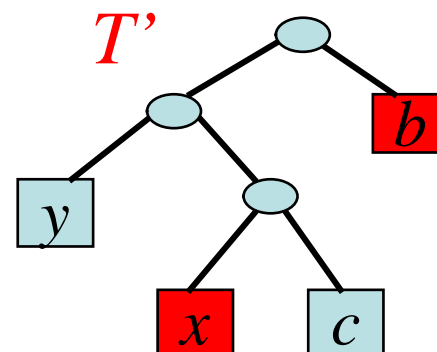
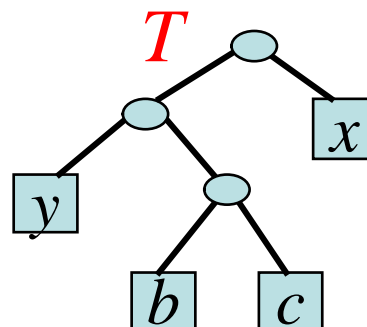
$\therefore B(T) - B(T') \geq 0, B(T) \geq B(T')$

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, $B(T) = B(T'')$, T'' 是 C 的最优化前缀编码树.

在 T'' 中, x 和 y 具有相同最大长度编码, 且仅最后位不同.





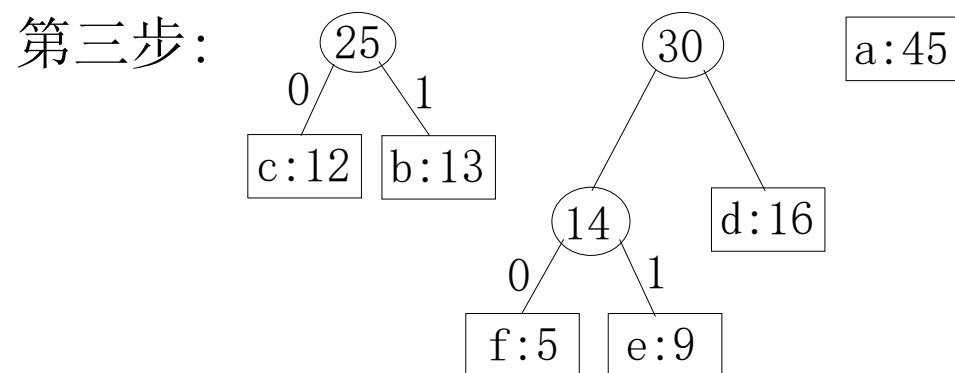
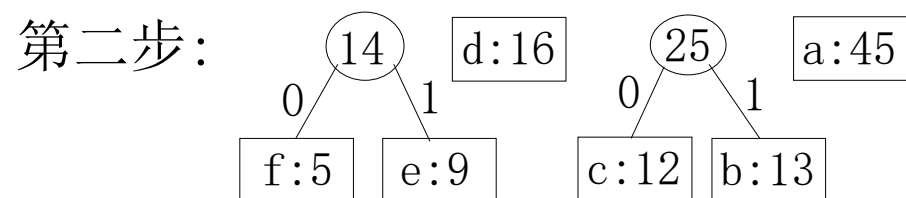
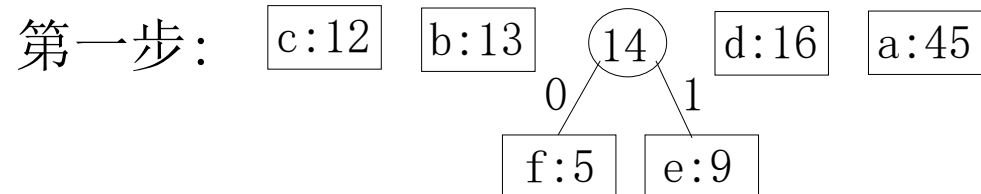
- 基本思想

- 循环地选择具有最低频数的两个结点，生成一棵子树，直至形成树
- 例子: $f:5, e:9, c:12, b:13, d:16, a:45$

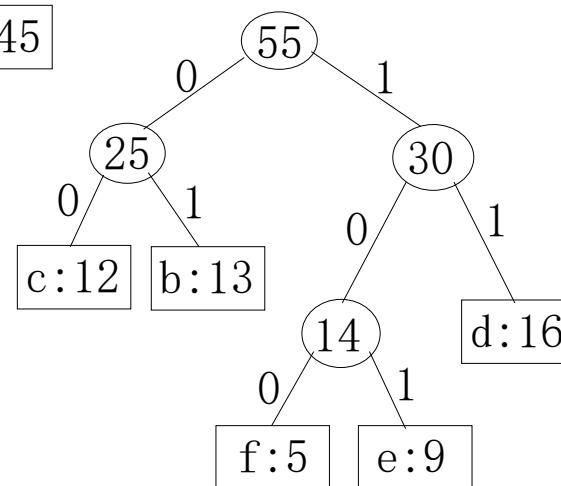


HIT

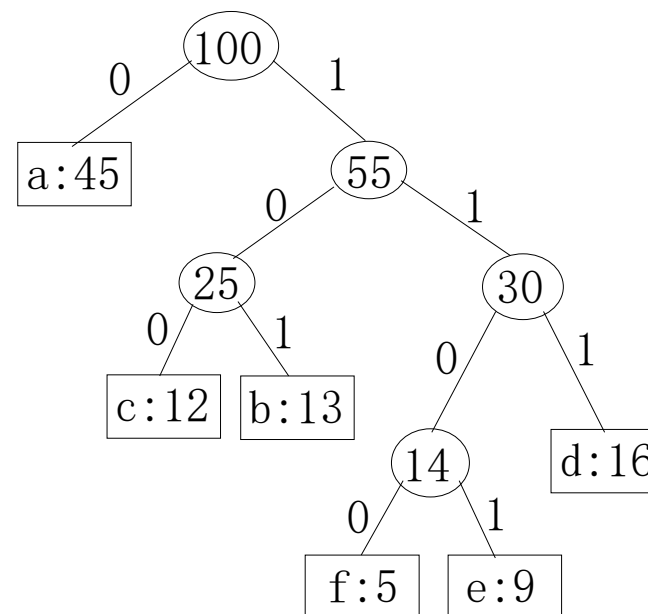
f:5 e:9 c:12 b:13 d:16 a:45



第四步: a:45



第五步:





- Greedy算法 (Q 是min-heap)

Huffman(C, F)

1. $n \leftarrow |C|$; $Q \leftarrow$ 根据 F 排序 C ; T 为一个空树;
2. FOR $i \leftarrow 1$ To $n-1$ DO
3. $z \leftarrow \text{Allocate-Node}()$;
4. $\text{left}[z] \leftarrow x \leftarrow \text{Extract-min}(Q)$ /* 从 Q 删除 x */;
5. $\text{right}[z] \leftarrow y \leftarrow \text{Extract-min}(Q)$ /* 从 Q 删除 y */;
6. $f(z) \leftarrow f(x) + f(y)$;
7. $\text{Insert}(Q, z, f(z))$;
8. Return $\text{Extract-min}(Q)$ /* 返回树的根 */



Huffman(C, F)

1. $n \leftarrow |C|$; $Q \leftarrow C$ (为 C 按照字符频数建立堆); T 为一个空树;
2. FOR $i \leftarrow 1$ To $n-1$ Do
3. $z \leftarrow \text{Allocate-Node}()$;
4. $\text{left}[z] \leftarrow x \leftarrow \text{Extract-min}(Q)$ /*从 Q 删除 x */;
5. $\text{right}[z] \leftarrow y \leftarrow \text{Extract-min}(Q)$ /*从 Q 删除 y */;
6. $f(z) \leftarrow f(x) + f(y)$;
7. $\text{Insert}(Q, z, f(z))$;
8. Return $\text{Extract-min}(Q)$ /* 返回树的根 */

- 第1步: 建堆 $O(n)$
- 每次循环: $O(\log n)$, 循环 $n-1$ 次: $O(n \log n)$
- $T(n) = O(n \log n)$



定理3. Huffman算法产生一个优化前缀编码树

证. 由于引理1、引理2成立，

且Huffman算法按照引理2的Greedy选择性确定的规则进行局部优化选择，所以Huffman算法产生一个优化前缀编码树。



5.4 最小生成树问题

- 问题定义
- 问题求解
 - 设计贪心选择法
 - 优化解结构分析
 - Greedy选择性证明
 - Kruskal算法
 - 算法复杂性



- 生成树

- 设 $G=(V, E)$ 是一个边加权无向连通图. G 的生成树是无向树 $T=(V, E_T)$, $E_T \subseteq E$.

- 生成树的权

- 如果 $W: E \rightarrow \{\text{实数}\}$ 是 G 的权函数, T 的权值定义为 $W(T) = \sum_{(u,v) \in T} W(u,v)$.

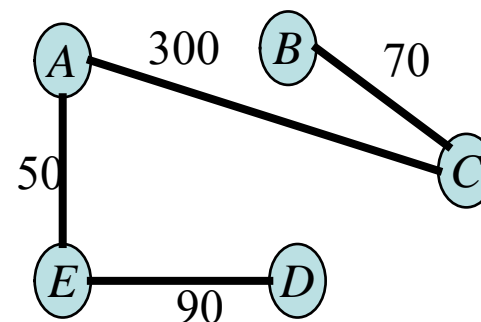
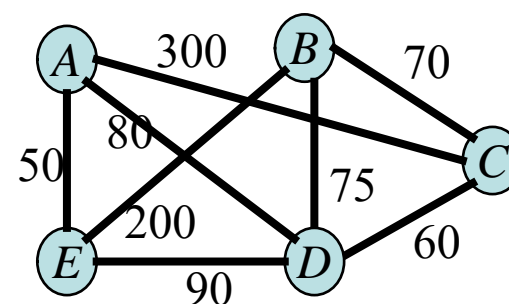
- 最小生成树

- G 的最小生成树是 $W(T)$ 最小的 G 之生成树.

- 问题的定义

- 输入: 无向连通图 $G=(V, E)$, 权函数 W

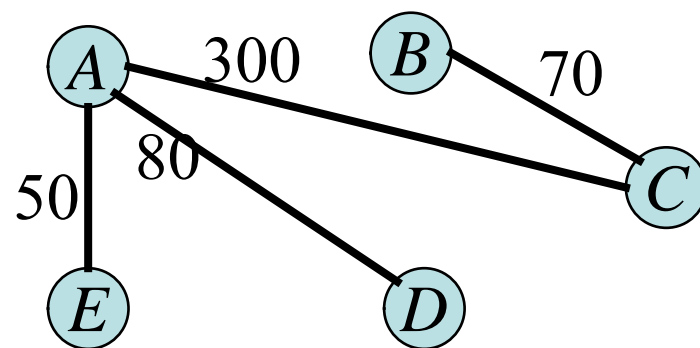
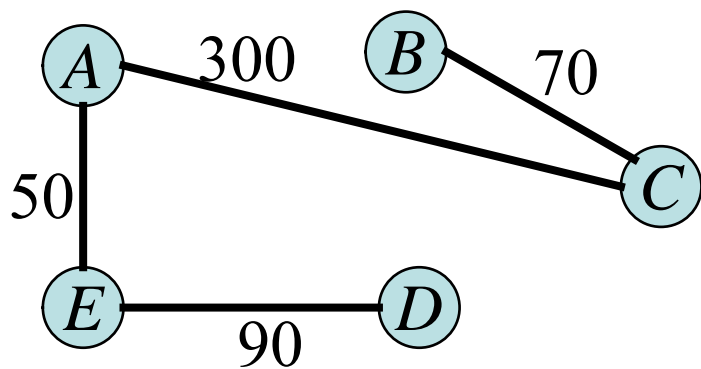
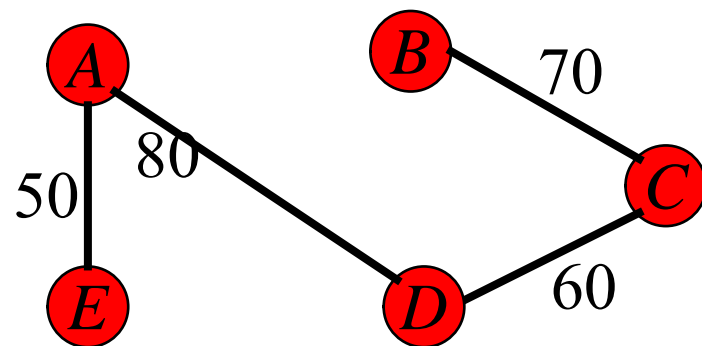
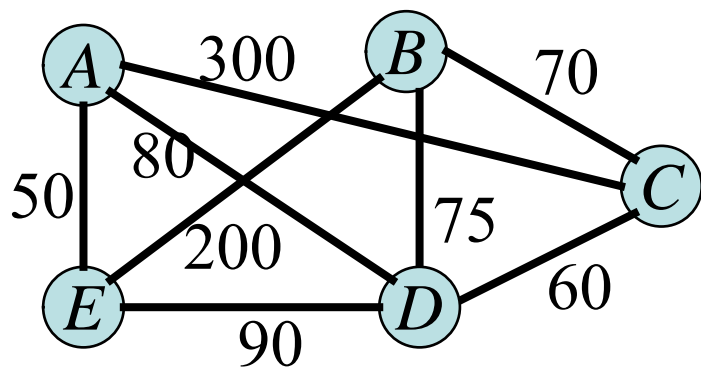
- 输出: G 的最小生成树





HIT

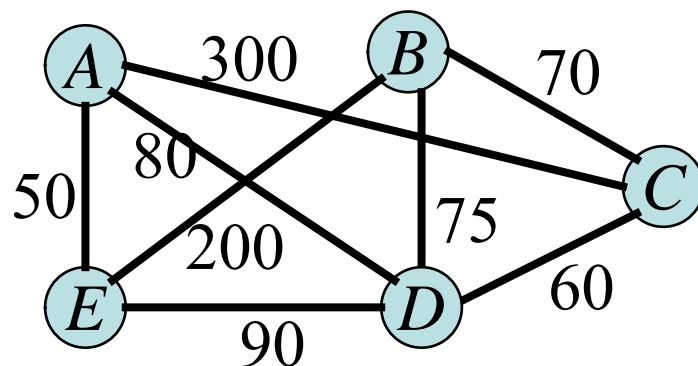
• 实例



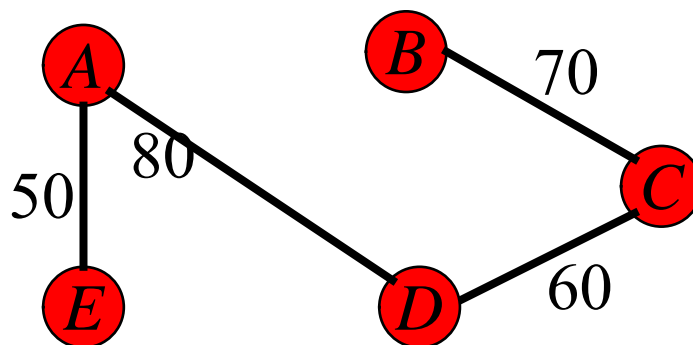


设计贪心选择方法

• 基本思想



- 初始: $A = \text{空}$; 构造森林 $G_A = (V, A)$;



设 A 是一个最小生成树的边子集合, 如果 $A \cup \{(u, v)\}$ 也是一个最小生成树的边子集合, 则 (u, v) 称为对 A 是安全的.

- 循环: 选择连接 G_A 中两棵树的最小安全边加入 A , 直至 G_A 是一棵生成树.



- 一般算法的定义

Generic-MST(G, W)

1. $A = \Phi$;

2. While A 不是生成树 Do

3. 寻找一个最小安全边 (u, v) ;

2. $A = A \cup \{ (u, v) \}$;

3. Return A

算法的关键!

设 A 是一个最小生成树的边子集合, 如果 $A \cup \{ (u, v) \}$ 仍是一个最小生成树的边子集合, 则 (u, v) 称为对 A 是安全的.

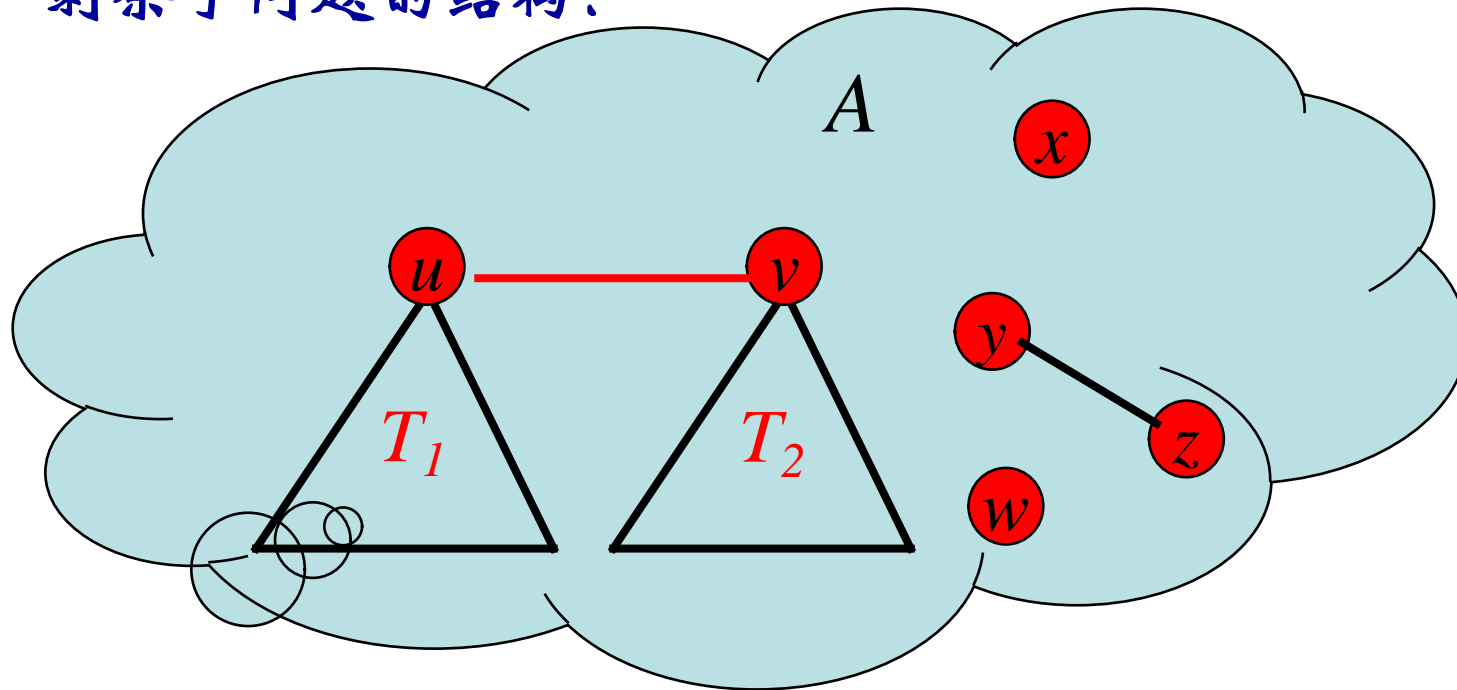


- 贪心选择方法

- 贪心选择:

- 从森林中选择一条最小安全边连接两棵子树为一棵子树，直至森林成为一棵树

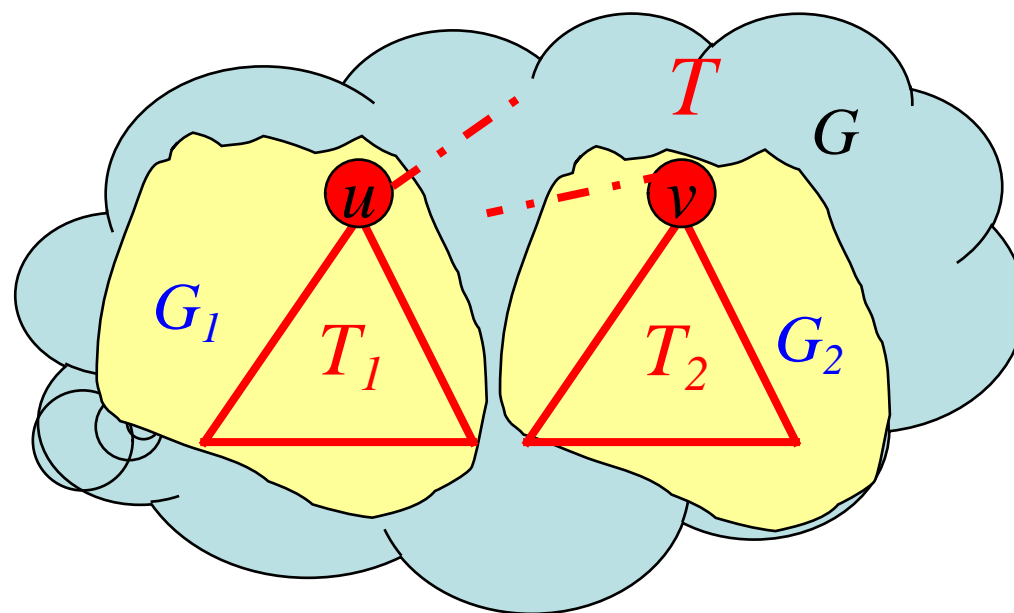
- 剩余子问题的结构:





定理1. 设 T 是 G 的最小生成树. 如果 T 包含子树 T_1 和 T_2 , T_1 是 G 的子连通图 G_1 的生成树, T_2 是 G 的子连通图 G_2 的生成树, 则 T_1 是 G_1 的最小生成树, T_2 是 G_2 的最小生成树.

证. (略)



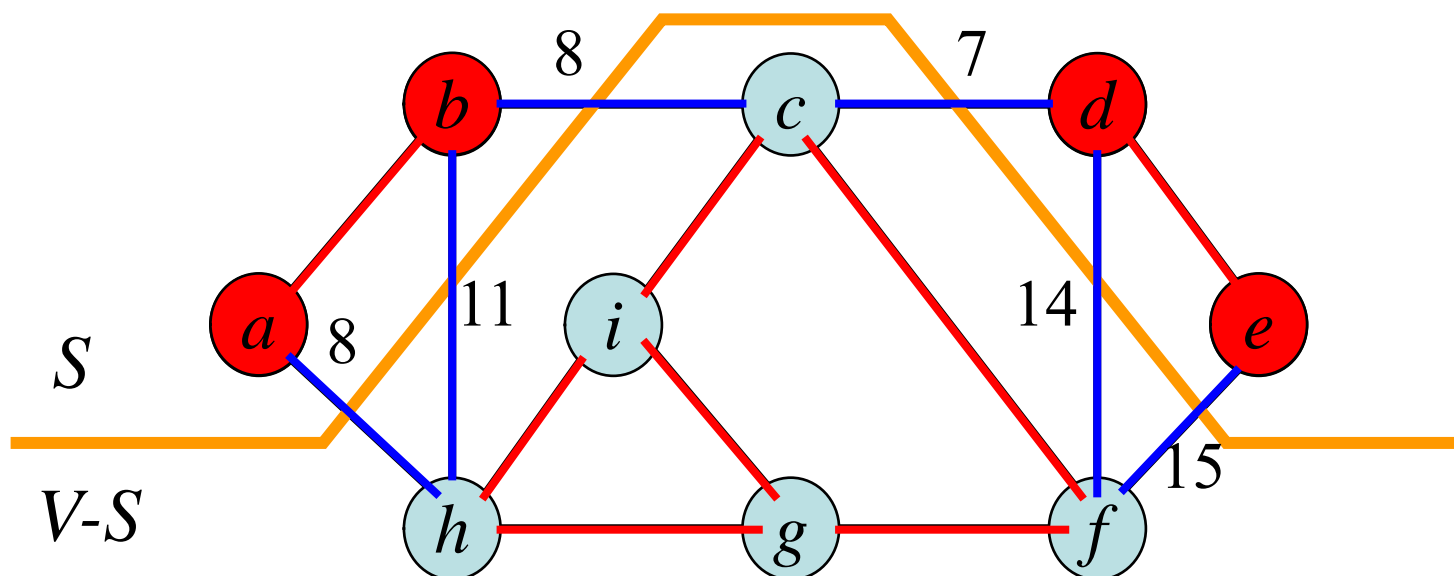


定义1. 无向图 $G=(V, E)$ 的一个划分是 V 的划分 $(S, V-S)$.

定义2. 若 $u \in S, v \in V-S$, 则 (u, v) 称为划分 $(S, V-S)$ 的交叉边.

定义3. 如果边集合 A 中没有边是划分 $(S, V-S)$ 的交叉边, 则称划分 $(S, V-S)$ 尊重 A .

定义4. 划分 $(S, V-S)$ 的交叉边 (u, v) 称为轻边, 如果在所有 $(S, V-S)$ 的交叉边中, (u, v) 的权值最小.





定理2. 设 $G=(V, E)$ 是具有边加权函数 W 的无向连通图, $A \subseteq E$ 是包含在 G 的某最小生成树中的边集合, $(S, V-S)$ 是 G 的尊重 A 的任意划分, (u, v) 是 $(S, V-S)$ 的交叉轻边, 则 (u, v) 对 A 是安全的.

证.

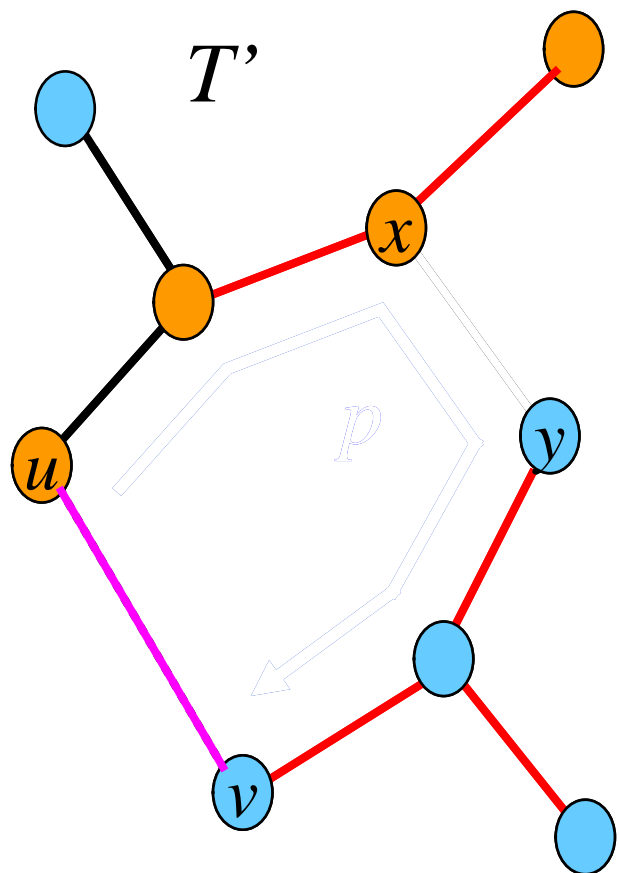
令 T 是包含 A 的最小生成树.

若 (u, v) 属于 T , 则 (u, v) 对 A 是安全的.

若 (u, v) 不属于 T .

我们构造一个 G 的最小生成树 T' , 使其包含 $A \cup \{(u, v)\}$, 从而证明 (u, v) 安全.

设 A 是一个最小生成树的边子集合, 如果 $A \cup \{(u, v)\}$ 也是一个最小生成树的边子集合, 则 (u, v) 称为**对 A 是安全的**.



- S = 黄 结点集合
- $V-S$ = 蓝 结点集合
- A = 红 边集合

由于 u 和 v 在划分 $(S, V-S)$ 的两边, T 至少存在一条交叉边在从 u 到 v 的路径 p 中, 设该交叉边为 (x, y) .

由于划分尊重 A , 故 (x, y) 不在 A 中. 删除 p 中的 (x, y) , 增加 (u, v) , 得到

$$T' = T - \{(x, y)\} \cup \{(u, v)\}.$$

往证 T' 是最小生成树.

因为 (u, v) 是交叉轻边, (x, y) 是交叉边, $W(u, v) \leq W(x, y)$.

$$W(T') = W(T) - W(x, y) + W(u, v) \leq W(T)$$

由于 T 是最小生成树, $W(T') = W(T)$.

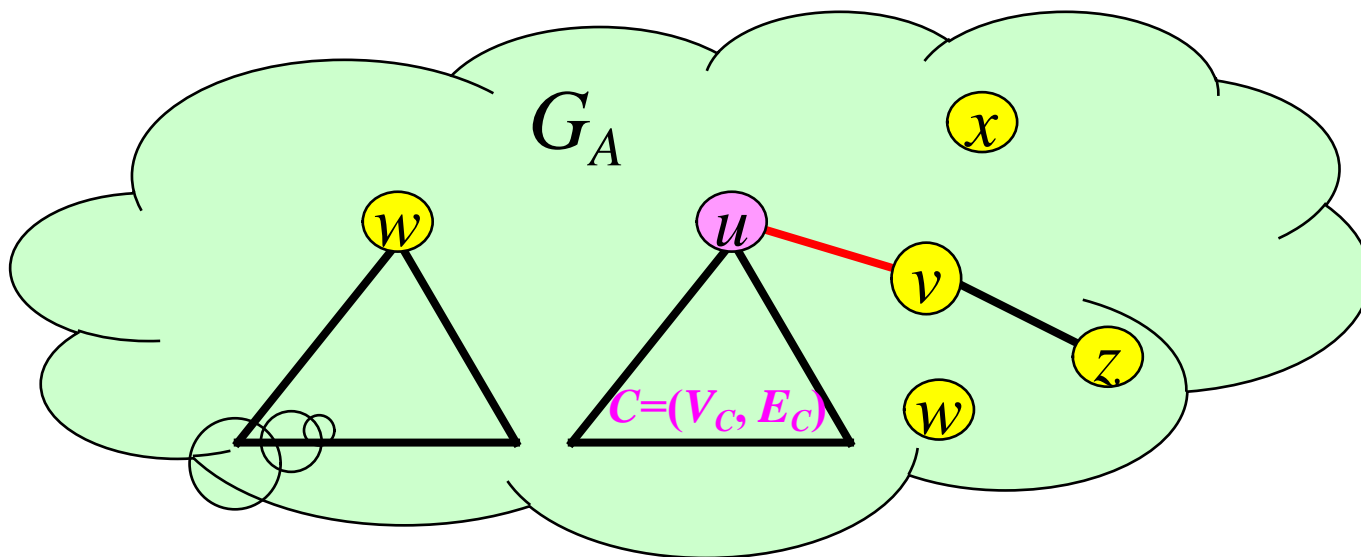
T' 是最小生成树, $A \cup \{(u, v)\} \subseteq T'$,

(u, v) 对于 A 是安全的.



推论1. 设 $G=(V, E)$ 是具有边加权函数 W 的无向连通图, $A \subseteq E$ 是包含在 G 的某个最小生成树中的边集合, $C=(V_C, E_C)$ 是森林 $G_A=(V, A)$ 中的树. 如果 (u, v) 是连接 C 和 G_A 中另一个树的交叉轻边, 则 (u, v) 对 A 是安全的.

证. 划分 $(V_C, V-V_C)$ 尊重 A , 因为 A 的边要么在 C 中, 要么在 $G_A=(V, A)$ 的另一个树中. (u, v) 是关于这个划分的交叉轻边. 于是, (u, v) 对 A 是安全的.





- 基本思想

找到连接森林 $G(V, A)$ 中两棵树的安全交叉轻边加入 A ，直至 $G(V, A)$ 成为一棵树。



MST-Kruskal(G, W)

1. $A = \Phi$;
2. For $\forall v \in V[G]$ Do
3. Make-Set(v); /* 建立只有 v 的集合 */
4. 按照 W 值的递增顺序排序 $E[G]$;
5. For $\forall (u, v) \in E[G]$ (按 W 值的递增顺序) Do
6. If Find-Set(u) \neq Find-Set(v)
7. Then $A = A \cup \{(u, v)\}$; Union(u, v);
8. Return A



MST-Kruskal(G, W)

1. $A = \Phi$;
2. For $\forall v \in V[G]$ Do
3. Make-Set(v); /* 建立只有 v 的集合 */
4. 按照 W 值的递增顺序排序 $E[G]$;
5. For $\forall (u, v) \in E[G]$ (按 W 值的递增顺序) Do
6. If Find-Set(u) \neq Find-Set(v)
7. Then $A = A \cup \{(u, v)\}$;
 Union(u, v);
8. Return A

- 令 $n = |V|$, $m = |E|$
- 第4步需要时间: $O(m \log m)$
- 第2-3步执行 $O(n)$ 个 Make-Set 操作
- 第5-7步执行 $O(m)$ 个 Find-Set 和 Union 操作
需要时间: $O(m \alpha(n))$
- $m \geq n - 1$ (因为 G 连通), 由 $\alpha(n) < \log n < \log m$
- 总时间复杂性: $O(m \log m)$

集合操作的复杂性见Intro. To Algo. 第21章 (498-509)



定理2. $\text{MST-Kruskal}(G, W)$ 算法能够产生图 G 的最小生成树.

证. 因为算法按照 Greedy 选择性进行局部优化选择, 并且每次选择的都是最小边.



5.5 贪心算法的基础理论

- Matroid (拟阵)
- Matroid 实例
- Matroid的性质
- 加权Matroid上的Greedy算法
- 算法的正确性



- Matroid定义

Matroid是一个序对 $M=(S, I)$ ，满足：

- (1) S 是一个有限非空集合.
- (2) $I \subseteq 2^S$, I 非空, I 中的子集称为 S 的独立子集.
- (3) 遗传性: 如果 $A \in I$, $B \subseteq A$, 则 $B \in I$
- (4) 交换性: 如果 $A \in I$, $B \in I$, $|A| < |B|$, 则
 $\exists x \in B - A$ 使得 $A \cup \{x\} \in I$.

遗传性和交换性是拟阵最根本的2条性质，
拟阵上的其他性质都是基于这2个性质发展出来的



- 实例(Graphic Matroid)

定义1. 设 $G = (V, E)$ 是一个无向连通图, $|E| \geq 1$. 由 G 确定的

$M_G = (S_G, I_G)$ 定义如下: S_G 是 G 的边集合 E ,

$I_G = \{A | A \subseteq E, G_A(V, A) \text{ 是森林}\}.$

定理1. 如果 G 是一个无向连通图, 则 $M_G = (S_G, I_G)$ 是一个拟阵.

证. (1) S_G 非空有限性: $S_G = E$ 是一个非空有限集合.

(2) I_G 的非空性: $\forall e \in E, G_e(V, \{e\})$ 是一个森林, $\{e\} \in I_G$,
于是, I_G 是 S_G 的非空集族.

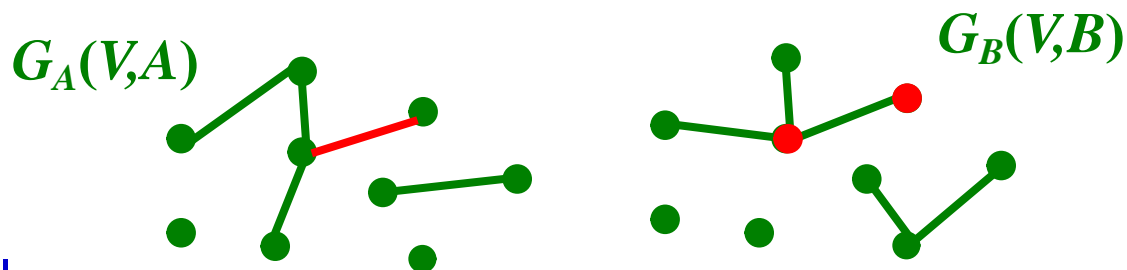
(3) M_G 满足遗传性:

如果 $A \in I_G, B \subseteq A$, 则 $G_B(V, B)$ 是一个森林.

于是, $B \in I_G, M_G$ 满足遗传性.



定理1. 如果 G 是一个无向连通图, 则 $M_G=(S_G, I_G)$ 是一个拟阵.
证(续).



(4) M_G 满足交换性:

设 $A \in I_G$, $B \in I_G$, $|A| < |B|$.

图论定理: 具有 k 条边的森林具有 $|V| - k$ 棵树.

$G_A=(V, A)$ 和 $G_B=(V, B)$ 分别具有 $|V|-|A|$ 和 $|V|-|B|$ 棵树.

由于 $|V|-|A| > |V|-|B|$, G_B 中必存在树 T , T 的节点在 G_A 的不同树中. 否则 G_A 中至少有一棵树的节点不在 G_B 中, 与 G_B 的节点集合为 V 矛盾.

由于 T 是连通的, T 必包含边 (u, v) : u, v 在 G_A 的不同树中.

于是, $(V, A \cup \{(u, v)\})$ 是森林, $A \cup \{(u, v)\} \in I_G$, $(u, v) \in B - A$, M_G 满足交换性.



• Matroid的性质

定义2. 设 $M=(S, I)$ 是一个Matroid, $A \in I$. 如果 $A \cup \{x\} \in I$, $x \notin A$, x 称为 A 的一个扩展(extension).

定义3. 设 $M=(S, I)$ 是Matroid, $A \in I$. 若 A 没有extension, 则称 A 为最大独立子集合.

定理2. 一个Matroid的所有最大独立子集合都具有相同大小.

证. 设 A 是Matroid M 的最大独立子集合, 而且存在 M 的另一个独立子集合 B , $|A| < |B|$.

根据 M 的交换性, $\exists x \in B - A$ 使 $A \cup \{x\} \in I$, 与 A 最大矛盾.



- Matroid的最优独立子集

实际背景:

很多可用Greedy算法获得最优解的问题可以归结为在加权Matroid中寻找最优独立子集问题, 即给定 $M=(S, I)$ 和权函数 W , 找到独立子集 $A \in I$, 使得 $W(A)$ 最大。



- 实例1: 最大生成森林问题

- 问题定义

输入: 无向图 $G=(V, E)$, $|E|>0$, 权函数 $W: E \rightarrow$ 正整数集

输出: 边子集合 $A \subseteq E$, $G_A(V, A)$ 是森林, $W(A)$ 最大

- 转换为加权Matroid上寻找最优独立子集问题

- 构造:

$M_G=(S_G, I_G)$, $S_G = E$, $I_G = \{A | A \subseteq E, G_A(V, A) \text{ 是森林}\}$, M_G 是拟阵

- M_G 的最优独立子集 A 满足:

- (V, A) 是森林

- $W(A)$ 最大



HIT
CS&E

- 加权Matroid最优独立子集问题

输入: Matroid $M=(S, I)$, 函数 $W: S \rightarrow$ 正数集

输出: M 的独立子集 $A \in I$, 使得 $W(A)$ 最大



- 贪心选择方法

从空集合 A 开始，每次选择权值最大的 $x \in S$ ，
扩展 A ，使 $A \cup \{x\} \in I$

引理1. 设 $M=(S, I)$ 是一个Matroid. 如果 $\{x\} \notin I$, 则 x 不是任何独立子集的元素.

证. 设 x 是独立子集 A 的元素, 即 $x \in A, A \in I$. 由 M 的遗传性, $\{x\} \in I$, 矛盾.

推论1. 任何元素一旦不能被初始选中, 则永远不会被选中.



贪心选择方法

从空集合 A 开始，每次选择权值最大的 $x \in S$ ，扩展 A ，使 $A \cup \{x\} \in I$

第一次选中 x 之后，剩余子问题：

$$S' = S - \{y \mid W(y) > W(x)\} - \{x\}$$

$$I' = \{B \subseteq S' \mid B \cup \{x\} \in I\}$$

M' 的权函数与 M 的权函数相同。

(如果元素 y 一旦没被选中， y 不会属于任何最优独立子集)



定理1. 设 $M=(S, I)$ 是一个加权 Matroid, W 是 M 的权函数, S 按 W 值递减排序. 设 x 是 S 中第一个满足 $\{x\} \in I$ 的元素, 则存在一个 M 的最优独立子集 A , $x \in A$.

证. 若存在最优独立子集 A 包含 x , 则定理得证.

否则, 设 B 是任意非空最优独立子集, $x \notin B$.

S 中 x 之前的元素 y 必不在 B 中, 否则由遗传性, $\{y\} \in I$, 与 “ x 是 S 中第一个满足 $\{x\} \in I$ 的元素” 矛盾.

显然, $\forall z \in B, W(x) \geq W(z)$. 如下构造含 x 的优化子集 A :

初始: $A = \{x\} \in I$;

用交换性: $\forall z \in B - A$, 若 $A \cup \{z\} \in I$, $A = A \cup \{z\}$, 直至 $|A| = |B|$.

显然, $\exists w \in B, A = (B - \{w\}) \cup \{x\}$.

于是, $W(A) = W(B) - W(w) + W(x) \geq W(B)$.

因为 B 是优化子集, 所以 $W(A) \leq W(B)$, $W(A) = W(B)$.

A 是优化子集, 且 $x \in A$.



定理2. 设 x 是第一个被Greedy算法选中的元素, A 是包含 x 的最优独立子集, $A'=A-\{x\}$ 是子问题 $M'=(S', I')$ 的最优独立子集.

$M'=(S', I')$ 定义如下:

$$S'=S-\{y \mid W(y) > W(x)\}-\{x\}$$

$$I'=\{B \subseteq S' \mid B \cup \{x\} \in I\}$$

M' 的权函数与 M 的权函数相同.

证.

因为 $A=A' \cup \{x\} \in I$, 所以 $A'=A-\{x\} \in I'$. (根据 I' 的定义)

若 A' 不是 M' 的最优独立子集, 则存在 M' 的一个最优独立子集 B 使得: $W(B) > W(A')$.

由于 $B \cup \{x\} \in I$, $W(A) = W(A') + W(x)$,

$$W(B \cup \{x\}) = W(B) + W(x) > W(A') + W(x) = W(A),$$

与 A 最优矛盾.



Greedy(M, W)

- 1 $A = \Phi$;
- 2 按权 W 值递减排序 S ;
- 3 For $\forall x \in S$ (按 $W(x)$ 递减顺序) DO
- 4 If $A \cup \{x\} \in I$ /* 选择目前 $W(x)$ 最大的 x */
- 5 Then $A = A \cup \{x\}$;
- 6 Return A .

时间复杂性

step 2: $O(|S| \log |S|)$

step 4: $O(f(|S|))$

$T(|S|) = O(|S| \log |S| + |S| f(|S|))$



引理2. Greedy 算法返回一个独立子集合.

证. 设 Greedy 返回集合 A , 对 $|A|$ 做数学归纳法.

当 $|A|=0$ 时, A 是空集, 由遗传性, A 是独立子集合.

设 $|A| \leq k$ 时, A 是独立子集.

当 $|A|=k+1$ 时, A 由第4-5步得到, 即 $A=A \cup \{x\}$.

第4步已判定 $A \cup \{x\} \in I$, $A=A \cup \{x\}$ 是独立子集.



定理3. 设 $M=(S, I)$ 是一个 Matroid, W 是 M 的权函数,
 $\text{Greedy}(M, W)$ 返回一个 M 的最优独立子集.

证. ①. 引理1说明, 任何没有被 Greedy 选中的 S 元素, 以后不会被选中, 可以不再考虑.

②. 算法每次循环都按照定理1给出方法进行贪心选择最大元素 x .

③. 算法每次循环都按照定理2的优化子结构求解子问题 M' 的最优独立子集的问题.

于是, $\text{Greedy}(M, W)$ 返回一个 M 的最优独立子集.



HIT
CS&E

5.6 任务调度问题

展示如何把一个可用Greedy方法求解的问题转换为求解拟阵最优子集问题



- 单位时间任务
需要一个单位时间就能够完成的任务
- 单位时间任务的调度问题

输入:

单位时间任务集 $S = \{1, 2, \dots, n\}$

正整数任务期限 $D = \{d_1, d_2, \dots, d_n\}$, 任务 i 须在 d_i 前完成

非负权集 $W = \{w_1, w_2, \dots, w_n\}$, 任务 i 不在 d_i 前完成罚款 w_i

输出:

S 的一个调度 (S 的一个执行顺序), 具有最小总罚款数



- 转换为加权Matroid的最优子集问题

定义1. 设 S 是一个任务调度.任务 i 在 S 中是迟任务如果它在 d_i 之后完成; 否则是早任务.

定义2. 如果在一个调度中,早任务总是先于迟任务,则称该调度具有早任务优先形式.

定义3. 如果一个调度具有早任务优先形式而且按期限单调递增顺序执行各个早任务,则称该调度具有规范化形式.

例如: 任务集合 $\{1,2,3,4\}$, 期限集合 $\{d_1=1, d_2=1, d_3=4, d_4=3\}$, $W=\{4,6,3,8\}$

可能的任务调度: 1,2,4,3 或 2,1,4,3 或 1,3,4,2 或 1,4,3,2.....

其中, 1,3,4,2 及 1,4,3,2 是早任务优先调度。

调度1,4,3,2具有规范化形式。



S不是独立集!

例如: 任务集合 $S=\{1,2,3,4\}$, 期限集合 $\{d_1=1, d_2=1, d_3=4, d_4=3\}$
 $N_1(S)=2, N_2(S)=2, N_3(S)=3, N_4(S)=4$.

若 $A=\{1,3,4\}$, 则 $N_1(A)=1, N_2(A)=1, N_3(A)=2, N_4(A)=3$;

定义4. 任务集合 A 称为独立集如果存在一个关于 A 的调度, 使得 A 中的任务皆为早任务.

例. 一个优化调度的早任务集合是独立集.

以下:

用 I 表示所有独立任务集合的集族

用 $N_t(A)$ 表示任务集合 A 中期限小于等于 t 的任务数



任务集合 $\{1,2,3,4\}$, 期限集合 $\{d_1=1, d_2=1, d_3=4, d_4=3\}$

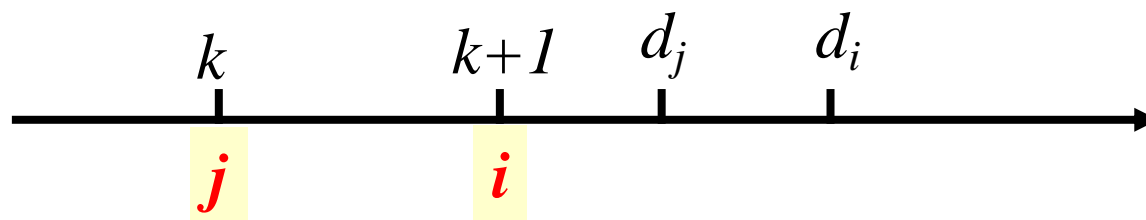
$1,2,4,3 \Rightarrow 1,4,3,2$ $1,3,4,2 \Rightarrow 1,4,3,2$

• 任务调度的规范化

第一步: 将调度安排成早任务优先形式, 即如果早任务 x 跟在迟任务 y 之后, 交换 x 和 y 的位置;

*** 早任务优先形式不改变任何任务的早或迟状态**

第二步: 如果任务 i 和 j 是早任务, 而且分别完成于时刻 k 和 $k+1$, $d_j < d_i$, 交换 i 和 j 的位置.



寻找最优调度问题成为寻找最优早任务集合 A 的问题. 一旦 A 被确定后, 就可以按期限单调递增序列出 A 中的所有元素, 然后按任意序列出迟任务(即 $-A$)



$N_t(A)$ = A 中期限小于等于 t 的任务数

引理1. 对于任意任务集合 A , 下面的命题等价:

1. A 是独立集合,
2. 对于 $t=1, 2, \dots, n$, $N_t(A) \leq t$,
3. 如果按照期限递增顺序调度 A 中任务, 则 A 中无迟任务.

证. 1→2. 如果 $N_t(A) > t$, 则有多于 t 个任务需要在 t 时间内完成, 不存在使得 A 中无迟任务的调度.

2→3. 若 A 中任务依其期限递增排列, 则2意味着排序后, 在时间 t 之前必须完成的 A 中任务数至多为 t . 于是, 按期限递增顺序调度 A 中任务, A 无迟任务.

3→1. 显然.



定理1. 若 S 是一个带期限的单位时间任务的集合, 且 I 为所有独立任务集构成的集族, 则 $M=(S, I)$ 是一个Matroid.

证明. 1. S 的非空有限性: 显然.

2. I 的非空性:

因为单个任务集属于 I , 所以 I 非空.

3. 遗传性:

若 $A \in I, B \subseteq A$, 则 $B \in I$.



$N_t(A)$ = A 中期限小于等于 t 的任务数

$A = \{1^1, 3^2, 4^3, 10^7, 8^8\}$, $B = \{5^1, 2^3, 4^3, 6^5, 10^7, 9^{10}, 7^{11}\}$

i^d 表示任务 i 期限为 d , 选 B 中哪个任务加入 A ?

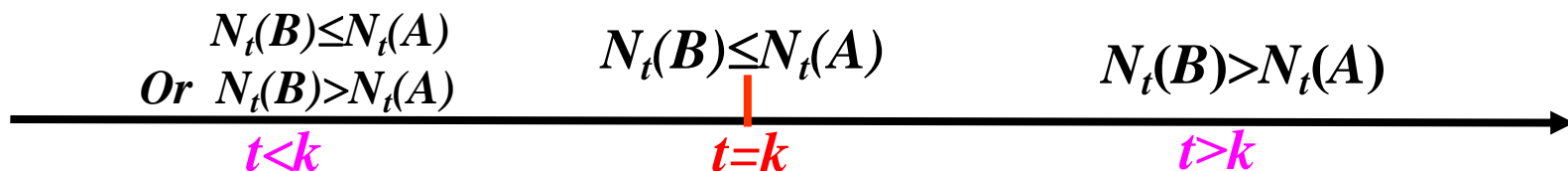
$A' = A \cup \{9^{10}\} = \{1^1, 3^2, 4^3, 10^7, 8^8, 9^{10}\}$

4. 交换性:

设 $A, B \in I$, $|B| > |A|$, 往证 $\exists x \in B - A$, $A \cup \{x\} \in I$.

由于 $|B| > |A|$, 不可能对于所有 t , $N_t(B) \leq N_t(A)$

令 $k = \max\{t \mid N_t(B) \leq N_t(A)\}$.



于是, B 中包含了比 A 中更多的具有期限 $k+1$ 的任务.

设 $x \in B - A$, x 具有期限 $k+1$.

令 $A' = A \cup \{x\}$. 往证 A' 是独立集.

对于 $1 \leq t \leq k$, $N_t(A') = N_t(A) \leq t$, 因为 A 是独立集.

对于 $k < t \leq n$, $N_t(A' = A \cup \{x\}) \leq N_t(B) \leq t$, 因为 B 是独立集.

于是, A' 是独立集.



- 为 $M=(S, I)$ 定义权值
 - 令 $w=\max\{w_1, w_2, \dots, w_n\}$
 - $\forall x_i \in S, W(x_i)=w_i.$
 - $\forall A \in I, W(A)=nw-\sum_{y \in S-A} W(y)$
- 任务调度问题转换为 $M=(S, I)$ 上寻找最大独立子集（或优化子集）问题
 - 若 A 是优化子集，则 $W(A)=nw-\sum_{y \in S-A} W(y)$ 最大
 - 即罚款 $\sum_{y \in S-A} W(y)$ 最小