

计算理论与 算法分析设计

教材:

[1][王] 王晓东, 计算机算法设计与分析(第4版), 电子工业.

[2][S] 唐常杰等译, Sipser著, 计算理论导引, 机械工业.

参考资料:

[3][C] 潘金贵等译, Cormen等著, 算法导论, 机械工业.

[4][M] 黄林鹏等译, Manber著, 算法引论-一种创造性方法, 电子.

[5][刘] 刘汝佳等, 算法艺术与信息学竞赛, 清华大学.

[6][L] Lewis等著, 计算理论基础, 清华大学.

第四章 贪心算法

Greedy

[王]本章内容

4.1 活动安排问题

4.2 贪心算法的基本要素

4.3 最优装载 见上一章课件

4.4 哈夫曼编码

4.5 单源最短路 见上一章课件

4.6 最小生成树

4.7 多机调度问题证明

第四章 贪心算法

Greedy

!!贪心不一定正确(0-1背包), 需要证明

1. 活动安排问题

2. 贪心算法基本要素

3. Huffman算法

4. 最小生成树

5. 多机调度问题

活动安排问题

n 个活动申请一个活动室, 各活动起始终止区间 (s_i, f_i)

- 输入: $n, (s_i, f_i), i = 1:n$

- 输出: 最大相容活动子集(无冲突, 活动个数最多)

举例: $\{(1,4), (3,5), (0,6), (5,7)\}$ //解不唯一

- 最优解有什么特点? 贪心, OSP($A[i]=\text{true}$ 或 false)

- 最优解可以包含 最早结束的 或 最晚开始的 //举例

先给出算法, 再证明正确性

1. 按终止时间排序 $f_1 \leq f_2 \leq \dots \leq f_n$. $A[1] = \text{true}; \text{pt} = 1;$

2. 对 $i = 2:n$

3. 若 $s[i] \geq f[\text{pt}]$, 则 $A[i] = \text{true}, \text{pt} = i,$

4. 否则 $A[i] = \text{false}$

活动安排算法正确性证明

n 个活动申请一个活动室, 活动起始终止区间 (s_i, f_i)

- 输入: $n, (s_i, f_i), i = 1:n$, 输出: 最大相容活动子集

- 第一步 证明贪心选择性质:

存在最优解包含相容最早结束的活动1

- 第二步 证明最优子结构性质:

若A是包含活动1的最大相容活动集(最优策略),

则A-{1}是区间 $[f_1, f_n]$ 上的最大相容活动集

- 由此用数学归纳法, 就可以证明算法正确性

第四章 贪心算法

Greedy

!!贪心不一定正确(0-1背包), 需要证明

1. 活动安排问题

2. 贪心算法基本要素

3. Huffman算法

4. 最小生成树

5. 多机调度问题

贪心算法的基本要素

对于一个具体问题, 怎么知道能否用贪心算法

- 贪心选择性质和最优子结构性质(比较矩阵连乘)
- 贪心选择性质

对比: 矩阵连乘, 0-1背包, 分数背包

贪心算法第一基本要素, 与DP主要区别

自顶向下计算

- **OSP**: 最优策略的子策略也是最优 //动规, 贪心
- 正确性证明一般过程:

贪心选择+OSP+数学归纳法

条件: 子问题与原问题类似, 相对独立

子问题的最优解和贪心选择联合得整体最优解

第四章 贪心算法

Greedy

!!贪心不一定正确(0-1背包), 需要证明

1. 活动安排问题
2. 贪心算法基本要素
3. Huffman算法
4. 最小生成树
5. 多机调度问题

哈夫曼编码

- **Huffman**编码广泛应用于数据文件压缩
- 压缩率通常在**20%~90%**之间
- 对文章先统计各字母出现频率, 再选择最优编码
- 输入: 字母集和各字母出现频率
- 输出: 各字母的最优编码(压缩率最大)
- 前缀码: 任何代码都不是其它代码的前缀

例 a:0, b:101, c:100, d:11

可以唯一解码011101100 = adbc

定长码是前缀码: a:00, b:01, c:10, d:11

不同编码举例

	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

使用定长码编码的存储空间:

$$3 \times (45 + 13 + 12 + 16 + 9 + 5) = 300\text{kb}$$

使用变长码编码的存储空间:

$$1 \times 45 + 3 \times (13 + 12 + 16) + 4 \times (9 + 5) = 224\text{kb}$$

节约空间约25%. 构造编码的原则?

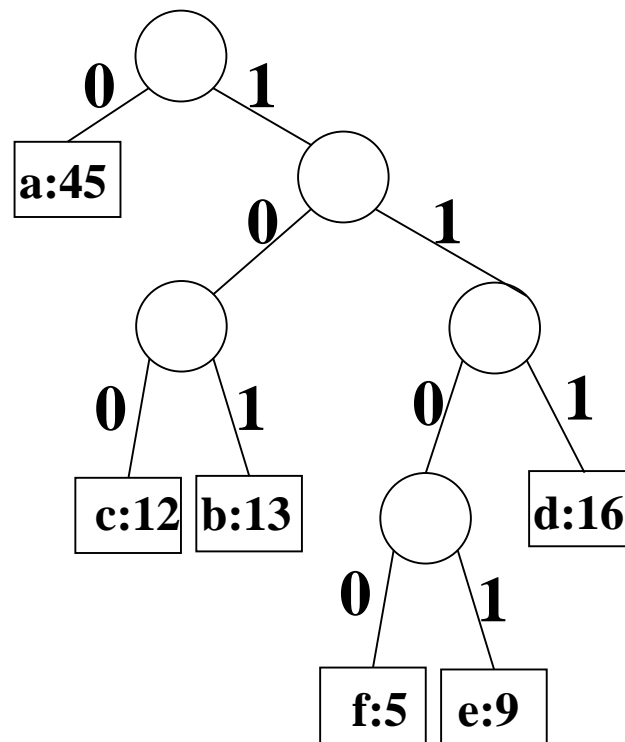
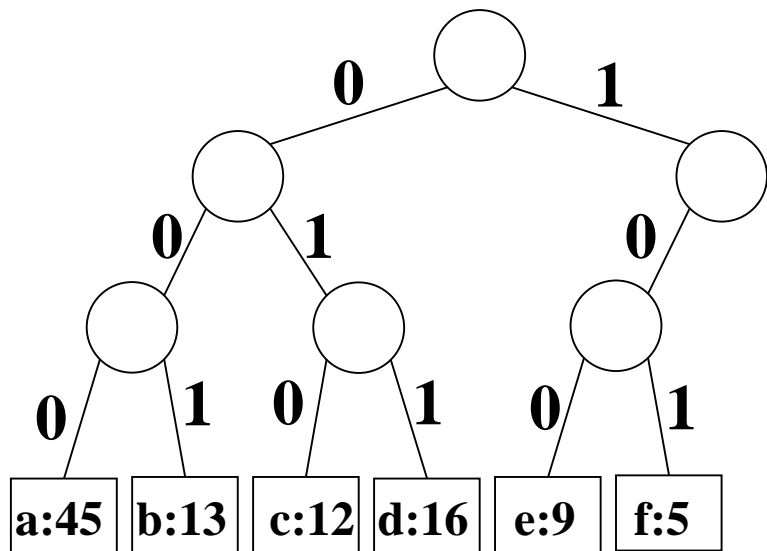
Huffman编码是最优前缀码

编码的数据结构

	a	b	c	d	e	f	存储
频率(千次)	45	13	12	16	9	5	100
定长码	000	001	010	011	100	101	300kb
变长码	0	101	100	111	1101	1100	224kb

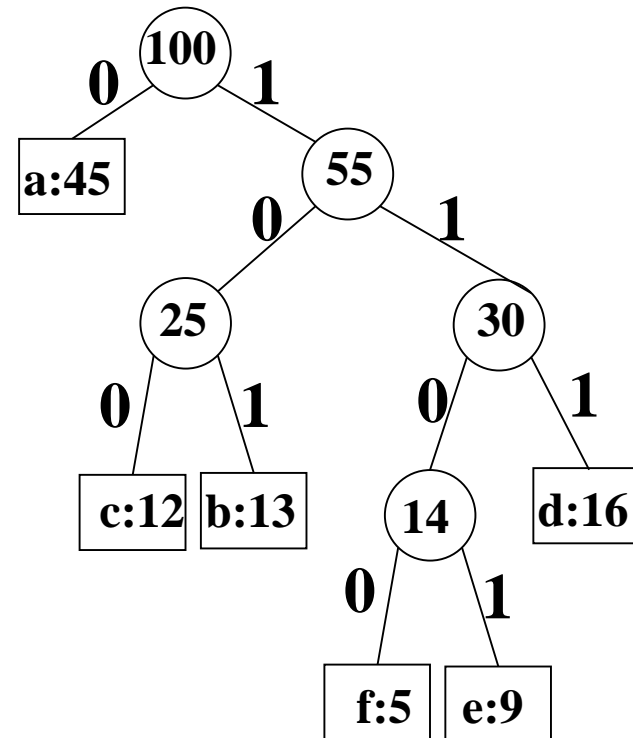
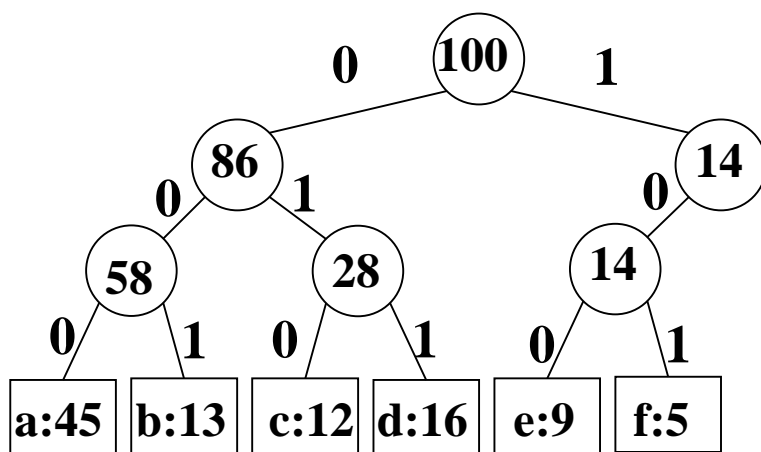
选择数据结构

二叉树:



好编码的特点

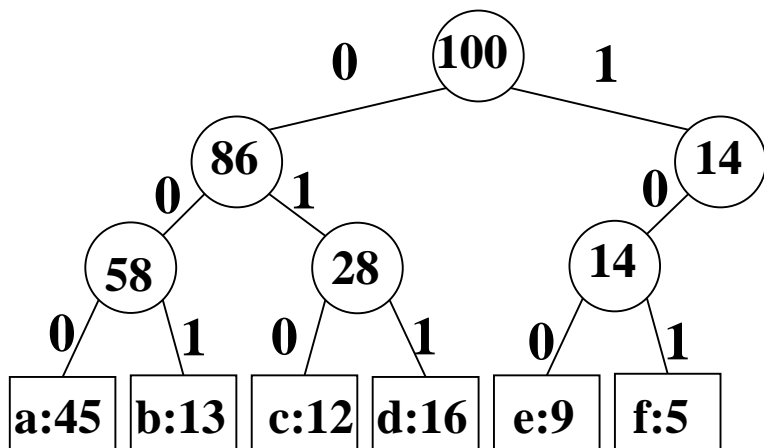
	a	b	c	d	e	f	存储
频率(千次)	45	13	12	16	9	5	100
定长码	000	001	010	011	100	101	300kb
变长码	0	101	100	111	1101	1100	224kb



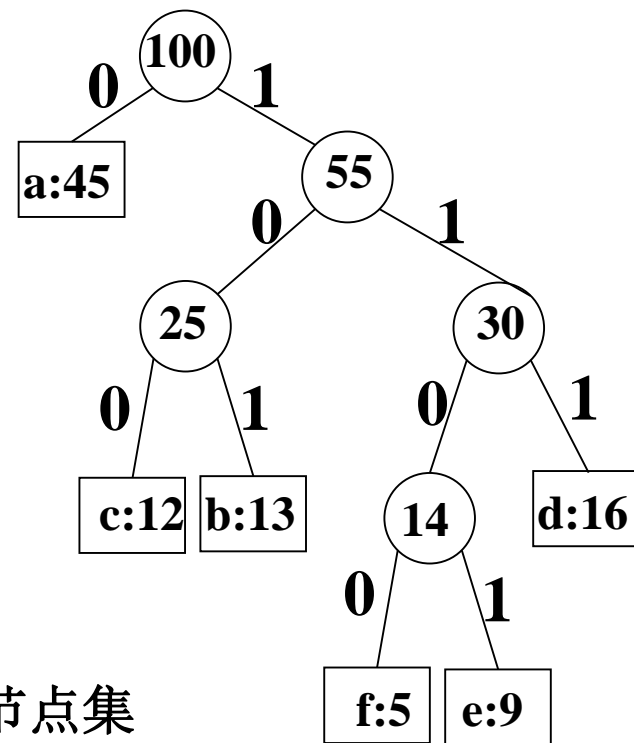
- 大频率编码短
- 小频率编码长
- 正则二叉树
- 相对平衡

构造优化模型

	a	b	c	d	e	f	存储
频率(千次)	45	13	12	16	9	5	100
定长码	000	001	010	011	100	101	300kb
变长码	0	101	100	111	1101	1100	224kb



叶节点深度
=
编码长度

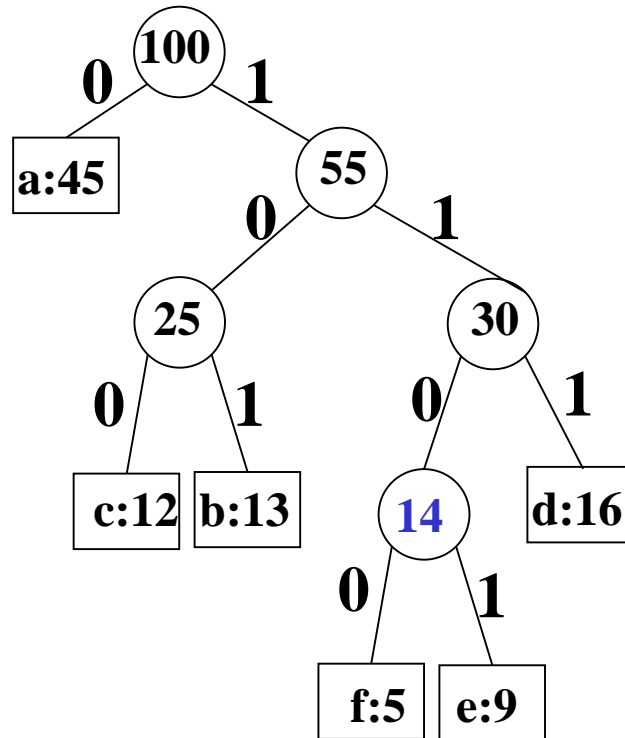
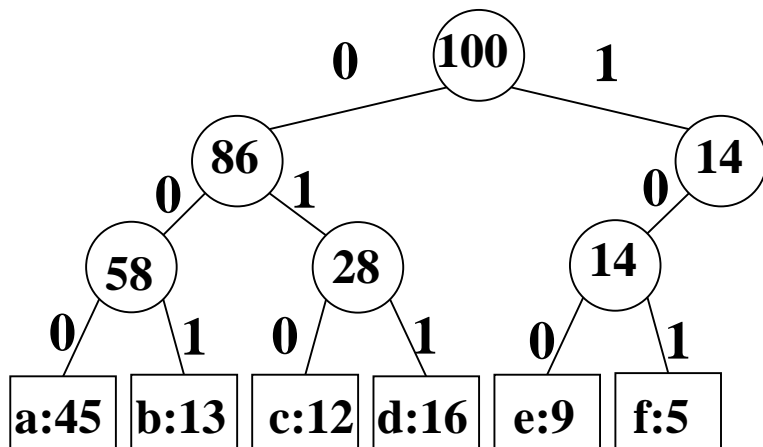


优化目标: 叶节点频率乘深度的和最小

$$\min_T \sum_{c \in C} f(c) d_T(c)$$

其中 $f(c)$ 是 c 的频率,
 $d_T(c)$ 是 c 的深度, C 是叶节点集

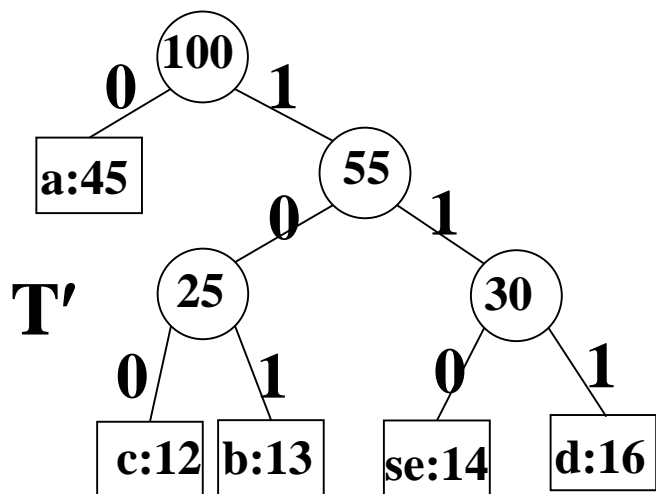
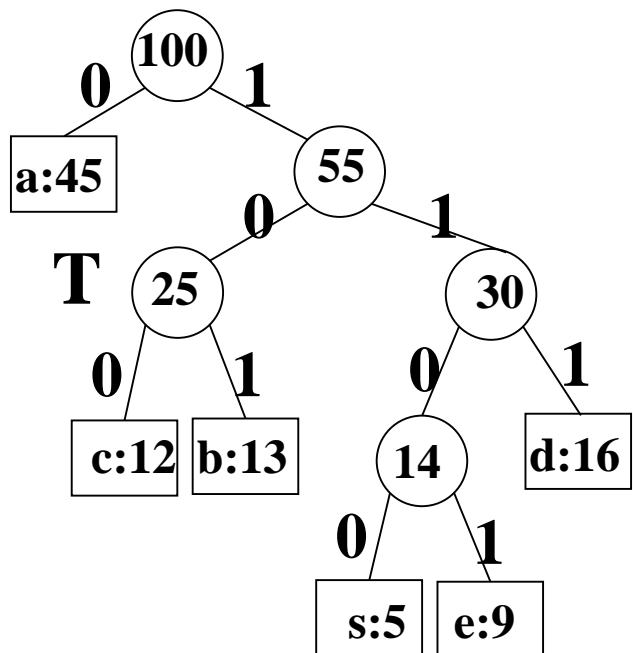
观察规律



$$\min_T \sum_{c \in C} f(c) d_T(c)$$

- 正则二叉树
- 相对平衡
- $f(c)$ 最小?
- $f(c)$ 最小两符号深度最大(贪心选择性质)
- 接下来怎么做? 构造子结构? **OSP**
- 合并fe得新二叉树(子结构)

最优编码T的性质



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

- 正则二叉树
- **贪心选择性质:**
若 $f(c)$ 最小, 则 $d_T(c)$ 最大
 $f(c)$ 最小的两个是兄弟(s,e)
- **OSP:**
构造T的子结构T'
合并T中 $f(c)$ 最小的两个得T'
 $B(T) = B(T') + f(s) + f(e)$?
若T最优, 则T'最优

Huffman算法

输入 $C[1:n]$, $f[1:n]$ //字母集 C 和频率 f

1. $Q=C$ //建优先队列 Q , $O(n)$

2. 对 $i = 1:n-1$ //循环 $n-1$ 次

3. 分配新节点 z

4. 取出 Q 中 f 最小的 x 作为 z 的左孩子 // $O(\log n)$

5. 取出 Q 中 f 最小的 y 作为 z 的右孩子 // $O(\log n)$

6. $f[z] = f[x] + f[y]$

7. 将 z 插入 Q 中 // $O(\log n)$

8. 取出 Q 中节点作为树根 //?

第四章 贪心算法

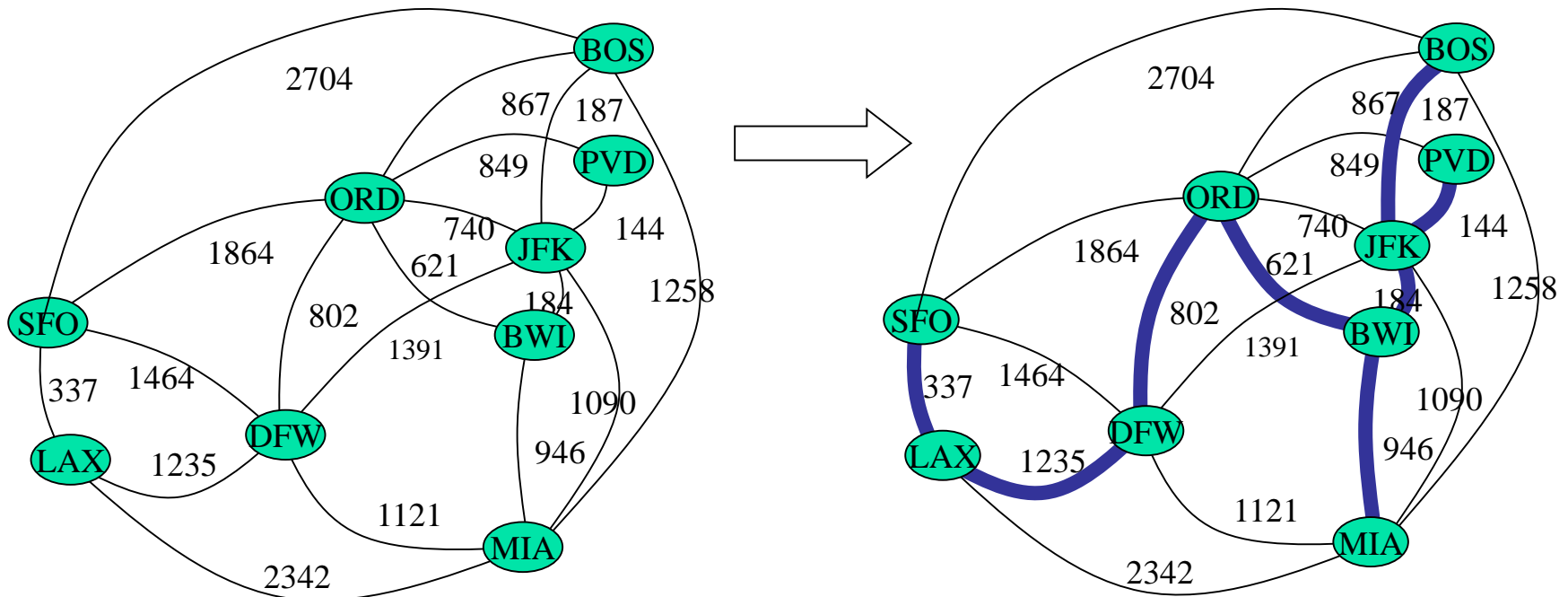
Greedy

!!贪心不一定正确(0-1背包), 需要证明

1. 活动安排问题
2. 贪心算法基本要素
3. Huffman算法
4. 最小生成树
5. 多机调度问题

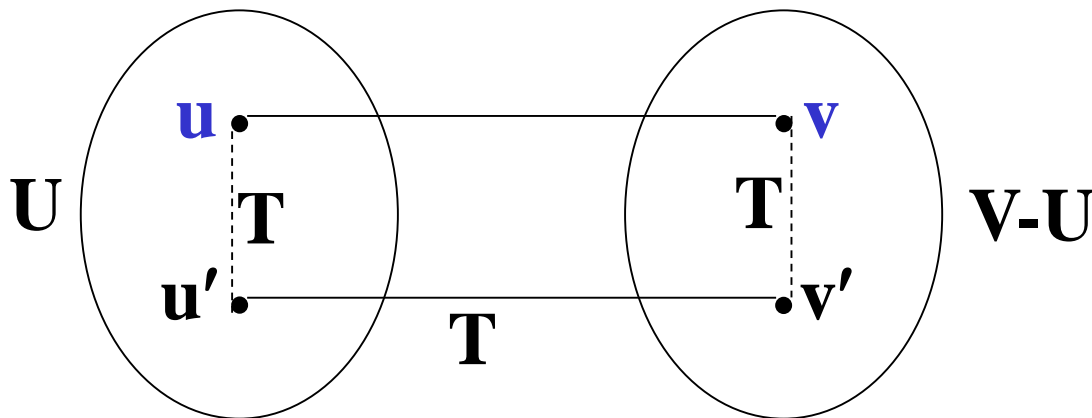
最小生成树(minimum spanning tree)

- 无向连通带权图 $G=(V,E,w)$.
- G 的生成树是 G 的包含所有顶点的一颗子树
- 若 G 的生成树 T' 在所有 G 的生成树中各边权总和最小, 则 T' 称为 G 的最小生成树(MST)



MST性质(贪心选择+归纳)

- 定理: 无向连通带权图 $G=(V,E,w)$. $U \subset V, A \subset E$.
若 A 是某MST的子集, A 没有边连接 U 和 $V-U$,
且 (u,v) 是 G 中连接 U 和 $V-U$ 的权最小的边,
则 G 存在MST包含 A 和边 (u,v) .
- 证明: 如图, 设 T 是包含 A 的MST, $(u,v) \notin T$, $(u',v') \in T$.
则必有 $w(u,v) \leq w(u',v')$. 由 T 去 (u',v') 添 (u,v) 得树 T' .



T' 也是MST
包含 A 和 (u,v)

MST算法正确性证明

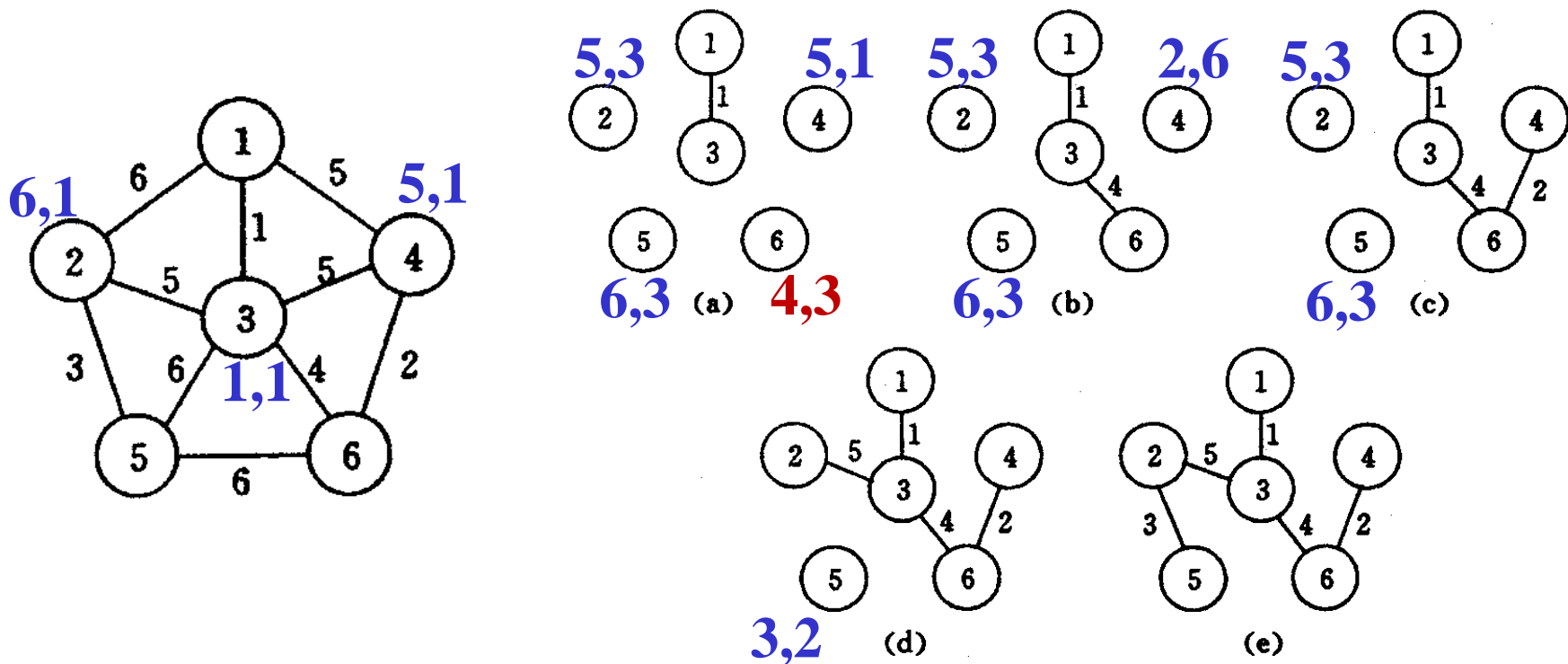
- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST
- MST性质: $U \subset V, A \subset E$. 若 A 是某MST的子集, $A \cap (U, V-U) = \emptyset$, (u,v) 在 $(U, V-U)$ 中权最小, 则存在MST包含 A 和 (u,v) .
- Prim: 维护 $U \subseteq V$, 初始任取 V 中一点
 1. 当 $U \neq V$, 取 U 到 $V-U$ 的最小权边 (u,v) , 将 v 加入 U

正确性证明: U, A 是已有边集
- Kruskal: 维护 $A, Q \subseteq E$, 初始 A 为空集, $Q=E$ 升序排列
 1. 当 Q 非空, 取 Q 最小权边 (u,v) ,
 2. 若 (u,v) 加入 A 会形成回路, 舍去; 否则加入 A .

正确性证明: $A, U=\{\text{已经挑出边中与 } u \text{ 连通的节点}\}$

Prim算法

- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST
- 维护 $U \subseteq V$, 初始任取 V 中一点
当 $U \neq V$, 取 U 到 $V-U$ 的最小权边 (u,v) , 将 v 加入 U .



Prim算法 $O(|V|^2)$

- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST

- 维护 $U \subseteq V$, 初始任取 V 中一点 r

当 $U \neq V$, 取 U 到 $V-U$ 的最小权边 (u,v) , 将 v 加入 U .

$key[u]$ 记 u 到 U 的最小距离, $G[u]$ 记 u 与 U 最小距离对应点

1. 初始 $G[u]=NIL$, $key[r]=0$, 其它 $key[u]=INF$, $Q=V$, U 空

2. 当 Q 非空

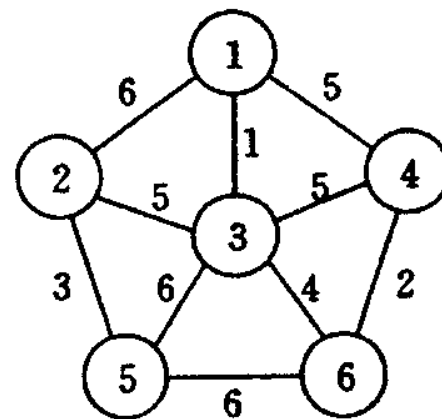
3. 取出 Q 中 u 使得 $key[u]$ 最小, 加入 U

4. 对 u 的每个邻居 v ,

5. 若 $v \in Q$ 且 $w[u,v] < key[v]$

6. 则 $key[v] = w[u,v]$, $G[v] = u$

Q 一般数组 $O(|V|^2 + |E|) = O(|V|^2)$

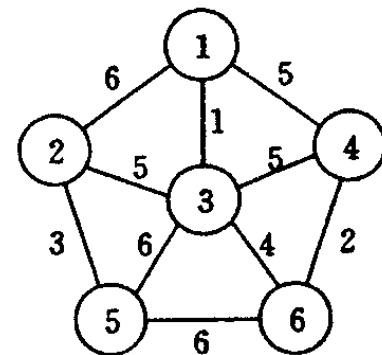


Prim算法 $O(|E|\log|V|)$

- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST

- 维护 $U \subseteq V$, 初始任取 V 中一点 r

当 $U \neq V$, 取 U 到 $V-U$ 的最小权边 (u,v) , 将 v 加入 U .



$key[u]$ 记 u 到 U 的最小距离, $G[u]$ 记 u 与 U 最小距离对应点

1. 初始 $G[u]=NIL$, $key[r]=0$, 其它 $key[u]=INF$, $Q=\{r\}$, U 空

2. 当优先队列(最小堆) Q 非空

3. 删除 Q 堆顶元素 u . 若 u 在 U 中, 则 **break**; 否则加入 U .

4. 对 u 的每个邻居 v ,

5. 若 $w[u,v] < key[v]$

6. 则 将 v 按键值 $key[v]=w[u,v]$ 插入 Q , $G[v] = u$

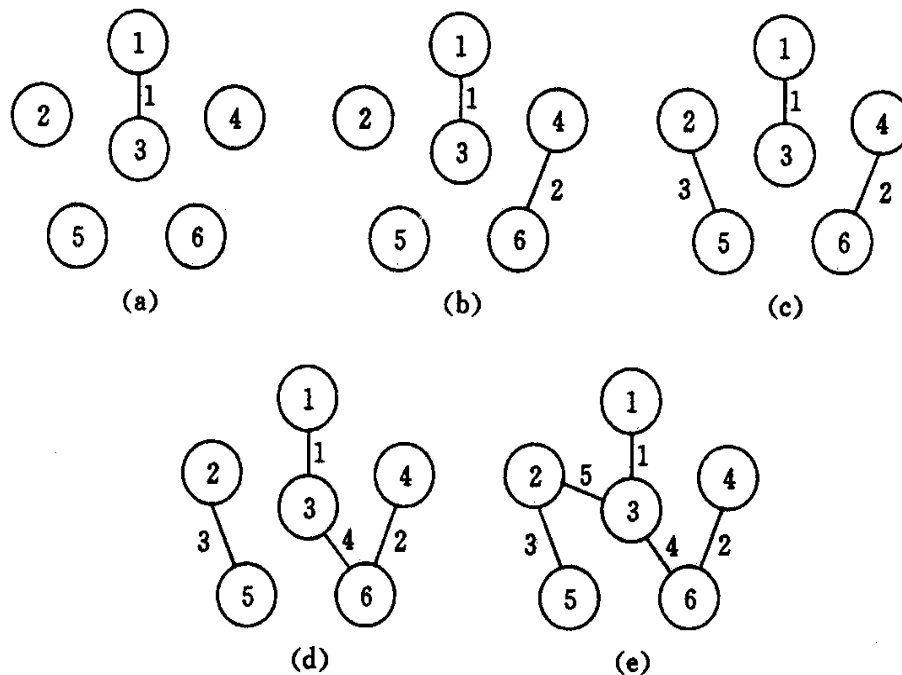
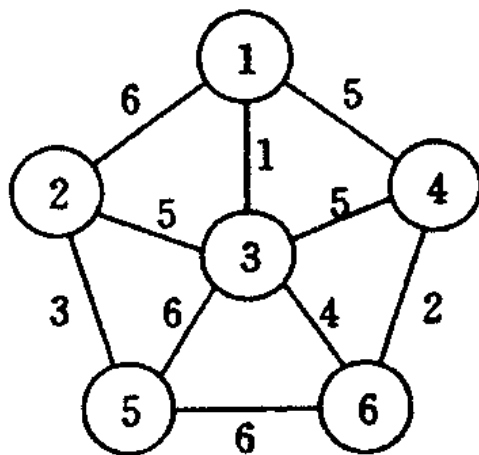
Q 优先队列 $O(|V|\log|E| + |E|\log|E|) = O(|E|\log|E|)$

Kruskal算法

- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST
- 维护 $A, Q \subseteq E$, 初始 A 为空集, $Q=E$ 升序

$Q = \{ 13, 46, 25, 36, 34, 23, 14, 12, 35, 56 \}$

当 Q 非空, 取 Q 最小权边 (u,v) , 若可行, 加入 A .



Kruskal算法细节 $O(|E|\log|V|)$

- 输入: $G=(V,E,w)$ 无向连通; 输出: G 的MST

- 维护 $A, Q \subseteq E$, 初始 A 为空集, $Q=E$ 升序

当 Q 非空, 取 Q 最小权边 (u,v) , 若可行, 加入 A .

1. A 为空, $Q=E$ 按边权升序排列

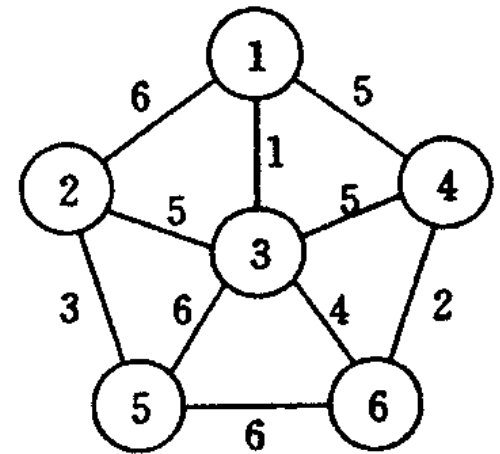
2. 当 Q 非空

3. 顺序取 Q 中边 (u,v)

4. 若 u,v 在不同连通分支(检查),

5. 则添 (u,v) 到 A , 合并 u,v 所在连通分支,

6. 输出 A



并查集算法处理连通分支, $O(|E|\log|E|+|E|\log^*|V|)$

并查集算法(Make-set, Find-set)

p[x]: x的父亲; **rank[x]**: x的阶; **Find(x)**: 找x的根

Make(x)

1 **p[x]=x**

2 **rank[x]=0**

Union(x,y)

1 **Link(Find(x), Find(y))**

Link(x,y) //合并两树根

1 若 **rank[x]>rank[y]**

2 则 **p[y]=x**

3 否则 **p[x]=y**

4 若 **rank[x]=rank[y]**

5 则 **rank[y]=rank[y]+1**

性质1:
只有树根
rank值
会增加

性质2:
树根阶
 $= r$
 \downarrow
树节点数
 $\geq 2^r$.

Find(x)

1 若 **x≠p[x]**, 则

2 **p[x] = Find(p[x])**

3 返回 **p[x]**

路径压缩(**pc**)技术

n个节点, m次操作. 不计**pc**: $O(m \log n)$; 计**pc**: $O(m \alpha(n))$

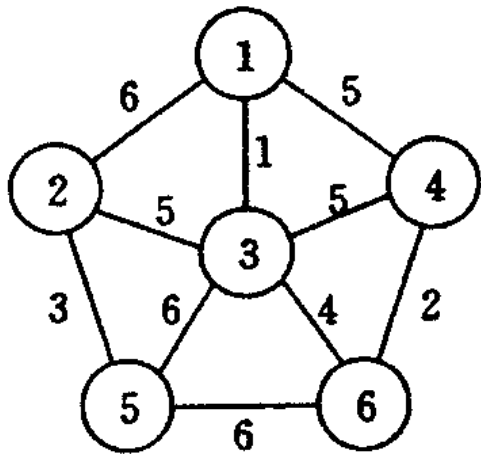
回顾定义($b>1$) $2^{^b}=2^{2^{(b-1)}}$, $2^{^1}=2$, $\log^*(2^{^b}) = b$

加入并查集结构的Kruskal算法

1. A为空, $Q=E$ 按边权升序排列, 每个点是一颗树
 2. 当Q非空
 3. 顺序取Q中边 (u,v)
 4. 若 u,v 在不同树中, 则添 (u,v) 到A, 合并 u,v 所在树,
 5. 输出A
-

1. A为空, $Q=E$ 按边权升序排列, $\forall x$ Make(x)
2. 当Q非空
3. 顺序取Q中边 (u,v)
4. 若 $\text{Find}(u) \neq \text{Find}(v)$, 则添 (u,v) 到A, Union(u,v),
5. 输出A

Kruskal: 取边, 查找, 合并



$Q = \{13, 46, 25, 36, 34, 23, 14, 12, 35, 56\}$

1:0 2:0 3:0 4:0 5:0 6:0

查找合并
(1,3) (4,6) (2,5)

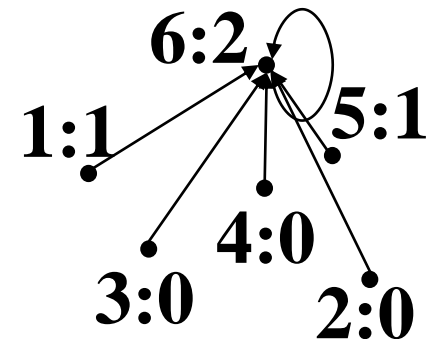
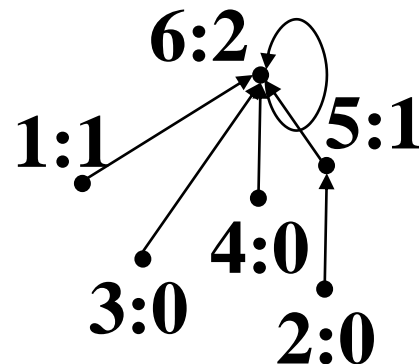
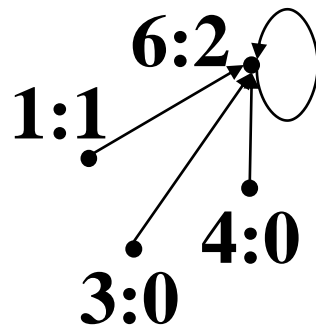
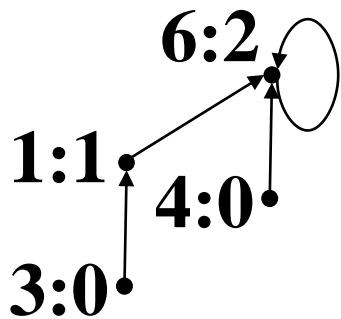
1:1 6:1 5:1
3:0 4:0 2:0

查并(3,6)

查(3,4)

查并(2,3)

查(1,2)



[刘]n个节点, m次操作, $O((n+m)\log^*n)$

n个节点分 \log^*n 块: 第b块= $\{x: 2^{b-1} < \text{rank}[x] \leq 2^b\}$

性质1: 只有树根rank会增加, 非树根rank不变.

性质2: 树根阶 = $r \Rightarrow$ 树节点数 $\geq 2^r$.

性质3: 至多有 $n/2^r$ 个rank为r的点, rank值不大于 $n-1$

性质4: 第b块至多 $n/(2^{b-1})$ 个点.

设Find(x_0)依次经过 x_0, x_1, \dots, x_{L-1} 到根节点 x_L . x_i 分四类

(1)根 (2)根的儿子 (3)非根或根的儿子, x_i 和父不在同块

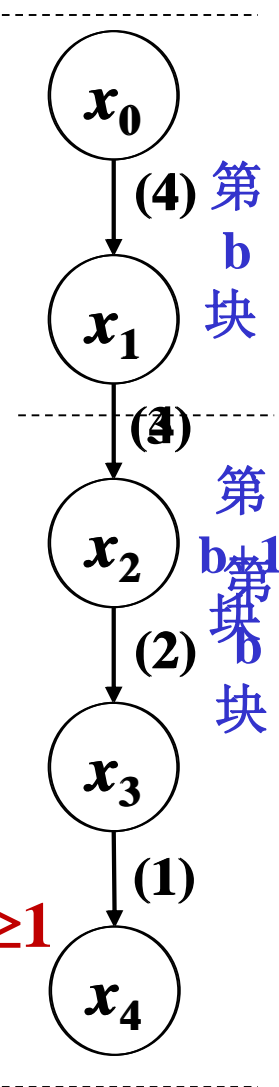
(4) 非根或根的儿子, x_i 和父在同块

Find(x_0): (1)贡献1次; (2)贡献1次; (3)贡献 $\leq \log^*n$ 次.

(4)贡献: 求m次操作总和的上界 // x_i : 阶不变, 父变, 父阶增 ≥ 1

• 设 x_i 和父都在第b块, 则 x_i 对(4)的贡献 $\leq 2^{b-1}$ 次;

• 第b块所有元素对(4)的贡献 $\leq (n/2^{b-1}) \times 2^{b-1} = n$ 次.



第四章 贪心算法

Greedy

!!贪心不一定正确(0-1背包), 需要证明

1. 活动安排问题
2. 贪心算法基本要素
3. Huffman算法
4. 最小生成树
5. 多机调度问题

多机调度问题

输入: $n, m, t[1:n]$, n 个作业, m 台机器

输出: 最早完成全部作业的调度方案.

要求: 作业不能拆分

Multiprocessor scheduling, NP完全, 现无有效的解法

近似算法: **LPT(Longest Processing Time)**算法

$n \leq m$: 一个机器一个作业

$n > m$: 作业按时长降序排列, 依次分配空闲机器

近似比例 $\leq 4/3 - 1/(3m)$. 例 $m=1, m=2$ (**11,10,7,7,7**)

PTAS: $\forall \epsilon > 0$, 有 $(1+\epsilon)$ -优化的多项式时间近似算法

APX: $\exists c > 1$, 有 c -优化的近似算法. 若 $P \neq NP$, $PTAS \subset APX$

Gone fishing

在一条水平路边, 有 n 个钓鱼湖, 从左到右编号 $1-n$.

John有 H 小时的时间, 希望钓到尽量多的鱼.

从湖1出发向右, 在一些湖边停留一定时间钓鱼.

求钓最多鱼的方案(多余时间都给湖1).

($i=1:n$)在第 i 个湖中钓鱼, 第1个5分钟可以钓到 $F[i]$ 条鱼,

以后每钓5分钟鱼, 可钓到的鱼数减少 $D[i]$.

($i=1:n-1$)从第 i 个湖到第 $i+1$ 个湖需要 $5 \times T[i]$ 分钟.

输入: n // $n=0$ 输入结束

H

$F[1:n]$

$D[1:n]$

$T[1:n-1]$

$n \dots$

2

1

10 1

2 5

2

45, 5

Number of fish expected: 31

1 2 3 4 5 6 7 8 9 10 11 12

1 1 1 1 1 1 1 1 1 → → 2

10 8 6 4 2 0 0 0 0 - - 1 合计31

本章小结

贪心不一定正确, 需要证明

1. 活动安排问题
2. 贪心算法基本要素
3. Huffman算法
4. 最小生成树
5. 多机调度问题

本章作业

1. 字符a~h出现的频率恰好是前8个Fibonacci数, 它们的Huffman编码是什么? 将结果推广到n个字符的频率恰好是前n个Fibonacci数的情形.
2. 若在0-1背包问题中, 各物品依重量递增排列时, 其价值恰好降序排列, 对这个特殊的0-1背包问题, 设计一个有效算法找出最优解, 并说明算法的正确性.
3. 将最优装载问题的贪心算法推广到2艘船的情形 贪心算法还能产生最优解吗?
4. 最优分解问题. 见下页

本章作业

4. 最优分解问题.

问题描述: 设 n 是一个正整数, 将 n 分解为若干互不相同的自然数之和, 且使这些自然数的乘积最大.

算法设计: 对于给定的正整数 n , 计算最优分解方案.

数据输入: 由文件input.txt提供输入数据.

文件只有一行, 是正整数 n .

结果输出: 将计算的最大乘积输出到文件output.txt

例如若 $n=10$, 则最优分解为 $2+3+5$, 最大乘积为30.

附录一

最小堆自底向上建堆 $\Theta(n)$

数据结构

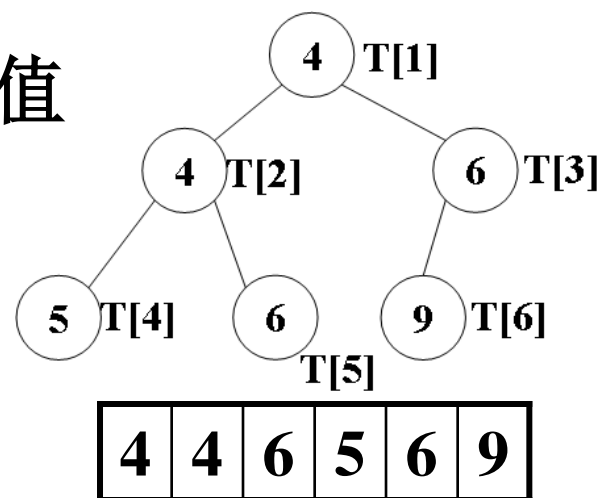
- 数组, 链表, 二叉树(堆和搜索树), hash表
- (最小)堆: 每个节点的关键值

小于或等于左右子节点的关键值

操作: 插入, 删除, $O(h)$

一般隐式存储, 即数组

$T[i]$ 的左孩子 $T[2i]$, 右孩子 $T[2i+1]$

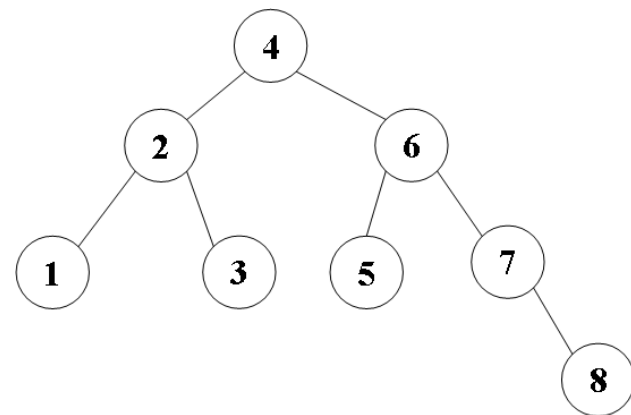


- 二叉搜索树: 每个节点的关键值

大于所有左分支节点关键值,

小于所有右分支节点关键值

操作: 查找, 插入, 删除, $O(h)$

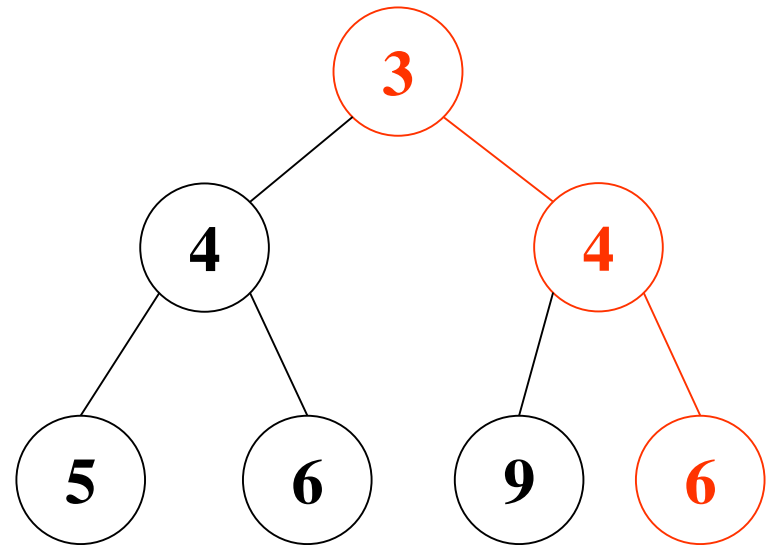
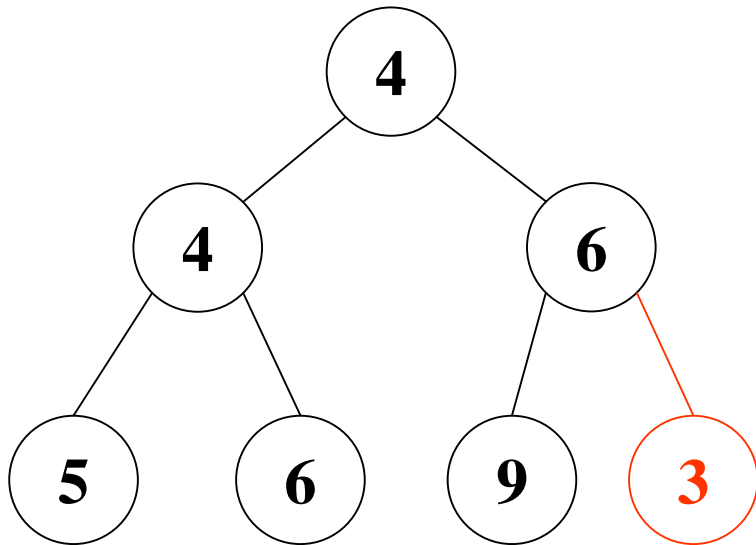


最小堆的插入操作

例如：在下面的堆中插入3

过程：先放在最后，再调整

1	2	3	4	5	6	7
4	4	6	5	6	9	3
4	4	3	5	6	9	6
3	4	4	5	6	9	6



堆的删除操作

例如: 对下面的堆中作删除

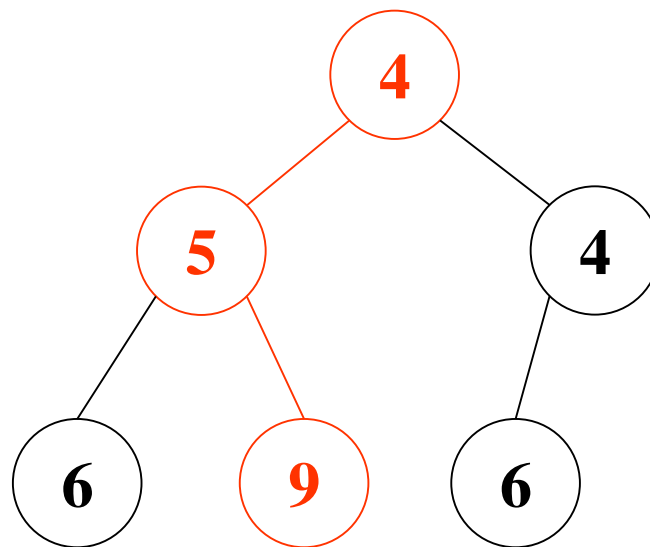
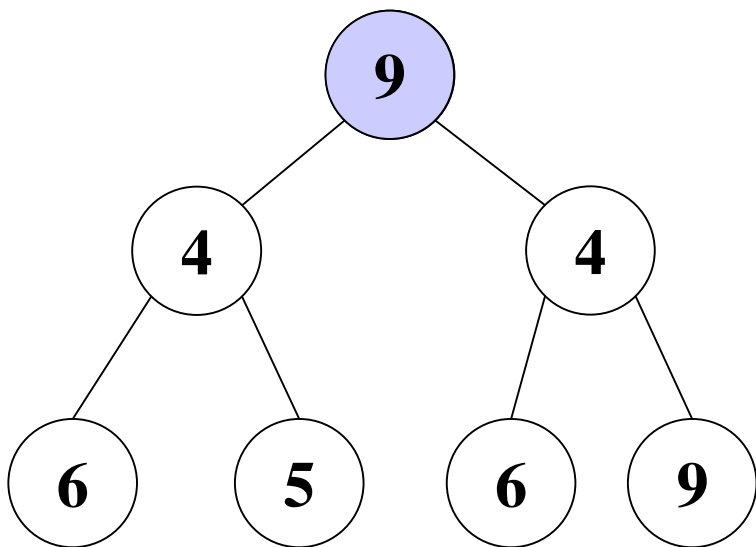
过程: 取出 $T[1]$,

将 $T[n]$ 赋给 $T[1]$, 调整

特点: 仍然是连续片段

- 堆顶以外元素位置很乱

1	2	3	4	5	6	7
3	4	4	6	5	6	9
9	4	4	6	5	6	
4	9	4	6	5	6	
4	5	4	6	9	6	



自顶向下建堆

	1	2	3	4	5	6	7	8
1	4	2	6	1	3	5	7	8
2	2	4	6	1	3	5	7	8
3	2	4	6	1	3	5	7	8
4	1	2	6	4	3	5	7	8
5	1	2	6	4	3	5	7	8
6	1	2	5	4	3	6	7	8
7	1	2	5	4	3	6	7	8
8	1	2	5	4	3	6	7	8

$$\sum_{i=1}^n \log i \geq \frac{n}{2} \log \frac{n}{2}$$

$$\sum_{i=1}^n \log i = \Theta(n \log n)$$

自底向上建堆

	1	2	3	4	5	6	7	8
1	4	2	6	1	3	5	7	8
2	4	2	6	1	3	5	7	8
3	4	2	6	1	3	5	7	8
4	4	2	6	1	3	5	7	8
5	4	2	6	1	3	5	7	8
6	4	2	5	1	3	6	7	8
7	4	1	5	2	3	6	7	8
8	1	2	5	4	3	6	7	8

$$T(n) = 2T(n/2) + \log n = \Theta(n)$$

本质上是分治(奇怪的分)

用数学归纳法容易证明:

$$T(n) \leq c_2 n - 2 \log n$$

$$T(n) \geq c_1 n$$

附录二

活动安排问题的动态规划算法

活动安排问题的动态规划算法

n个活动申请一个活动室, 活动起始终止区间(s_i, f_i)

- 输入: $n, (s_i, f_i), i = 1:n$, 输出: 最大相容活动子集
- 设 $f_1 \leq f_2 \leq \dots \leq f_n$, 添加 $f_0 = 0, s_{n+1} = \infty$
- $S[i,j] = \{\text{在活动}i\text{结束后和活动}j\text{开始前能添加的活动}\}$
- $c[i,j] = S[i,j]$ 中的最大相容活动子集
- 输出 = $c[0, n+1]$
- 若 $i \geq j$, 则 $S[i,j] = \emptyset$.

$$c[i,j] = \begin{cases} 0 & \text{若 } S[i,j] = \emptyset \\ \max_{i < k < j} \{c[i,k] + c[k,j] + 1\} & \text{否则} \end{cases}$$