

数据结构与算法分析

华中科技大学软件学院

2017年秋

大纲

- 1 树的表示
- 2 二叉树
- 3 二叉查找树
- 4 平衡二叉树
- 5 B树

课程计划

- 已经学习了
 - 线性表的数组与链表实现
 - 桶式排序与基数排序
 - 堆栈及其应用
 - 队列

课程计划

- 已经学习了
 - 线性表的数组与链表实现
 - 桶式排序与基数排序
 - 堆栈及其应用
 - 队列
- 即将学习
 - 树的表示
 - 二叉树，树的遍历，决策树
 - 二叉查找树
 - 平衡二叉树：AVL树、红黑树
 - B-树，B+树

树



Roadmap

1 树的表示

2 二叉树

3 二叉查找树

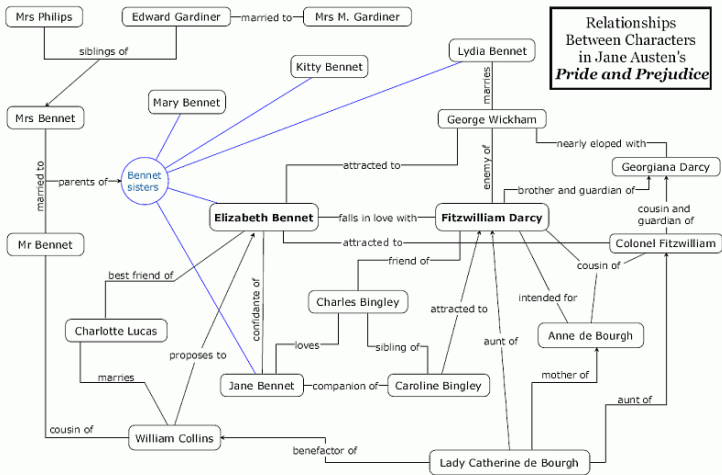
4 平衡二叉树

5 B树

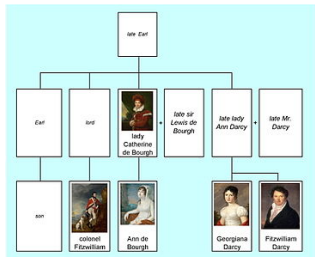
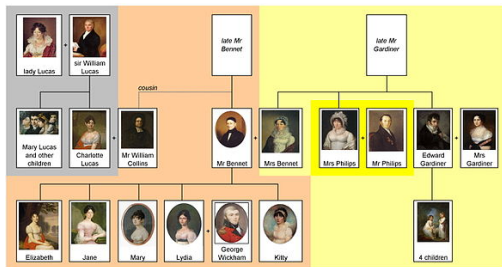
树

- 已经学习了有关树的计算/循环
- 其他的树：
 - 文件夹/子域名的层次结构
 - 表示中缀表达式
 - XML（可扩展标记语言）
 - 系谱图，组织架构图
 - 字典搜索
- 树是一系列用特定方式连结的节点

关系图表



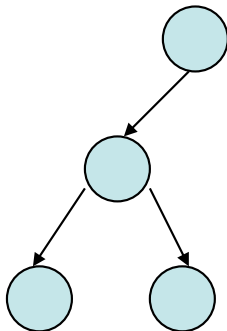
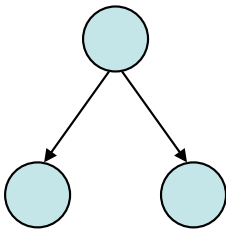
系谱图



树的定义

- 递归定义：
 - $\{\}$ 是一棵树 (空树);
 - 一个结点 r , 有 ≥ 0 非空子树 T_i , 并且有直接的、从 r 到 T_i 的边缘 (连线)
- 画法对于自然的树是颠倒的
 - 顶部的结点 = 树根 (root)
 - 底部的结点 (子结点) = leaves
- 树 (或一个结点) 有一个高度 (height), 即从根出发的最长路径长度
 - 从根到底 (或结点) 的连线数最大值
 - $\text{Height}(\text{null}) = -1$, $h(\text{root}) = 0$

树的递归定义

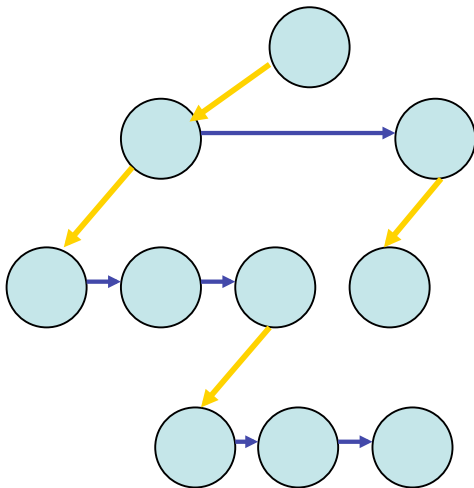


树的实现

```
typedef struct TreeNode *PtrToNode;  
typedef struct treeNode  
{  
    Elementype element;  
    PtrToNode firstChild;  
    PtrToNode nextSibling;  
} TreeNode;
```

- Tree是TreeNode的一个带有成员的实例
- 如果R是树根，那么R.nextSib是什么？

父子兄弟树



文件结构

- 文件/目录(File/directory) 列表
- 将目录视为一种特殊的树

```
void listDir (DirectoryOrFile D, int depth)
{
    printName (D, depth);
    if (isDirectory (D))
    {
        for each child C in dir D
            listDir (C, depth+1);
    }
}

void listAll (DirectoryOrFile D)
{
    listDir (D, 0);
}
```

文件结构举例

Desktop(桌面)

- My documents (我的文档)
 - My music
 - My pictures
- My computer
 - C
 - Documents and settings
 - Gang Shen
 - desktop
 - Program files
 - D

Roadmap

1 树的表示

2 二叉树

3 二叉查找树

4 平衡二叉树

5 B树

二叉树

一般地，限制容许的子结点个数为2以内（可为0或2）

- 每个结点都有2个子结点，可能为空
- 空值（Nulls）不画

```
typedef struct binTreeNode
{
    ElementType element;
    struct binTreeNode *left;
    struct binTreeNode *right;
} BinTreeNode;
```

二叉树的遍历

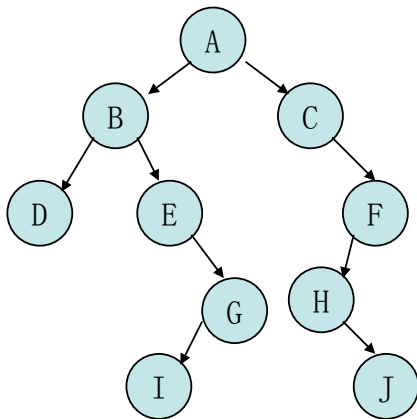
- 三种遍历二进制树的所有结点方法（深搜优先式）
 - 先序
 - 后序
 - 中序
- 取决于访问父结点的次序：先/后/在它的子结点之间 以先序，举例
 - 访问当前结点
 - 访问左子树
 - 访问右子树

二叉树的先序遍历

```
void pre_order (BinTreeNode *root)
{
    if (root == NULL)
        return;

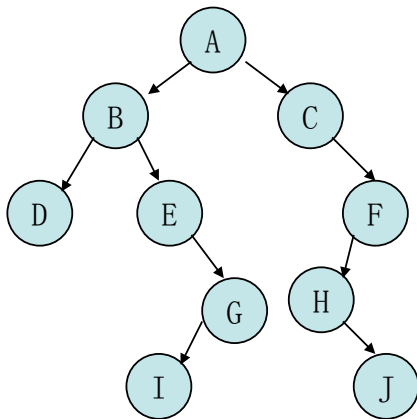
    print ("%c, □", root->element);
    pre_order (root->left);
    pre_order (root->right);
}
```

遍历



- 先序(preorder):
- 中序(inorder):
- 后序(postorder):

遍历

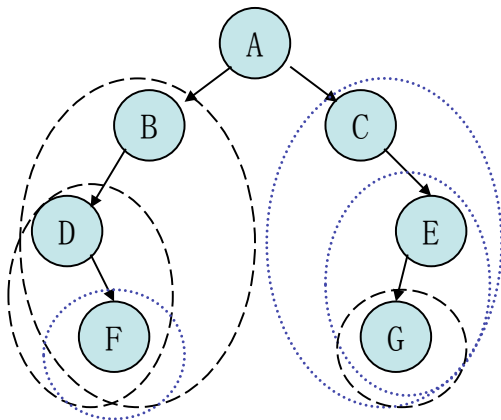


- 先序(preorder): A, B, D, E, G, I, C, F, H, J,
- 中序(inorder): D, B, E, I, G, A, C, H, J, F,
- 后序(postorder): D, I, G, E, B, J, H, F, C,

通过遍历结果构造二叉树

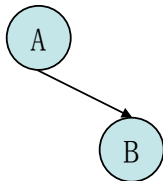
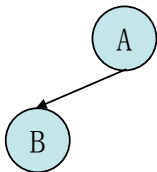
- 如果我们只有一棵树的横向结果(traversal results)，我们失去了什么信息？
- 一次移动(traversal)足以重建一棵树吗？
 - 组合(combos)怎么样？
- 考虑一次中序移动和一次前/后序移动，我们应如何构造一棵树？
- 中序：D, F, B, A, C, G, E,
- 相比先序：A, B, D, F, C, E, G,

构造一棵树



- In-order: **D**, **F**, **B**, A, **C**, **G**, **E**,
- Compare Pre-order: A, **B**, **D**, **F**, **C**, **E**, **G**,

先序+后序？



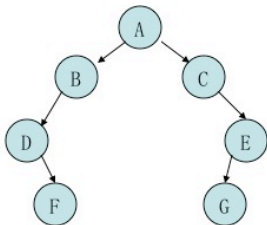
两棵不同的树享有同样的移动方式

- 先序Pre-order: A, B,
- 后序Post-order: B, A,

用数组表示二叉树

Parent->children

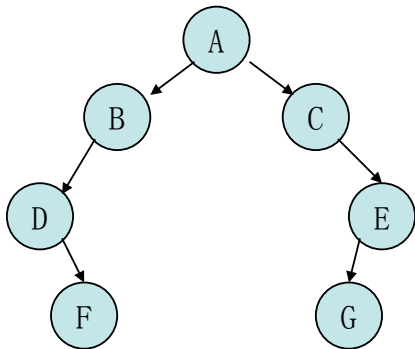
A	1	2
B	3	-1
C	-1	4
D	-1	5
E	6	-1
F	-1	-1
G	-1	-1



Children->parent

A	-1
B	0
C	100
D	1
E	102
F	103
G	4

用数组表示二叉树



$A[i]$: 当前结点, $A[2*i]$:左子树, $A[2*i+1]$:右子树

	A	B	C	D			E	F					G	
--	---	---	---	---	--	--	---	---	--	--	--	--	---	--

表达式树

- 代数表达式可以用树呈现
- $(a + (b * c)) + (((d * e) + f) * g)$
 - 难以直接描述、分析(parse)
 - 先前看到过的：可被转变为后缀，然后描述
- 构想(idea)：以树呈现
 - 运算数在树叶上
 - 运算在树根、内结点上(internal nodes)
- 然后移动来打印后缀
- 对于每个结点，打印左侧、右侧，然后元素

打印后缀表达式

采用后序遍历

```
void postfix (BinTreeNode *root)
{
    if (root != NULL)
    {
        postfix (root->left);
        postfix (root->right);
        print "(%c, ", root->element);
    }
}
```

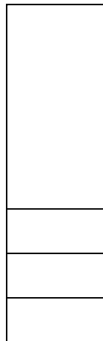
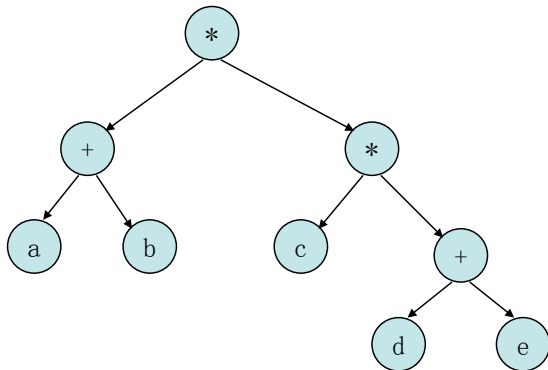
构造表达式树

- 对指定的一棵树, 很容易打印后序式. 如何得到树?
- 类似于从中缀infix→后缀postfix 的堆栈运算

```
while (!EOF)
{
    read variables push asTreeNode
    read operations pop t1, t2 to
        for a TreeNode and push new
        TreeNode (op,t1,t2)
}
```

Example: ab+cde+**

表达式树



二叉树的高度

$\text{Height}(T)$ = 从树根到树叶的最大距离

Theorem

二进制树的结点数 $n(T) \leq 2^{h(T)+1} - 1$.

证明：归纳法。基例，单结点： $n(T) = 1$,
 $2^{h(T)+1} - 1 = 2^{0+1} - 1 = 1$. 假定树高 $\leq H$ 为真。令
 $h(T) = H+1$, 那么 $n(T) = 1 + n(T_1) + n(T_2)$
 $\leq 1 + 2^{h(T_1)+1} - 1 + 2^{h(T_2)+1} - 1$ (归纳)
 $\leq 2 * \max(2^{h(T_1)+1}, 2^{h(T_2)+1}) - 1$
 $= 2 * 2^{\max(h(T_1)+1, h(T_2)+1)} - 1$
 $= 2 * 2^{h(T)} - 1$ (高度的定义)
 $= 2^{h(T)+1} - 1$

二叉树的高度

- 如果知道结点个数 # , 可求得最大高度
- 8 结点 高度至少为2
- $n(T) \leq 2^{h(T)+1} - 1, \log(n(T)) \leq h(T) + 1 - \log(1), \log(n(T)) - 1 \leq h(T)$
- 完全二进制树: 过半数结点为树叶
- 因为除最后一层外, 所有的层级都是内部的 (internal), 幂的和

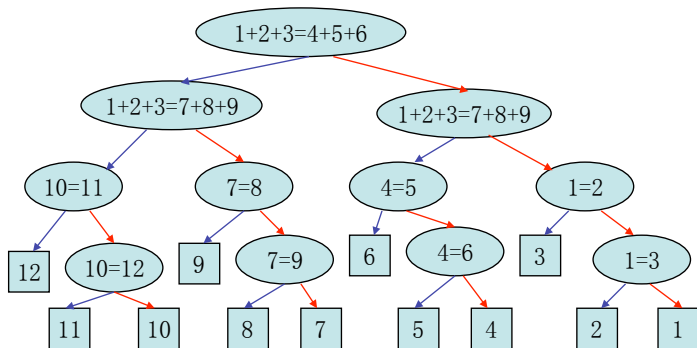
决策树

- 可以把很多算法看成决策树，叶子结点为决策点，中间节点为状态节点
 - 在每一层的状态节点进行判断
 - 根据测试结果进入某条支路
 - 最终停止
- 举个例子：有一个天平和12个球，其中一个球有缺陷。假设 这个有缺陷的球的重量与其他球不同，请问至少需要多少次比较才能找到？
- 画一个有12个叶子节点的决策树

决策树

- 分而治之：
 - 平凡情况：3个或4个球，两次称量可以得到结果
 - 归纳步骤：将问题结果划分成不同的组，每次称量都能减小问题规模
- 用数字标记每一个球：1, 2, 3, ...
- 决策树的叶子表示不合格球的12种可能结果
- 树高： $h(T) = \log_{12}$ ，表示决策轮次

决策树



→ \neq

→ $=$

Roadmap

1 树的表示

2 二叉树

3 二叉查找树

4 平衡二叉树

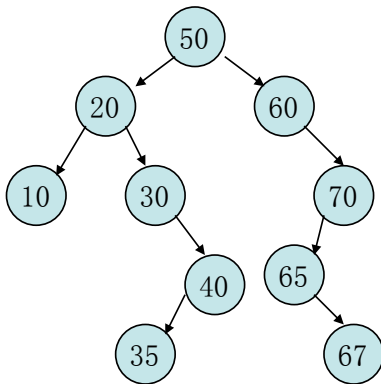
5 B树

二叉查找树

- 搜索过程快速，而插入/删除缓慢的数据结构
 - 大小有限的有序数组
- 插入/删除过程快速，而搜索缓慢的数据结构
 - 链表
- 二叉查找树：插入/删除过程快速，并且搜索过程快速
 - 没有大小限制
 - 如果平衡，所有操作的复杂度均为 $\log(n)$

- 二叉查找树 = 二叉树 + 一个新的属性
 - 对于一个节点n, n的左子树上所有结点的值均小于n的值; 右子树上所有结点的值均大于n的值
 - 对两棵子树上的所有节点均适用, 并非只针对两个直接儿子节点
- 但这不像听起来那么难: 在插入时完成
- 除此之外, 不需要副本记录, 也不需要计数器
- 回忆: $2^{h(T)+1} - 1 = N(T)$

BST 举例



- 先序:
- 中序:
- 后序:

BST中的查找

类似于折半查找

```
typedef struct BinTreeNode *SearchTree;
```

```
Position find (ElementType x, SearchTree T)
{
    if (T == NULL)
        return (NULL);

    if (x < T->element)
        return find (x, T->left);
    else if (x > T->element)
        return find (x, T->right);
    else return (T);
}
```

改写为迭代

仍然类似于折半查找

```
Position find (ElementType x, SearchTree T)
{
    while (T != NULL)
    {
        if (x < T->element)
            T = T->left;
        else if (x > T->element)
            T = T->right;
        else
            return (T);
    }
    return (NULL);
}
```

BST findmin/findmax

寻找最小元素：任给T，T的左儿子总是更小，右儿子总是更大

```
Position findMin (SearchTree T)
{
    if (T == NULL)
        return (NULL);

    if (T->left == NULL)
        return (T);

    return (findMin (T->left));
}
```

time = depth of tree $\approx \log n$

非递归形式代码

尾递归——在程序结尾处递归调用, 容易转换为迭代
想法: 一直向左走, 直到叶节点, 然后返回该节点的属性

把if(exp)改为while(! exp)

```
Position findMin (SearchTree T)
{
    if (T == NULL)
        return (NULL);

    while (T->left != NULL)
        T = T->left;

    return (T);
}
```

BST插入

想法:创建新节点, 或者返回x已存在

```
SearchTree insert (ElementType x, SearchTree T)
{
    if (T == NULL)
    {
        T = malloc (sizeof (TreeNode));
        T->element = x;
        T->left = NULL;
        T->right = NULL;
    }
    else if (x < T->element)
        T->left = insert (x, T->left);
    else if (x > T->element)
        T->right = insert (x, T->right);

    return (T);
}
```

非递归插入

```
void insert (ElementType x, SearchTree T)
{
    SearchTree *P;

    while (T != NULL)
    {
        if (x < T->element)
        {
            P = &(T->left);
            T= T-> left
        }
        else if (x > T->element)
        {
            P = &(T->right);
            T= T-> right
        }
        else
            return;          /* already existing */
    }
    *P = (SearchTree)malloc (sizeof (TreeNode));
    *P->element = x;
    *P->left = NULL; *P->right = NULL;
}
```

BST删除

- 通常，最困难的情况是：插入总是发生在底部，但是删除可能发生在任何地方
- 首先，找到该节点，如果子节点个数 ≤ 1 ，那么就很容易了
 - 0：直接删除
 - 1：用子节点直接取代该节点
- 如果有两个子节点，就相对复杂
 - 策略：用右子树最小的元素(或左子树中最大的元素)替换已删除的节点
 - 对于这个节点，我们必须删除它(简单)

删除代码

```
SearchTree Delete (ElementType x, SearchTree T)
{
    Position tmpCell;

    if (T == null)
        return (NULL);

    if (x < T->element)
        T->left = delete (x, T->left);
    else if (x > T->element)
        T->right = delete (x, T->right);
    else if (T->left != NULL && T->right != NULL)
    {
        tmpCell = findMin (T->right);
        T->element = tmpCell->element;
        T->right = delete (T->element, T->right);
    }
    else T = T->left != null ? T->left : T->right;

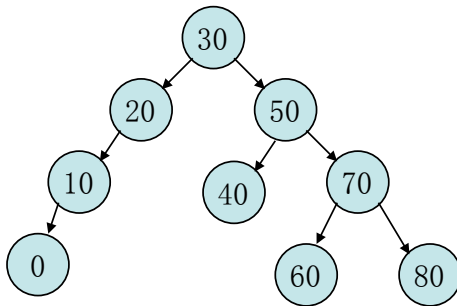
    return (T);
}
```


插入与删除

- 插入 30, 20, 50, 70, 10, 40, 60, 80, 0
- 删除 0, 20, 30
- 然而删除是困难的, 一个解决办法是: 懒惰删除
- 不真正地删除它, 只是将它标记为删除状态
 - 如果重复, 就减小计数(count)
 - 常用于数据库中

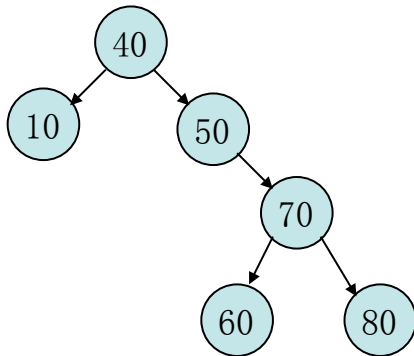
插入

Insert 30, 20, 50, 70, 10, 40, 60, 80, 0



删除

Delete 0, 20, 30



BST平均复杂度

- 树每下降一层，就将剩余节点分为相同规模的两部分（假设每个非叶节点都有两个子节点）
- 这表明：对于二叉查找树，深度 = $O(\log n)$
- 但是这对所有情况都适用吗？
- 给定向右的 n 次连续插入（插入的数要满足什么条件？），树的深度将为 n

BST最坏复杂度

- BST属性必须始终保持
- 有多少棵树包含节点1 2 3?
 - 值集合 (Value set) 不能决定树
 - 但值列表 (Value list) 可以
- 如果很多节点都只有一个子节点呢?
 - 如果BST不平衡呢?
 - 如果插入1 2 3 4 5, 或5 4 3 2 1呢?
- 在更糟糕的情况下, BST退化成一个链表
 - 搜索: $O(n)$
 - 插入: $O(n)$
 - 删除: $O(n)$
- 实际上, BST的表现可能会很糟糕

实验三：约瑟夫问题

- N人围成一个圈，每经过K个人就消灭一个人，谁将是最后的幸存者？
- 将每个人依次标记为 $0, 1, 2, \dots, n - 1$ ，用 $J(n, k)$ 表示当共有 n 人时幸存者的标号
- 首先删除标记为 $k - 1$ 的人，把最初标记为 k 的人标记为 0 ，并从 0 开始，重新标记其余的标签。

0	1	2	3	...	$k-2$	$k-1$	k	$k+1$...	$n-1$
					...	$n-2$		0	1	...

- 动态规划

$$J(n, k) = (J(n - 1, k) + k) \% n, \text{ if } n > 1,$$
$$J(1, k) = 0$$

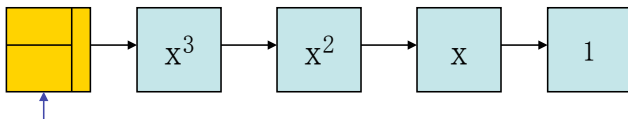
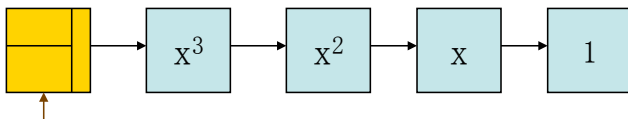
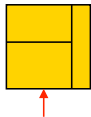
- 我们可以在 $O(n)$ 时间内完成迭代计算

实验三：多项式乘法

- 如果已排序好，我们知道在哪里插入以后的项目（在当前节点之后），搜索花费为 $O(n^2)$
- 举个例子，
 $f(x) = x^3 + x^2 + x + 1, g(x) = x^3 + x^2 + x + 1$, to get $f(x)*g(x)$
- 首先计算 $x^3 * g(x)$, 有 $H \rightarrow x^6 \rightarrow x^5 \rightarrow x^4 \rightarrow x^3$, 总是插入到末尾，进行n次插入
- 然后计算 $x^2 * g(x)$ 并与之前结果相加，
 $H \rightarrow x^6 \rightarrow 2x^5 \rightarrow 2x^4 \rightarrow 2x^3 \rightarrow x^2$, n+1次 添加/插入
- 有 $n + (n + 1) + \dots + (2n - 1) = O(n^2)$

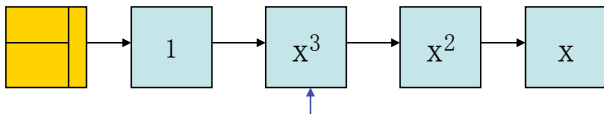
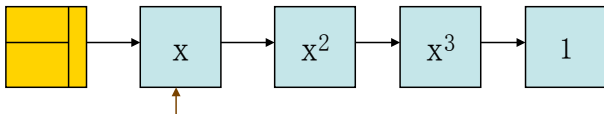
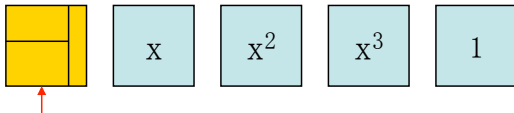
排序

通过记住最后一个节点来连续插入



未排序

搜索必须从头部开始



复杂度分析

- 令 $f(n) = \sum c_i x_i$, $g(n) = \sum d_i x_i$, for $i = 0..n$, 并用 $T(n)$ 表示 $f(n) * g(n)$ 的时间复杂度
- 因为 $f(n) * g(n) = c_n d_n x^{2n} + x^n * (c_n g(n-1) + d_n f(n-1)) + f(n-1) * g(n-1)$
- 实际上, 这完全取决于如何插入第一个 $n + 1$ 项
 - 平凡情况 $T(1) = O(1)$
 - 已排序情况: $T(n) = T(n-1) + O(n)$
 - 未排序情况: $T(n) = T(n-1) + O(n^2)$

实验5

- 给出一棵二叉树的先序（或后序）遍历结果，以及中序遍历结果，如何构造这棵树？假定遍历结果以数组方式输入，请写出相应函数，判断是否存在生成同样遍历结果的树，如果存在，构造这棵树。
- 二叉树的层序遍历。使用队列作为辅助存储，按树的结点的深度，从根开始依次访问所有结点。

Roadmap

1 树的表示

2 二叉树

3 二叉查找树

4 平衡二叉树

5 B树

AVL树

- 自平衡BST的最早形式
- Adelson, Velsky , E. M. Landis, 1962, 一种信息组织的算法
- BST的问题： 可能不平衡
- 解决方法: 保持树平衡，必要时使之重新平衡
- 平衡这个词是什么意思？
 - 根的左子树和右子树都有相同的高度？
 - 不, 这还不够！ 或者说: 每个节点都需要满足此条件？
 - 不, 条件太强了。只需要在 $2^{h+1} - 1$ 范围内满足条件
- 中间: 对于每个节点, $h(L)$ 和 $h(R)$ 的差值 ≤ 1

AVL树

- 不存储节点深度，而存储各子树的高度
- 可以存储实际高度，或者仅仅储存“高度位” (height bits)
 - = - 相同
 - / - 左侧比右侧大1
 - \ - 右侧比左侧大1
- $H(\text{AVL})$ 保证 $\leq \sqrt{2} \log n$ ，搜索自然为 $\log n$

AVL插入结点

- 插入：可能会破坏AVL的特性
- AVL树需要在插入节点后使用“旋转”来使自己再次达到平衡
- 插入可能会使树不平衡
 - 树的哪部分会不平衡？
 - 根节点
 - 插入节点的父节点
 - 插入节点到根节点的路径上
- 我们重新平衡最深的节点

AVL平衡

- 假设我们必须平衡节点a: a的两个子树的高度差必须 ≥ 2
 - 插入节点后平衡
 - 插入节点后子树高度差为2
- 产生不平衡的四种情况
 - 1 对a的左儿子的左子树进行一次插入
 - 2 对a的左儿子的右子树进行一次插入
 - 3 对a的右儿子的左子树进行一次插入
 - 4 对a的右儿子的右子树进行一次插入
- 对称: 1 \sim 4, 2 \sim 3

旋转

- 情形 1 和 4 相对简单
- 情形 2 和 3 较难：
 - 单旋转
 - 双旋转
- 把单旋转当成一次操作
- 双旋转为两次单旋转

单旋转代码

```
Position rotateWithLeftChild (Position k2)
{
    Position k1 = k2->left;

    k2->left = k1->right;
    k1->right = k2;
    k2->height = MAX (height (k2->left),
                      height (k2->right)) + 1;
    k1->height = MAX (height (k1->left),
                     k2->height) + 1;

    return (k1);
}
```

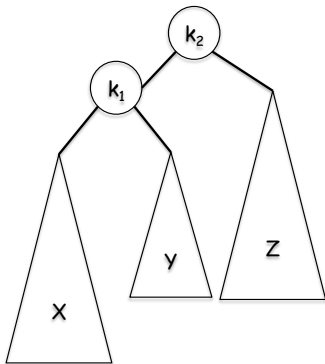
双旋转代码

```
position doubleRotateWithLeftChild (Position k3)
{
    k3->left = rotateWithRightChild (k3->left);

    return (rotateWithLeftChild (k3));
}
```

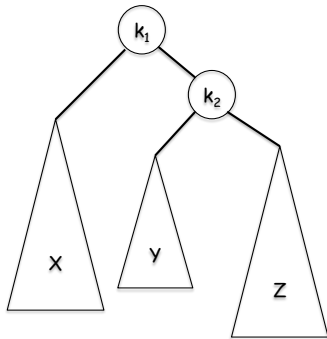
举例：插入子树X

- 插入节点后， k_1 不满足平衡特性导致 k_2 不满足平衡特性
- 插入前高度差必须在1以内，所以 $h(X) = h(Y) + 2$
- 在向X中插入节点后，有 $h(X) = m + 2$ ， $h(Y) = m$ ， $h(Z) = m$
- 插入节点前，有 $h(X) = m + 1$ ， $h(Y) = m$ ， $h(Z) = m$

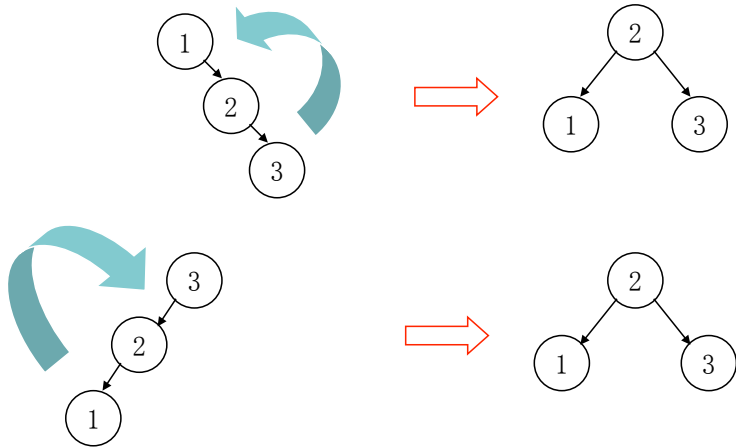


旋转

- 为了重新平衡，我们需要上移 X ，下移 Z
- 让 k_1 成为新的根节点， k_2 保持和 k_1 与 Z 连接
- 但是 Y 成为了 k_2 的左子树
- 这还是二叉树吗？
是的

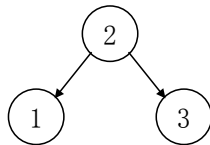
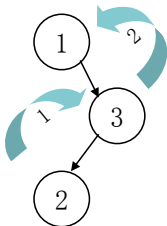
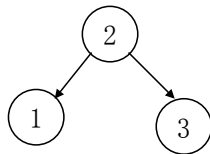
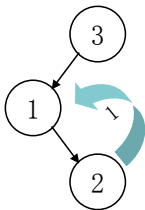


单旋转



双旋转

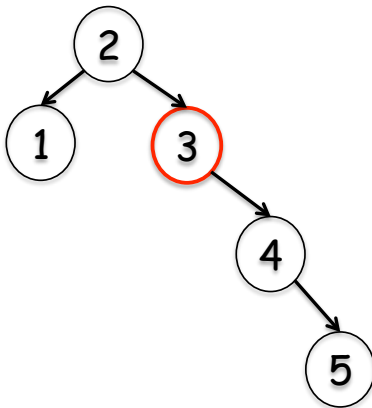
插入顺序: 9, 8, 7, 6, 5, 4, 3, 2, 1
和 9, 7, 8, 4, 6, 5, 3, 1, 2



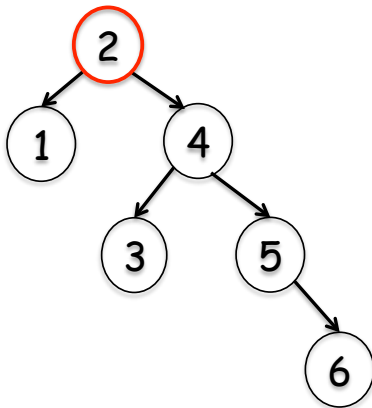
动态生成AVL树



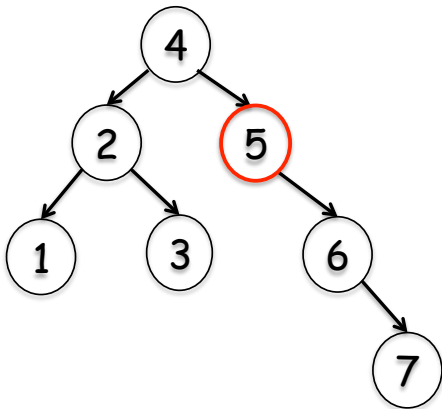
调整



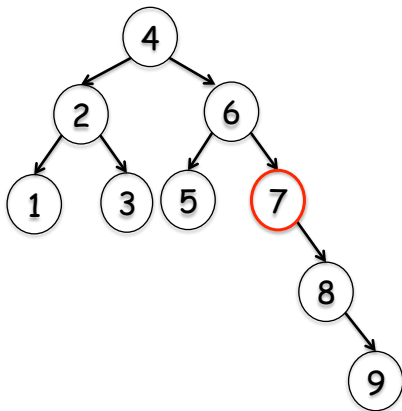
调整



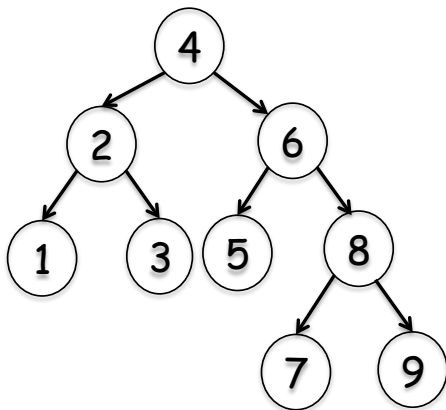
调整



调整



调整



删除

- 删除比插入要复杂
- 进行二叉排序树的删除
- 检查删除后是否平衡
- 如果不平衡，使用旋转来使二叉树达到平衡
- 或者使用懒惰删除法

AVL树的性能

- 最坏情形: $A_h = A_{h-1} + A_{h-2} + \text{root}$, $A_0 = \text{root}$
- 最少节点数: 0 1 2 4 7 12 20
- 是不是眼熟?
- 回忆斐波那契数列: 1 2 3 5 8 13 21
- 一般的: $N_h = N_{h-1} + N_{h-2} + 1$, 改写为:
 $(N_h + 1) = (N_{h-1} + 1) + (N_{h-2} + 1)$
- 使 $F_h = N_h + 1$, 我们有 $F_h = F_{h-1} + F_{h-2}$

AVL树的高度

- 斐波那契 $F_h = O(\phi^h)$, $\phi = \frac{1+\sqrt{5}}{2} \approx 1.618$
- 因此 $N_h = O(\phi^h)$
- N_h 代表 AVL树中的节点数
- 回忆斐波那契数列: 1 2 3 5 8 13 21
- $\log N_h = \log 1.618^h = h \log 1.618 = h * .69$
- $h \approx 1.44 \log N_h$

Roadmap

1 树的表示

2 二叉树

3 二叉查找树

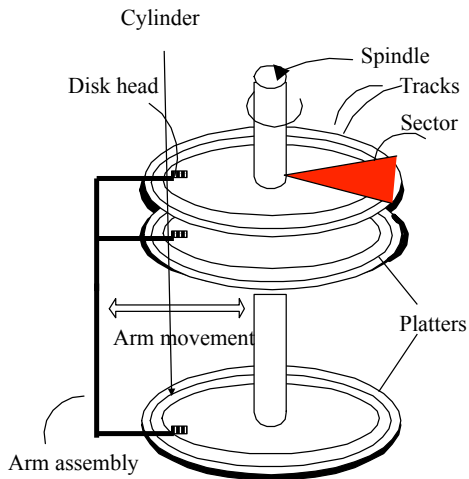
4 平衡二叉树

5 B树

B树的引入

- 二叉排序树在平衡情况下很实用
- 可以在对数(\log)时间内访问(最好情况下)
- 如果 n 变大会如何?
- 如果数据库中有数十亿、数万亿的数据会如何?
- 没有足够的内存, 必须放入二级存储器。规模较大, 存取时间较慢。
- 在复杂性分析中对RAM模型的假设可能不再有效。

硬盘的结构



磁盘与内存

- $7200/\text{M} \sim 1/120\text{s}$ 每转 = 8.3ms 每转
- 如果中途离开, $8.3/2 = 4.15\text{ms}$
- 但是头部也要移动, 所以: $\approx 10\text{ms} = .01\text{s} \sim 100$ 转/秒
- 比较RAM平均存取时间: $10\text{ns} = .00000001\text{s} \sim$ 一亿次访问/秒
- $100,000,000/100 = 1,000,000 \times$ faster

内存与CPU

- 典型的处理器，Pentium/Core/Tegra 3/A7 1GHz以上，每秒十亿次操作？
- 10亿个“循环” ~ 操作每秒
- 比RAM快，但不是那么快
- 与磁盘比较， $1,000,000,000/100 = 10,000,000 \times$ faster
- 一次操作或内存访问不等同于一次磁盘访问
- 磁盘大小增加迅速，但磁盘访问速度却增长缓慢

计算模型

- 在主存中：CPU时间可以用大O表示法给出。
- 在数据库中，时间是由I/O成本决定的
 - 依旧使用大O表示法，但是是对I/O而言
 - 大O经常成一个常数
- 计算的I/O模型将“简单操作”重新定义为“磁盘访问”
- 结果：需要重新设计某些算法和数据结构

数据库中的查找

- 大改进: $\log_2 1\text{MB} = 20$
 - 每个操作将剩下的范围分成两半!
- 但是注意: 主要问题是磁盘访问
- 20 比 2^{20} 好得多, 但是访问磁盘很慢

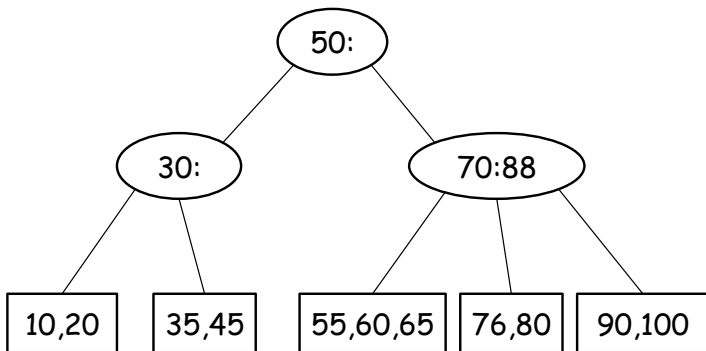
从BST到B树

- 像BSTs，每个节点映射到一块
 - 分支 $\gg 2$
 - 每次访问将剩余范围划分开，比如划分为300个
 - B-trees = BSTs + blocks
 - B+树是B树的一种变体
- 数据只存储在叶子中
 - 叶子形成一个（排序）链表
 - 更好地支持范围查询
- 结果：
 - 短得多的深度会使磁盘读取次数更少
 - 必须在节点内找到元素
 - 用 CPU/RAM 时间等价磁盘读取时间

B树的定义

- 阶为M的B树满足如下特性
 - 根节点要么是一片树叶，要么其儿子树在2和M之间
 - 所有的节点最多有m个子节点
 - 所有内部节点最少有 $m/2$ 个子节点
 - 所有内部节点，有k个儿子则有k-1个键
 - 所有的树叶都在相同的深度上
- 叶子存储数据
- 内部节点的键：分离子树
- 搜索复杂度：amortized(摊还) to $\log_{m/2} n$ to $\log_m n$

B树举例



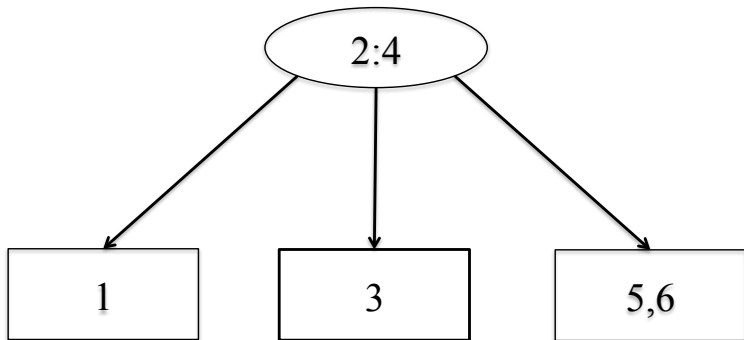
B树的操作

- 搜索 - 类似于BST，从根开始，递归向下，由键指示
- 插入 - 如果插入到一个完整的节点，可能会导致溢出，在这种情况下需要重新平衡，用中间节点分割节点
- 删除 - 删除并重组树，恢复其不变量

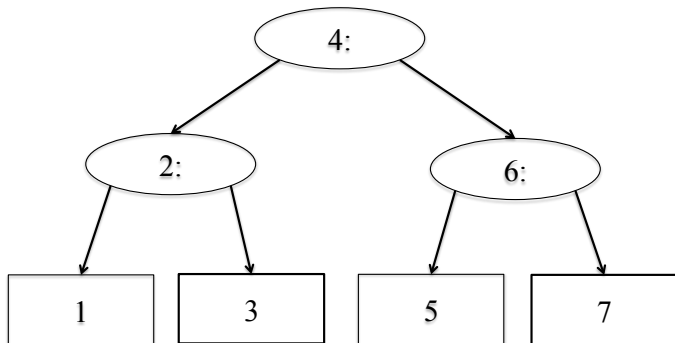
插入

- 首先找到正确的插入位置 t
- 试着插入一个叶子，如果它没有溢出，那么操作完成
- 否则，选择中值将节点分为两类；插入分离的节点可能会增加高度。自底向上的重复，直到满足B-树的性质：不要把所有鸡蛋放进一个篮子里，留一些空间

B树插入



B树插入



B树查找的效率

- 参数：
 - block=4kb
 - integer = 4bytes
 - pointer = 8bytes
- 最大的n满足 $4n + 8(n + 1) \leq 4096$ $n=340$
- 每个节点有 170-340 个键, 取平均值 $(170+340)/2=255$
- 然后：
 - 255 rows at depth = 1
 - $255^2 = 64k$ rows at depth = 2
 - $255^3 = 16M$ rows at depth = 3
 - $255^4 = 4G$ rows at depth = 4

小结

- 树可以表示结点与结点之间的非线性关系，因此在遍历时会丢失部分偏序信息
- 二叉查找树在最好情况下可以达到查找、插入运算的对数时间复杂度
- 平衡二叉树可以通过重新平衡的操作保持树的高度与结点数目之间的对数关系
- 由于磁盘I/O操作耗时，通过建立B树可以在查找时减少对磁盘的访问

实验6

- 4.42 判定二叉树的同构，交换部分（或全部）节点的左右儿子后得到的树称为同构。请写出相应函数，判断两棵给定的树是否同构。
- 4.46 2-d树。使用两个关键词的二叉查找树。完成问题(b)和(c)。