

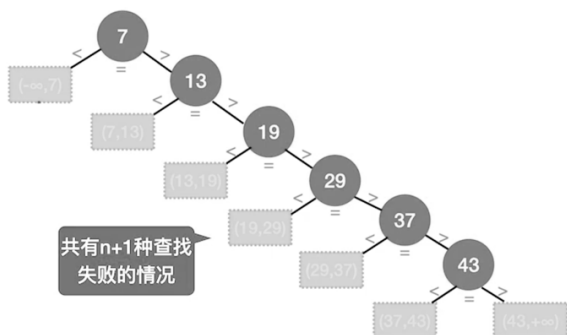
第六章、查找

- 1、掌握静态查找表——顺序表、有序表、索引表的查找算法;理解算法复杂性的分析过程;熟悉算法特点。
- 2、掌握动态查找表——二叉排序树和平衡二叉树的概念、基本操作及其实现。
- 3、理解 B 树的概念和特点。
- 4、熟练掌握哈希查找思想、哈希冲突解决方法、哈希查找性能。

掌握静态查找表——顺序表、有序表、索引表的查找算法;理解算法复杂性的分析过程;熟悉算法特点

动态查找表 { 静态查找表 { ①检索某个特定的数据元素是否在查找表中
②检索某个特定的数据元素的各种属性
③在查找表中插入/删除数据元素

	算法评价
<pre> typedef float KeyType; typedef int KeyType; type char *keytype; #define EQ(a, b) ((a) == (b)) //-----上面是定义 typedef struct { ElemType elem; //数据元素存储空间基址,0号单元留空 int length; //表长度 } SSTable; int Search.Seq(SSTable ST, KeyType key) { //若找到,函数值为该元素在表中的位置,否则为0 ST.elem[0].key = key; //哨兵 for (int i = ST.length; ST.elem[i] != key; --i) //后往前 { return i; } } </pre> <p>顺序查找</p>	<p>ASL { 成功: $1 + 2 + 3 + \dots + n =$ 失败: $n + 1$</p> <p>缺点: 平均查找长度大, 当 n 很大时, 查找效率低</p> <p>优点: 简单、适用面广, 对表的结构没有任何要求</p>
<p>顺序查找的二叉判定树</p>	<p>仅针对有序表</p> <p>失败 ASL:</p>



$$\frac{1 + 2 + \dots + n + n}{n + 1} = \frac{n}{2} + \frac{n}{n + 1}$$

折半查找

```
int Search.Bin(SSTable ST, KeyType key)
{
    low = 1, high = ST.length;
    while (low <= high)
    {
        mid = (low + high) / 2; // mid向下取整
        if (key == ST.elem[mid].key)
        {
            return mid;
        }
        if (key <= ST.elem[mid].key)
        {
            high = mid - 1; // 比中值小, 在前半部分继续查
        }
        else
        {
            low = mid + 1; // 比中值大, 在后半部分继续查
        }
    }
    return 0;
}
```

仅针对有序表和顺序存储结构

每个记录查找概率相同

设有序表长度

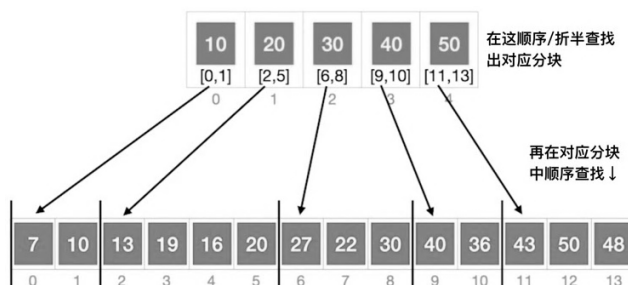
$n = 2^{h-1}$, 成功 ASL

$$= \frac{1}{n} \sum_{i=1}^h j * 2^{j-1}$$

$$= \frac{n+1}{n} \log_2(n+1) - 1$$

$$= \log_2(n+1) - 1 (n > 50)$$

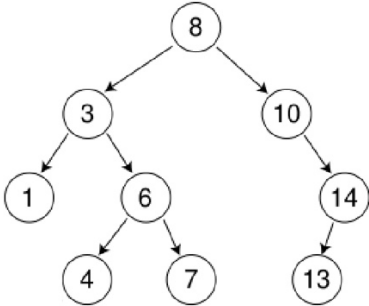
分块 / 索引顺序查找



用于索引表表示的静态查找表

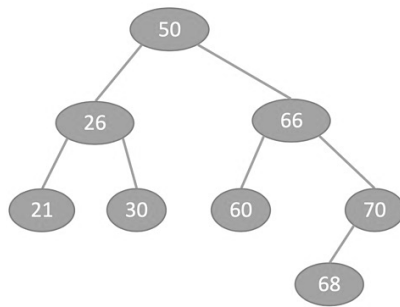
ASL { 成功: 复杂, 按情况计算
失败: 复杂, 不考

二叉排序树/二叉查找树/二叉搜索树 BST

	代码	算法评价
定义	<p>递归表现为左子树<根<右子树</p>  <pre> graph TD 8((8)) --> 3((3)) 8 --> 10((10)) 3 --> 1((1)) 3 --> 6((6)) 6 --> 4((4)) 6 --> 7((7)) 10 --> 14((14)) 14 --> 13((13)) </pre>	
查找	<p>法 1: 非递归实现</p> <p>//在二叉排序树中查找值为 key 的结点</p> <pre> BSTNode *BST_Search(BSTree T,int key){ while(T!=NULL&&key!=T->key){ //若树空或等于根结点值, 则结束循环 if(key<T->key) T=T->lchild; //小于, 则在左子树上查找 else T=T->rchild; //大于, 则在右子树上查找 } return T; } </pre> <p>法二: 递归实现</p> <p>//在二叉排序树中查找值为 key 的结点 (递归实现)</p> <pre> BSTNode *BSTSearch(BSTree T,int key){ if (T==NULL) return NULL; //查找失败 if (key==T->key) return T; //查找成功 else if (key < T->key) return BSTSearch(T->lchild, key); //在左子树中找 else return BSTSearch(T->rchild, key); //在右子树中找 } </pre>	<p>空间复杂度</p> <ul style="list-style-type: none"> 非递归: $O(1)$ 递归: $O(h)$ <p>成功的 ASL:</p> <ul style="list-style-type: none"> 最好: $O(\log_2 n)$ 最坏: $O(n)$ <p>最好时为满二叉 最坏时为单支树</p>
插入	<p>插入的结点一定是叶子结点, 且是查找失败时的查找路径上的最后一个结点的左孩子 or 右孩子</p> <p>//在二叉排序树插入关键字为k的新结点 (递归实现)</p> <pre> int BST_Insert(BSTree &T, int k){ if(T==NULL){ //原树为空, 新插入的结点为根结点 T=(BSTree)malloc(sizeof(BSTNode)); T->key=k; T->lchild=T->rchild=NULL; return 1; //返回1, 插入成功 } else if(k==T->key) //树中存在相同关键字的结点, 插入失败 return 0; else if(k<T->key) //插入到T的左子树 return BST_Insert(T->lchild,k); else //插入到T的右子树 return BST_Insert(T->rchild,k); } </pre>	

构造

按照序列从左到右的顺序，大的往右边放，小的往左边放。比如，{50, 66, 60, 26, 21, 30, 70, 68} 的 BST



```
//按照 str[] 中的关键字序列建立二叉排序树
void Creat_BST(BSTree &T,int str[],int n){
    T=NULL;          //初始时T为空树
    int i=0;
    while(i<n){      //依次将每个关键字插入到二叉排序树中
        BST_Insert(T,str[i]);
        i++;
    }
}
```

不同关键字序列可能得到同款二叉排序树，也可能得到不同款二叉排序树

删除

step1. 搜索找到目标结点

step2. $\left\{ \begin{array}{l} \text{叶子结点} \rightarrow \text{直接删} \\ \text{只有左子树/右子树} \rightarrow \text{子树替代掉它的位置} \\ \text{左右子树都有} \left\{ \begin{array}{l} \text{用直接后继替代, 后继原位执行删除} \\ \text{用直接前驱替代, 前驱原位执行删除} \end{array} \right. \end{array} \right.$

平衡二叉树 AVL

● 相关定义

结点的平衡因子：该结点左子树的高-右子树的高

平衡二叉树：平衡二叉树结点的平衡因子只可能是-1、0 或 1

● 查找效率

AVL 高一层，就多分析一次，所以分析查找效率就是分析 AVL 的高。

平衡二叉树的最大深度 $O(\log_2 n)$ \rightarrow 平均查找长度/时间复杂度 $O(\log_2 n)$

● 高为 h 的平衡二叉树最少有几个结点

此时所有非叶子结点的平衡因子都为 1

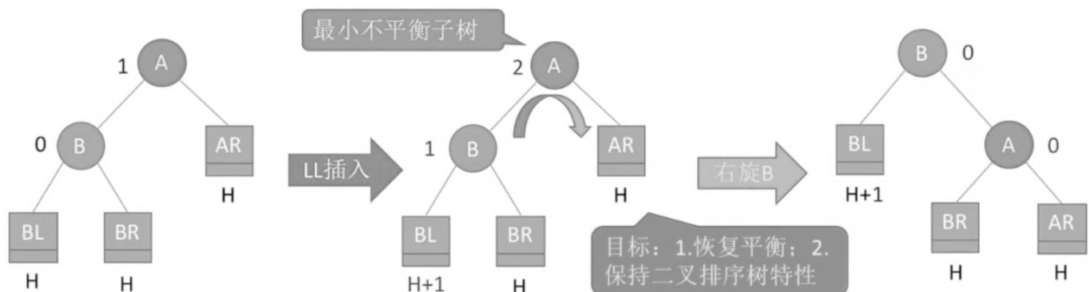
设 n_h 为高度为 h 的 AVL 最少有的结点数，如 $n_0 = 0$, $n_1=1$, $n_2 = 2$

$$n_h = n_{h-1} + n_{h-2} + 1$$

- 插入结点后如何构造一个新的AVL

{	LL → A的左孩子的左子树中插入导致不平衡	
	RR → A的右	右
	LR → A的左	右
	RR → A的右	左

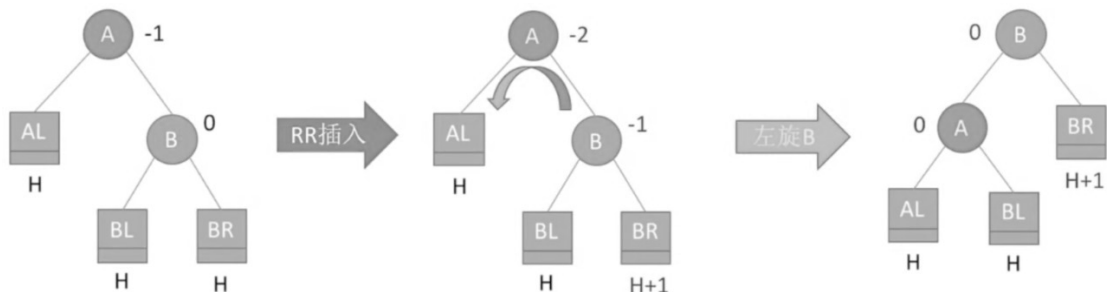
LL



//代码实现

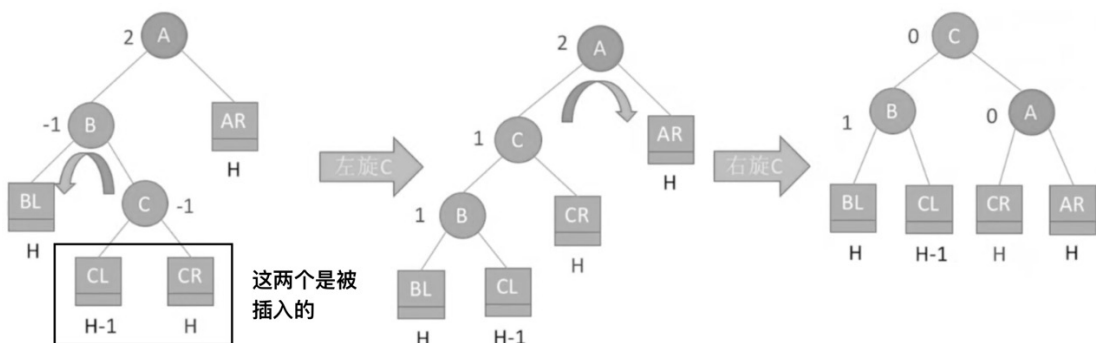
- ① A->lchild=B->rchild
- ② B->rchild=A
- ③ A的父结点->lchild/rchild=B

RR

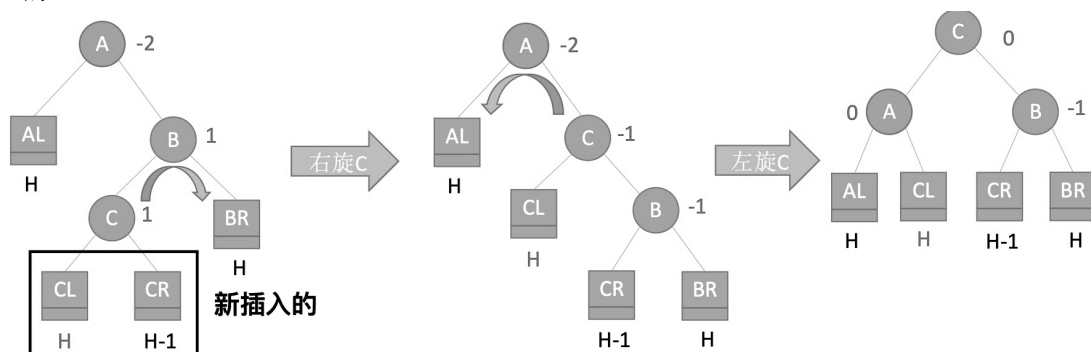


- ① A->rchild=B->lchild
- ② B->lchild=A
- ③ A的父结点->lchild/rchild=B

A 的 LR



A 的 RL



理解 B 树/多路平衡查找树的概念和特点（不要求代码）

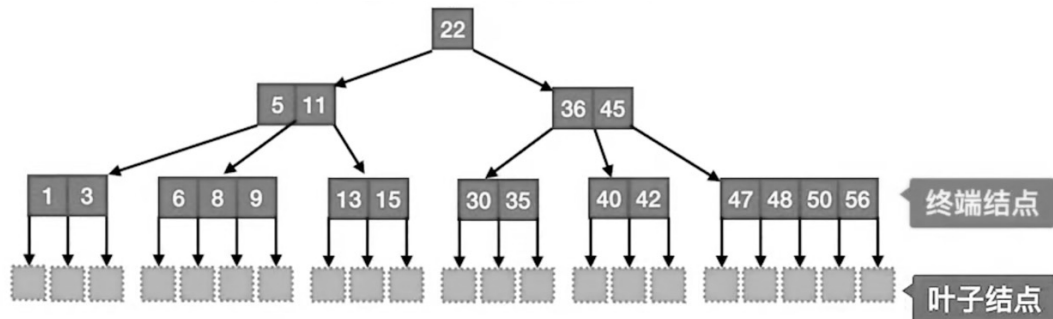
B 树的阶：B 树中所有结点的孩子个数的最大值，通常用 m 表示。

B 树概念

一棵 m 阶 B 树要么是空树，要么是满足如下特性的 m 叉树：

- 根节点的子树 $\in [2, m]$ ，关键字数 $\in [1, m - 1]$ 。
- 其他结点的子树 $\in [\lceil m/2 \rceil, m]$ ，关键字数 $\in (\lceil m/2 \rceil - 1, m - 1]$)
- 对任一结点，其所有子树高度都相同
- 关键字的值：子树 $0 < \text{关键字 } 1 < \text{子树 } 1 < \text{关键字 } 2 < \text{子树 } 2 < \dots$ (类比二叉查找树 左 $<$ 中 $<$ 右)

likethis



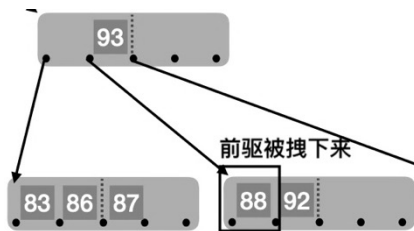
B 树的特点

n 个关键字的 B 树必有 $n+1$ 个叶子结点

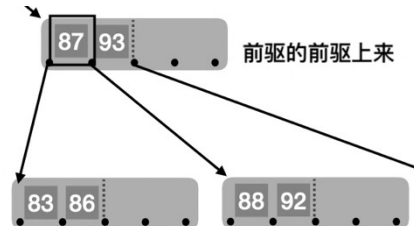
含 n 个关键字的 m 阶 B 树的高 $\begin{cases} \text{最小: } \log_m(n+1) \\ \text{最大: } \log_{\lceil m/2 \rceil} \frac{n+1}{2} + 1 \end{cases}$

推理：

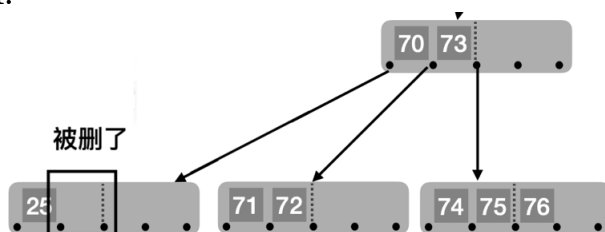
最大高度——=求：多高的 B 树至少有 n 个结点。让各层的分叉尽可能的少，即根节点只有 2 个分叉，其他结点只有 $\lceil m/2 \rceil$ 个分叉，各层结点至少有：第一层 1、第二层 2、第三层 $2\lceil m/2 \rceil \dots$ 第 h 层 $2(\lceil m/2 \rceil)^{h-2}$ 、第 $h+1$ 层共有叶子结点 (失败结点) $2(\lceil m/2 \rceil)^{h-1}$ 个



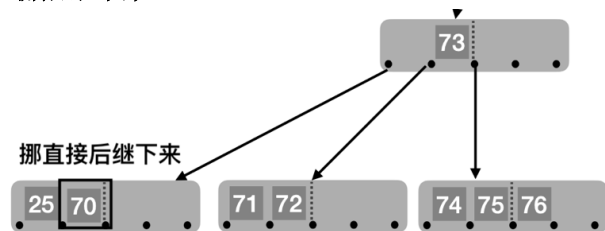
step3. 再把前驱的前驱搬上去



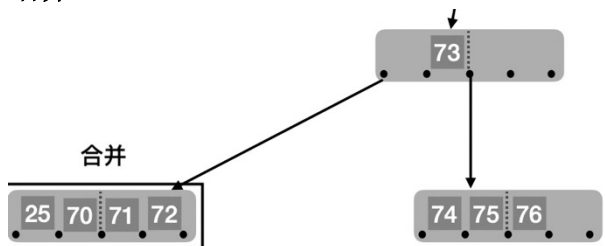
step1.



step2. 挪后继下来

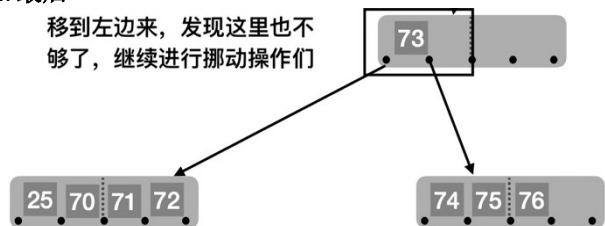


step3. 合并



step4. 最后

移到左边来，发现这里也不
够了，继续进行挪动操作们



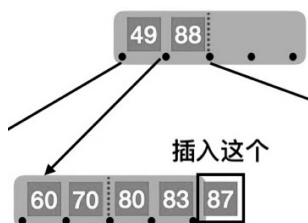
终端结点
↓
删除后低于下限
↓
兄弟不够用
↓
抠父结点然后合并

B 树的插入

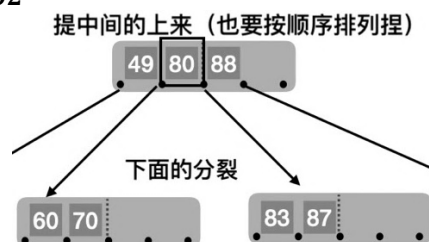
只有新元素一定是插入底层的终端结点，用查找来确定插入位置

插入后超过最大限制？
否 → 直接插
是 → 提中间关键字到父结点

step1



step2

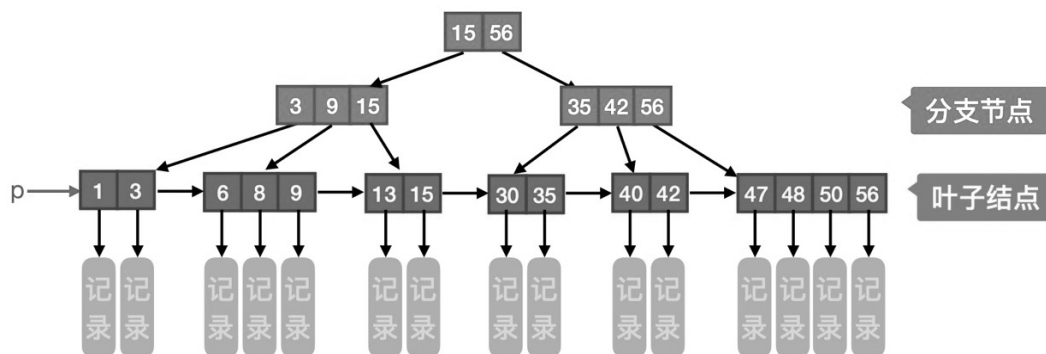


B+树概念

一颗 m 阶的 B+树满足：

- 每个分支结点最多有 m 棵子树 (孩子结点)。
- 非叶根结点至少有两棵子树，其他每个分支结点至少有 $\lceil m/2 \rceil$ 棵子树。
- 结点的子树个数与关键字个数相等。
- 所有叶结点包含全部关键字及指向相应记录的指针，叶结点中将关键字按大小顺序排列，并且相邻叶结点按大小顺序相互链接起来。
- 所有分支结点中仅包含它的各个子结点中关键字的最大值及指向其子结点的指针。

likethis



注：以上是一棵4阶B+树

B+树查找

查找方式		特点
多路查找	<p>查找目标：9</p>	论查找成功与否，最终一定都要走到最下面一层结点
顺序查找	<p>p 指针顺序查找</p>	

B/B- VS B+

	B/B-	B+
来源	二叉查找树的进化——>m 叉查找树	分块查找的进化——>多级分块查找
关键字	n 个关键字对应 n+1 个分叉(子树)	n 个关键字对应 n 个分叉
结点包含信息	所有结点中都包含记录的信息	只有最下层叶子结点才包含记录的信息(可使树更矮)
查找方式	不支持顺序查找。查找成功可能停在任一层结点，查找速度不稳定	支持顺序查找。查找成功或失败都会到达最下一层结点，查找速度稳定
相同点	<ul style="list-style-type: none"> ●除根节点外，最少 $\lceil m/2 \rceil$ 个分叉(确保结点不要太“空”) ●任何一个结点的子树都要一样高(确保“绝对平衡”) 	

哈希查找思想、哈希查找性能

定义们

散列表/哈希表：一种数据结构，数据元素的关键字与其存储地址直接相关

同义词：若不同的关键字通过散列函数映射到同一个值，则称为同义词

冲突：通过散列函数确定的位置已经存放了其他元素，称为冲突

查找长度：需要对比关键字的次数称为查找长度

装填因子 α = 表中记录数/散列表长度

常见散列函数

名称	定义	适用情况	举例
除留余数	散列表表长为 m ，取一个不大于 m 但最接近或等于 m 的质数 p $H(\text{key}) = \text{key} \% p$		散列表表长 15，散列函数 $H(\text{key}) = \text{key} \% 13$
直接定址	$H(\text{key}) = \text{key}$ 或 $H(\text{key}) = a * \text{key} + b$ 其中， a 和 b 是常数。这种方法计算最简单，且不会产生冲突。	关键字的分布基本连续。 (若不连续，则空位多浪费存储空间)	存储同一个班级的学生信息，班内学生学号为 (1120112176~1120112205) $H(\text{key}) = \text{key} - 1120112176$
数字分析	选取数码分布较为均匀的若干位作为散列地址。	关键字的每一位出现的频率不同	手机号码，前三位大多一样。 设计长度为 10000 的散列表，以手机号后四位作为散列地址
平方取中	取关键字的平方值的中间几位作为散列地址。 (取多少位要视实际情况而定)	关键字的每位取值都不够均匀或均小于散列地址所需的位数。	以身份证号为关键字存储学校的学生信息，设计散列函数。学生的生日、地址等大多集中在某几个数字

哈希冲突解决方法

{	开放定址法	{	线性探测再散列
		二次/平方探测再散列	
		伪随机探测再散列	
链地址法			
再哈希法			

		效率
再哈希法	除了原始的散列函数 $H(\text{key})$ 之外，多准备几个散列函数，当散列函数冲突时，用下一个散列函数计算一个新地址，直到不冲突为止	
链地址法	<p>如图，关键字为 {19, 14, 23, 1, 68, 20, 84, 27, 55, 11, 10, 79}，散列函数 $H(\text{key})=\text{key}\%13$</p> <p>$\%13$ 后为 {6, 1, 10, 1, 3, 7, 6, 1, 3, 11, 10, 1}</p>	
定义	<p>指可存放新表项的空闲地址既向它的同义词表项开放，又向它的非同义词表项开放。</p> <p>采用开放定址法时，删除结点不能只将结点置为空，需要做一个删除标记。</p>	
线性探测	<p>发生冲突时，每次往后探测相邻的下一个单元是否为空</p> <p>如 {19, 14, 23, 1, 68, 20, 84, 27, 55, 11, 10, 79}，散列函数 $\text{key}\%13$</p> <p>step1. 插入 1 时，1 位置已有关键字，因此向后搜索</p> <p>step2. 顺序往后一个就找到空了</p>	
开放定址法		
二次 / 平方探测	<p>发生冲突，向右一个，向左一个，向右四个，向左四个... 到头了就从头开始</p> <p>$d_i = 0^2, 1^2, -1^2, 2^2, -2^2, \dots, k^2, -k^2$</p> <p>如图，插入 84。上面的黑字是比较的次数</p>	
伪随机探测	<p>$d_i = 0, 5, 24, 11, \dots$</p>	

