

哈希表：空间换时间。

哈希函数的设计：'键'的分布越均匀越好

(一) 整型

小范围整数直接使用

小范围负整数进行偏移 $\sim 100 - 100 \rightarrow 0 \sim 200$

大整数：取模，ex: 身份证号: 5172219802096612
 $\text{int mod } 10^6 \rightarrow 096612 \rightarrow \text{分布不均}$

解决办法：模一个素数。

单精度：32bit (float)

① sign
 ② exponent (8-bit)
 ③ fraction (23-bit)

64bit (double)

① sign
 ② exponent (11-bit)
 ③ fraction (52-bit)

直接当整型处理

String

$$\text{code} = c \times 26^3 + o \times 26^2 + d \times 26 + e \times 26^0$$

$$= c \cdot B^3 + o \cdot B^2 + d \cdot B + e \cdot B^0$$

$$\text{hash}(\text{code}) = (c \cdot B^3 + o \cdot B^2 + d \cdot B + e) \% M$$

$$= [(c \cdot (c \cdot B + o) \cdot B + d) \cdot B + e] \% M$$

$$= [(((c \% M) \cdot B + o) \% M \cdot B + d) \% M \cdot B + e] \% M$$

$\Rightarrow \text{int hash} = 0$
 for(int i=0; i<S.length(); i++)
 hash = chash * B + S.charAt(i) % M

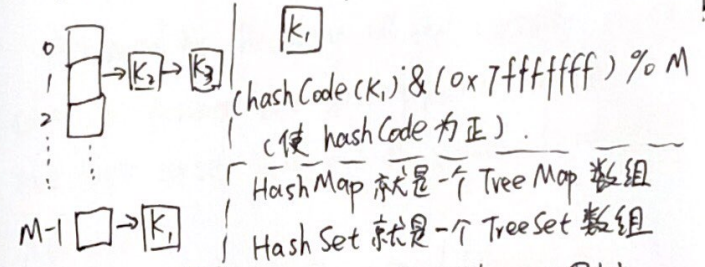
复合类型 Date: year, month, day

$$\text{hash}(\text{date}) = ((\text{date.year} \% M) \cdot B + \text{date.month}) \% M + \text{date.day} \% M$$

int a = 42; ((Integer)a).hashCode() $\rightarrow 42$
 int b = -42; ((Integer)a).hashCode() $\rightarrow -42$

double \Rightarrow (Double).hashCode()

哈希冲突的处理方法: Separate Chaining



当哈希冲突达到一定程度，每个位置从链表 \rightarrow 红黑树 (TreeMap)

链表: $O(N/M)$ 红黑树: $O(\log(N/M))$

平均每个地址承载的元素多过一定程度，扩容。

$$N/M > \text{upperTol}$$

平均每个地址承载的元素少过一定程度，缩容。

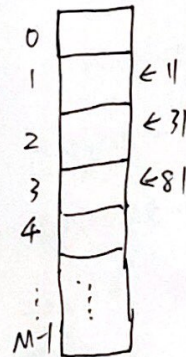
$$N/M < \text{lowerTol}$$

resize 平均复杂度为 $O(1)$

$M \geq 2 \times M$ 不是素数
 在 table 中取素数

哈希冲突更多的处理方法:

开放地址: $\text{if: hash}(x) = x \% 10$



1) 线性探测，遇到冲突找下一个空的索引。

2) 平方探测，+1, +4, +9, +16.

3) 二次哈希，用第2个 + hash2(key).

4) rehashing, 用其他的。

并查集 Union Find

union(p, q)
 isConnected(p, q)

都是 $O(h)$ 复杂度。

Union 是把 P 的 root 指向 Q 的 root.

isConnected(p, q) 判断 P, q 的 root 是否相等。

基于 ~~union~~ 的优化: union 的时候: 选择深度小的树对指向深度高的。

基于 size 的优化: node 少的指向 node 多的。

路径压缩: Path-compression: parent(p) = parent(parent(p)) (在 find 过程中)