

笔记本： data-structure
创建时间： 2018/8/24 10:53
标签： skiplist
URL： <https://blog.csdn.net/u013011841/article/details/39158585>

原 深度理解跳跃链表：一种基于概率选择的平衡树

2014年09月09日 17:29:22 阅读数：1966 标签： 跳跃链表 更多

版权声明：本文为博主原创文章，未经博主允许不得转载。 <https://blog.csdn.net/u013011841/article/details/39158585>

跳跃链表：一种基于概率选择的平衡树

作者：林子

Blog: <http://blog.csdn.net/u013011841>

时间：2014年9月

声明：欢迎指出错误，转载不要去掉出处

跳跃链表简介

二叉树是一种常见的数据结构。它支持包括查找、插入、删除等一系列操作。但它有一个致命的弱点，就是当数据的随机性不够时，会导致其树形结构的不平衡，从而直接影响算法的效率。

跳跃链表 (Skip List) 是1987年才诞生的一种崭新的数据结构，它在进行查找、插入、删除等操作时的期望时间复杂度均为 $O(\log n)$ ，有着近乎替代平衡树的本领。而且最重要的一点，就是它的编程复杂度较同类的AVL树，红黑树等要低得多，这使得其无论是在理解还是在推广性上，都有着十分明显的优势。

跳跃链表的最大优势在于无论是查找、插入和删除都是 $O(\log n)$ ，不过由于跳跃链表的操作是基于概率形成的，那么它操作复杂度大于 $O(\log n)$ 的概率为，可以看出当 n 越大的时候失败的概率越小。

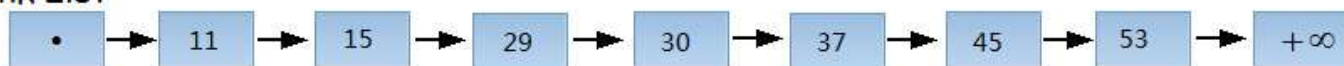
另外跳跃链表的实现也十分简单，在平衡树中是最易实现的一种结构。例如像复杂的红黑树，你很难在不依靠工具书的帮助下实现该算法，但是跳跃链表不一样，你可以很容易在半个小时内就完成其实现。

跳跃链表的空间复杂度的期望为 $O(n)$ ，链表的层数期望为 $O(\log n)$ 。

如何改进普通的链表？

我们先看看一个普通的链表

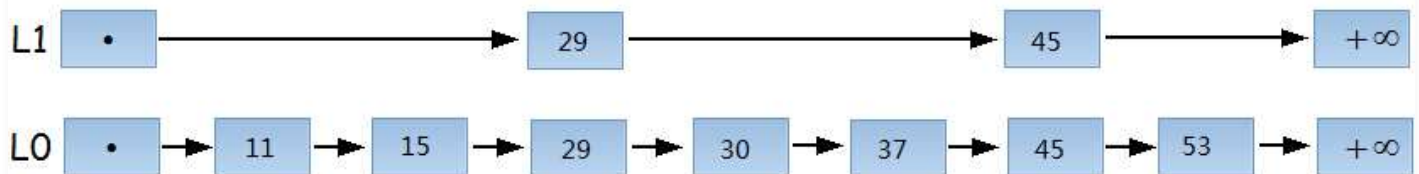
Link List



可以看出查询这个链表 $O(n)$ ，插入和删除也是 $O(n)$ 。因此链表这种结构虽然节省空间，但是效率不高，那有

没有什么办法可以改进呢？

我们可以增加一条链表做为快速通道。这样我们使用均匀分布，从图中可以看出L1层充当L0层的快速通道，底层的结点每隔固定的几个结点出现在上面一层。



我们这里主要以查找操作来介绍，因为插入和删除操作主要的复杂度也是取决于查找，那么两条链表查找的最好的时间复杂度是多少呢？

一次查找操作首先要在上层遍历 $\leq |L_1|$ 次操作,然后在下层遍历 $\leq (L_0/L_1)$ 次操作,至多要经历

$$|L_1| + \frac{|L_0|}{|L_1|}$$
$$= |L_1| + \frac{n}{|L_1|}$$

次操作，其中 $|L_1|$ 为L1的长度,n为L0的长度.

那么最好的时间复杂度,也就怎么设置间隔距离才能使查找次数最少有

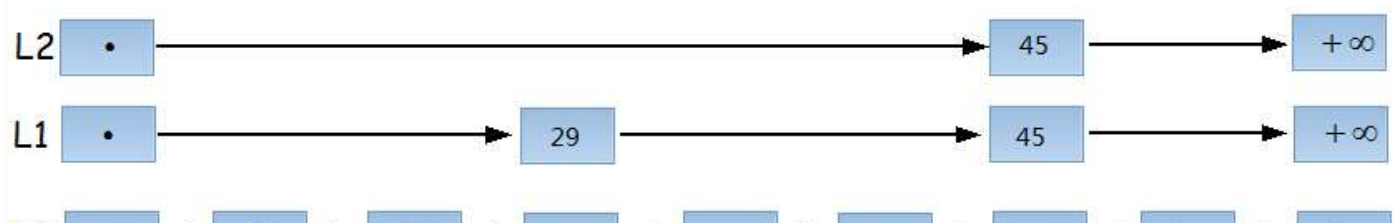
$$\min mize \quad |L_1| + \frac{n}{|L_1|}$$

我们对 $|L_1|$ 的长度求导得

$$|L_1| = \sqrt{n}$$

把上式代入函数,查找次数最小也就是 $2\sqrt{n}$. 这意味着下层每隔 \sqrt{n} 个结点在上层就有一个结点作为快速跑道。

那么三条链表呢





同理那么我们让 $L2/L1=L1/L0$,然后同样列出方程，求导可得 $L2=\sqrt[3]{n}$ ，查找次数为 $3*\sqrt[3]{n}$

第k条链条.....查找次数为 $k\sqrt[k]{n}$

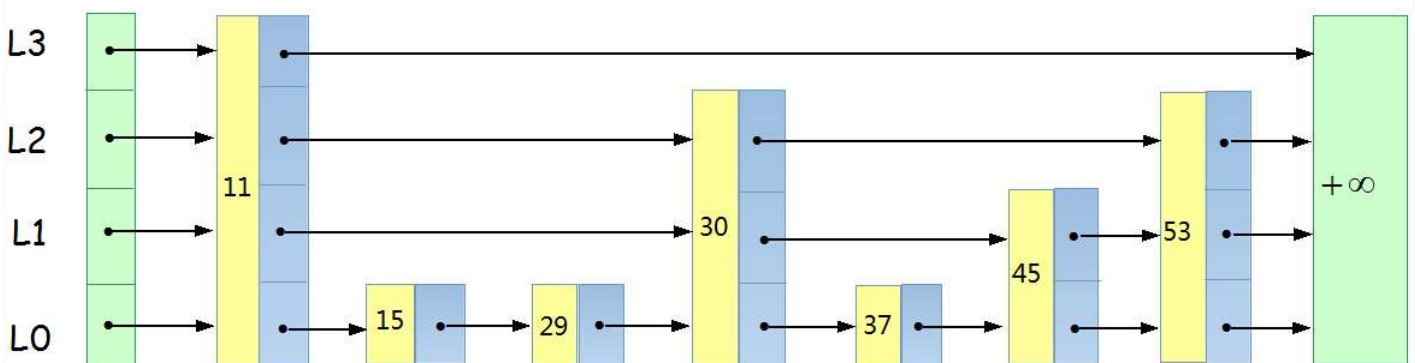
我们这里取 $k=\log n$ ，代入的查找次数为 $2\log n$.

到此为主，我们应该知道了，期望上最好的层数是 $\log n$ 层,而且上下层结点数比为2，这样查找次数常数最小,复杂度保持在 $O(\log n)$ 。

跳跃链表的结构

跳跃表由多条链构成 ($L0, L1, L2, \dots, Lh$)，且满足如下三个条件：

- 每条链必须包含两个特殊元素： $+\infty$ 和 $-\infty$ (其实不需要)
- $L0$ 包含所有的元素，并且所有链中的元素按照升序排列。
- 每条链中的元素集合必须包含于序数较小的链的元素集合。



结点结构源代码

```
1 struct node
2 {
3     int key;
4     struct node *forward[MAXlevel];
5 };
```

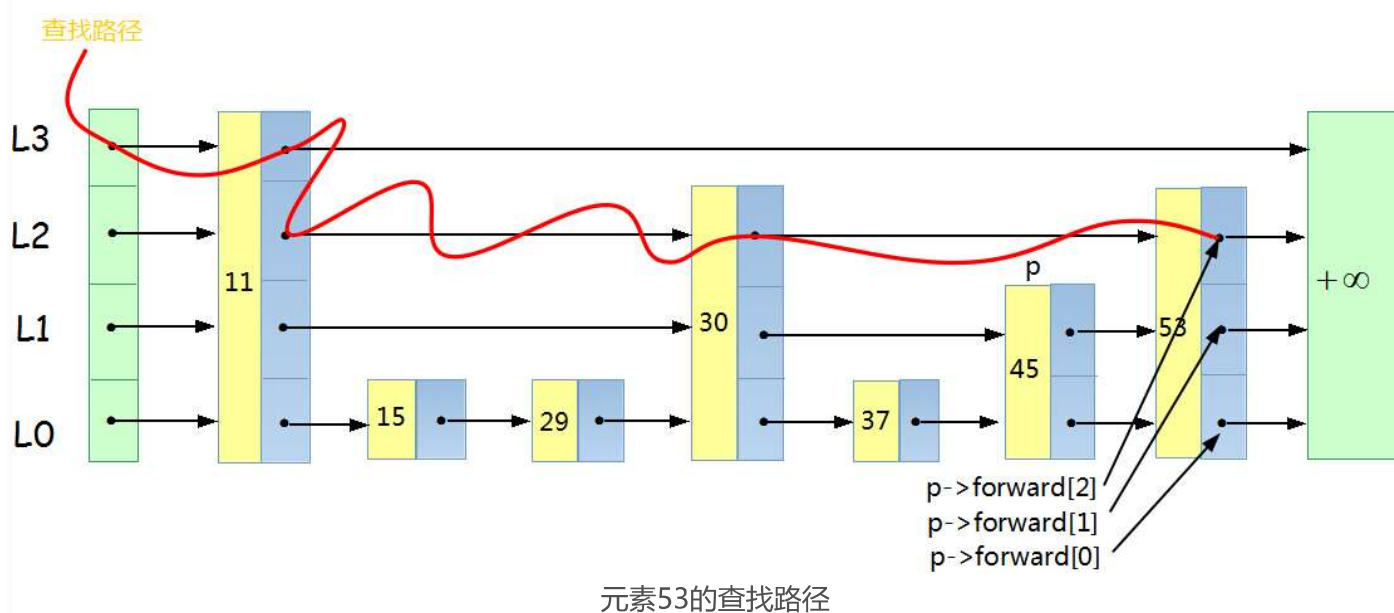
MAXlevel可以是 $\log(n)/\log(2)$

跳跃链表查找操作

目的：在跳跃表中查找一个元素x

在跳跃表中查找一个元素x，按照如下几个步骤进行：

1. 从最上层的链（Lh）的开头开始
2. 假设当前位置为p，它向右指向的节点为q（p与q不一定相邻），且q的值为y。将y与x作比较
 - (1) $x=y$ 输出查询成功及相关信息
 - (2) $x>y$ 从p向右移动到q的位置
 - (3) $x<y$ 从p向下移动一格
3. 如果当前位置在最底层的链中（L0），且还要往下移动的话，则输出查询失败



```
1 struct node* search(struct node*head, int key, int level)
2 {
3     int i;
4     struct node *p;
5     p=head;
6     for(i=level;i>=0;i--)
7         while((p->forward[i]!=0)&&(p->forward[i]->key<key))
8             p=p->forward[i];
9     p==p->forward[0]; //回到0级链，当前p或者空或者指向比搜索关键字小
10    的前一个结点
11    if(p==0) //这时若p为空则推出检索，返回0
12        return 0;
13    else if(p->key==key) //找到，返回成功
14        return p;
15    else
16        return 0; //否则仍然检索失败，返回0
17 }
```

跳跃链表插入操作

目的：向跳跃表中插入一个元素x

首先明确，向跳跃表中插入一个元素，相当于在表中插入一列从S0中某一位置出发向上的连续一段元素。有两个参数需要确定，即插入列的位置以及它的“高度”。

关于插入的位置，我们先利用跳跃表的查找功能，找到比x小的最大的数y。根据跳跃表中所有链均是递增序列的原则，x必然就插在y的后面。

而插入列的“高度”较前者来说显得更加重要，也更加难以确定。由于它的不确定性，使得不同的决策可能会导致截然不同的算法效率。为了使插入数据之后，保持该数据结构进行各种操作均为 $O(\log n)$ 复杂度的性质，我们引入随机化算法（Randomized Algorithms）。

我们定义一个随机决策模块，它的大致内容如下：

产生一个0到1的随机数r如果r小于一个常数p(通常取0.25或0.5)，则执行方案A，否则，执行方案B.

```
1  int randX(int &level)
2  {
3      int i,j,t;
4      t=rand();
5      for(i=0;j=2;i<MAXlevel;i++,j+=j)
6          if(t>RAND_MAX/j)
7              break;
8      if(i>level)
9          level=i;
10     return i;
11 }
```

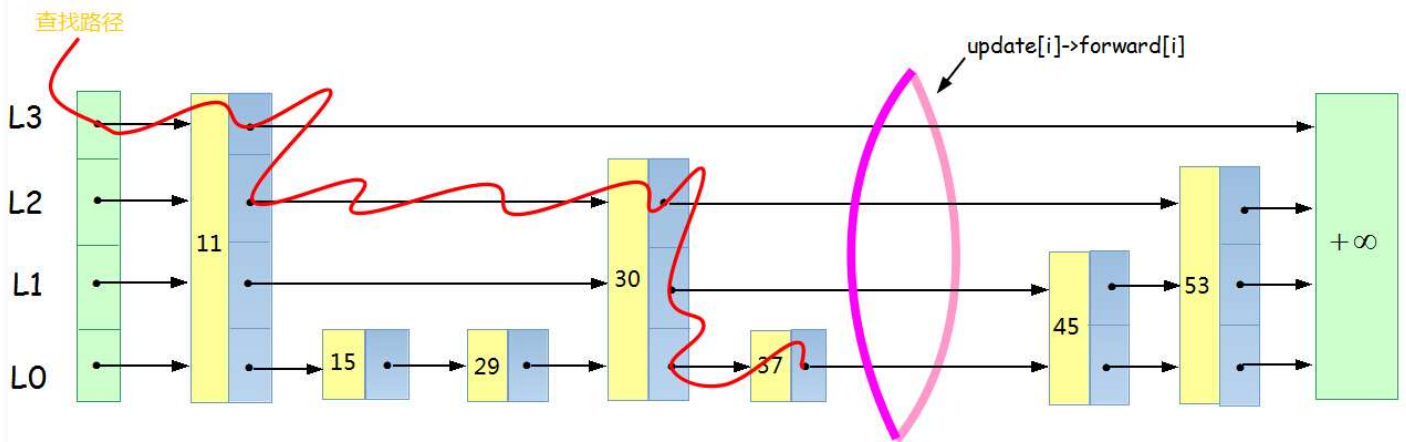
RAND_MAX为随机函数rand()能随机到的最大值,每个结点的高度都由该随机函数决定

初始时列高为1。插入元素时，不停地执行随机决策模块。如果要求执行的是A操作，则将列的高度加1，并且继续反复执行随机决策模块。直到第i次，模块要求执行的是B操作，我们结束决策，并向跳跃表中插入一个高度为i的列。

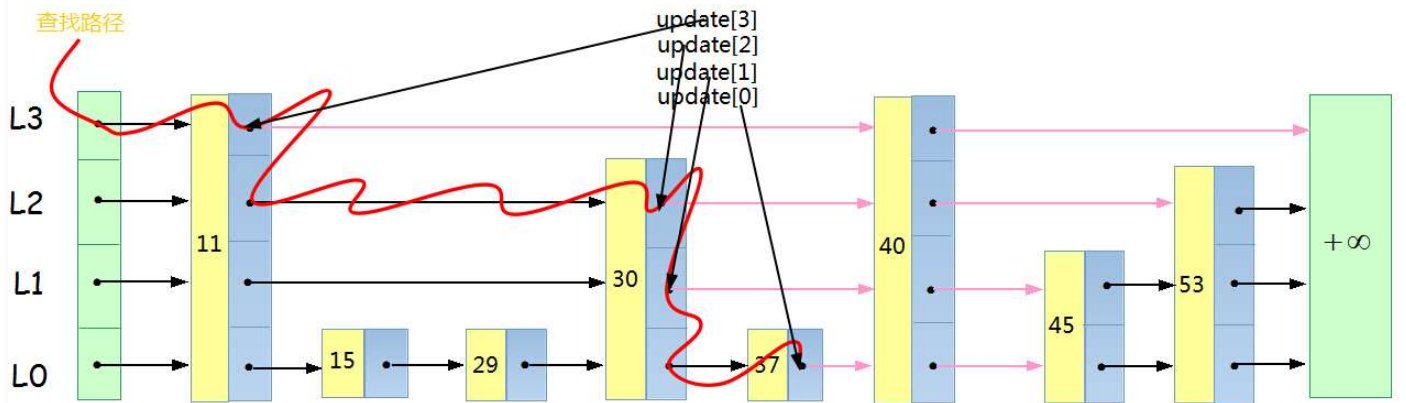
我们来看一个例子：

假设当前我们要插入元素“40”，且在执行了随机决策模块后得到高度为4

步骤一：找到表中比40小的最大的数，确定插入位置



步骤二：插入高度为4的列，并维护跳跃表的结构



紫色的箭头表示更新过的指针

```
1 void insert(struct node*head,int key,int &level)
2 {
3     struct node*p,*update[MAXlevel];
4     int i,newlevel;
5     p=head;
6     newlevel=randX(level);
7     for(i=level;i>=0;i--)
8     {
9         while((p->forward[i]!=0)&&(p->forward[i]->key<key))
10             p=p->forward[i];
11         update[i]=p;//update[i] 记录了搜索过程中在各层中走过的最大的结点位置
12     }// 设置新结点
13     p=new(struct node);
14     p->key=key;
15     for(i=0;i<MAXlevel;i++)
16         p->forward[i]=0;
17     for(i=0;i<=newlevel;i++)// 插入是从最高的newLevel层直至0层
```

```

18     {
19         p->forward[i]=update[i]->forward[i]; //插入到分配到的层数
20         update[i]->forward[i]=p;
21     }
22 }

```

跳跃链表的删除

目的：从跳跃表中删除一个元素x

删除操作分为以下三个步骤：

在跳跃表中查找到这个元素的位置，如果未找到，则退出

将该元素所在整列从表中删除

将多余的“空链”删除

删除的过程即为插入的逆操作

```

1  int deletenode(struct node*head,int &level)
2  {
3      int delkey;
4      struct node*r;
5      cout<<"请输入要删除的数字:";
6      cin>>delkey;
7      r=search(head,delkey,level);
8      if(r)
9      {
10         int i=level;
11         struct node*p,*q;
12         while(i>=0)
13         {
14             p=q=head;
15             while(p!=r&&p!=null)
16             {
17                 q=p;
18                 p=p->forward[i];
19             }
20             if(p)
21             {

```

```

21         {
22             if(i==level&&q=head&&p->forward[i]==0)
23                 level--;
24             else
25                 q->forward[i]=p->forward[i];
26         }
27         i--;
28     }
29     delete r;
30     return 1;
31 }
32 else
33     return 0;
34 }

```

跳跃链表的搜索时间复杂度为 $O(\log n)$

定理： n 个元素的跳跃链表的每一次搜索的时间复杂度有很高的概率为 $O(\log n)$.

高概率：事件 E 以很高的概率发生意味着对于 $a > 1$, 存在一个合适的常数使得事件 E 发生的概率 $\Pr\{E\} > 1 - O(1/n^a)$.

其中 a 是任意选择的一个数，不同的 a 影响搜索时间复杂度的常数，即 $a \cdot O(\log n)$, 这个在后面介绍.

我们要证跳跃链表的时间复杂度，不能只是证明一次搜索的复杂度为 $O(\log n)$, 是要证明全部的搜索都是 $O(\log n)$, 因为这是基于概率的算法，如果光一次有效率并没有多大作用.

我们定义时间 E_i 为某一次搜索失败的概率，那么假设 k 次搜索, 我们先假定失败的概率为 $O(1/n^a)$, 其中至少有一次失败的概率为

$$\begin{aligned}
 & \Pr\{E_1 \cup E_2 \cup \dots \cup E_k\} \\
 & \leq \Pr\{E_1\} + \Pr\{E_2\} + \dots + \Pr\{E_k\} \\
 & = k \times \frac{1}{n^a} \\
 & \underline{\underline{\text{令 } k = n^c \frac{1}{n^{a-c}}}}
 \end{aligned}$$

可以估算出 k 次有一次失败的概率为 $1/n^{(a-c)}$, 那么我们只要让 $a > c+1$ 或者 a 取无穷大，就可以证明每一次搜

索都具有高概率成功。

跳跃链表的层数有高概率为 $O(\log n)$

类似上面的方法，对于 n 个元素，如果有一个层数超过 $O(\log n)$ 就算失败。那么对于某一个元素超过 $c \log n$ 层，即失败的概率为 $\Pr\{E\}$ 。那么对于一次搜索失败的概率为

$$\begin{aligned} & n \times \Pr\{\text{某一结点抬升的层数} \geq c \cdot \log_2 n\} \\ & \leq n \cdot \left(\frac{1}{2}\right)^{c \cdot \log_2 n} \\ & = n \cdot \left(\frac{1}{2^{\log_2 n^c}}\right) \\ & = \frac{n}{n^c} \\ & = \frac{1}{n^{c-1}} \end{aligned}$$

令 $a=c-1$,则只要 $a \geq 1$ 时，就有高概率的可能使得层数为 $O(\log n)$

跳跃链单次查找复杂度大于 $O(\log n)$ 的概率

每完成一次查找，都肯定要从最顶层移动到最下面一层，这每改变一次层数是由概率选择时候的 p 处于扔硬币中的正面决定的。既然上面知道层数高概率为 $c \log n$ 层,那么扔正面的次数为 $c \log n - 1$ 次。

我们假设扔了 $c \log n$ 个正面，超过 $10 \cdot c \log n$ 次是反面，则有

$$\begin{aligned} & \Pr\{\text{扔了 } c \log n \text{ 次正面, 超过 } 10c \log n \text{ 次反面}\} \\ & \leq \binom{10c \log n}{c \log n} \left(\frac{1}{2}\right)^{9c \log n}, \text{ 因为 } \binom{y}{x} \leq \left(e \frac{y}{x}\right)^x \end{aligned}$$

$$\leq e \left(\frac{10c \log n}{c \log n} \right)^{c \log n} \left(\frac{1}{2} \right)^{9 \log n}$$

$$= \frac{(10e)^{c \log n}}{2^{9c \log n}}$$

$$= \frac{2^{\log(10e)c \log n}}{2^{9c \log n}}$$

$$= 2^{[\log(10e)-9]c \log n}$$

$$= \frac{1}{2^{[9-\log(10e)]c \log n}}$$

$$\underline{\underline{\text{令 } \alpha = [9 - \log(10e)]c}} \quad \frac{1}{n^\alpha}$$

因为9-log(10e)中的9是线性增长，要远远大于log(10e)中的对数增长，因此超过10clogn的概率随着10的增长变得越来越小.所以全部的操作都在10logn以内，我们使用k替代10作为常数，即查找次数为klogn，为O(logn).