



数据结构与算法

Data Structure and Algorithm

极夜酱

目录

1	计算复杂性理论	1
1.1	时间/空间复杂度	1
1.2	递推方程	6
1.3	$P=NP?$	10
2	数组	13
2.1	数组	13
2.2	查找算法	14
2.3	数组元素插入与删除	16
3	链表	18
3.1	链表	18
3.2	链表的增删改查	20
3.3	带头结点的链表	24
3.4	倒数第 k 个结点	26
3.5	环形链表	27
3.6	反转链表	30
4	栈	32
4.1	栈	32
4.2	入栈与出栈	35
4.3	最小栈	36
4.4	括号匹配	38
4.5	表达式求值	39
5	队列	46
5.1	队列	46
5.2	循环队列	48
5.3	栈实现队列	50
5.4	队列实现栈	52

5.5	双端队列	54
6	哈希表	56
6.1	哈希表	56
6.2	哈希函数	58
6.3	冲突处理	63
6.4	性能分析	67
7	分治法	68
7.1	分治法	68
7.2	大整数加法	70
8	排序算法	73
8.1	排序算法	73
8.2	冒泡排序	75
8.3	选择排序	80
8.4	插入排序	83
8.5	鸡尾酒排序	86
8.6	归并排序	90
8.7	快速排序	92
8.8	计数排序	96
8.9	桶排序	98
8.10	基数排序	101
8.11	珠排序	103
8.12	猴子排序	105
9	树	106
9.1	树	106
9.2	二叉树	108
9.3	二叉树的遍历	112
9.4	二叉搜索树	117
9.5	哈夫曼树	121
9.6	哈夫曼编码	126

10 图	128
10.1 图	128
10.2 图的表示	130
10.3 图的遍历	135
10.4 连通图	138
10.5 最短路径	142
10.6 最小生成树	150
10.7 拓扑排序	156
11 分治法	160
11.1 最近点对	160
11.2 凸包	162
11.3 芯片测试	166
11.4 大整数乘法	169
11.5 快速幂	174
11.6 矩阵乘法	177
12 排序算法	179
12.1 希尔排序	179
12.2 归并排序	185
12.3 快速排序	189
12.4 堆排序	196
13 贪心算法	205
13.1 贪心算法	205
13.2 活动安排	209
13.3 部分背包	211
13.4 贪心算法局限性	213
14 动态规划	215
14.1 动态规划	215
14.2 硬币找零	220
14.3 路径问题	224
14.4 跳跃游戏	226

14.5 0-1 背包	228
15 模式匹配	231
15.1 BF	231
15.2 Sunday	234
15.3 RK	237
15.4 BM	239
15.5 KMP	242
16 树	246
16.1 AVL 树	246
16.2 红黑树	261
16.3 B 树	282
16.4 B+ 树	285
16.5 并查集	289

Chapter 1 计算复杂性理论

1.1 时间/空间复杂度

1.1.1 算法 (Algorithm)

算法是解决问题的一种方法，由一系列的步骤组成。算法有 5 个特点：

1. 有穷性 (finiteness)：算法必须在有限个步骤后终止。
2. 确定性 (definiteness)：算法的每个步骤必须有确切的定义。
3. 输入项 (input)：一个算法有 0 个或多个输入。
4. 输出项 (output)：一个算法有 1 个或多个输出，没有输出的算法是毫无意义的。
5. 可行性 (effectiveness)：算法中执行的任何步骤都可以被分解为基本的可执行操作。

1.1.2 时间复杂度 (Time Complexity)

算法有高效的，也有拙劣的。衡量算法好坏的标准有时间复杂度和空间复杂度。

想象一个场景：老板让小灰和大黄完成一个需求，两人都完成并交付了各自的代码，代码的功能是一样的。但是，大黄的代码运行一次要花 100ms，占用 5MB 内存；小灰的代码运行一次要花 100s，占用 500MB 内存。

“小灰，收拾东西走人，明天不用来上班了！”

小灰虽然也按照老板的要求实现了功能，但他的代码存在很严重的问题：运行时间长、占用空间大。

算法的时间复杂度是指算法中基本运算的执行次数，其中基本运算指的是加减法、交换、比较等操作。

算法的效率应该取决于算法本身，与机器无关。分析算法运行效率时应该考虑的是运行时间与输入规模之间的关系。通过计算算法中基本运算的执行次数，可以得到一个关于输入规模 n 的函数。

时间复杂度一般采用大 O 表示法，表示算法的运行时间与输入规模之间的增长关系。常见的时间复杂度包括 $O(1)$ 、 $O(\log n)$ 、 $O(n)$ 、 $O(n \log n)$ 、 $O(n^2)$ 、 $O(2^n)$ 、 $O(n!)$ 等。

时间复杂度需要满足以下原则：

1. 如果运行时间是常量级，则用 $O(1)$ 表示。
2. 只保留时间函数中的最高阶项。
3. 如果最高阶项存在，则省去最高阶项前的系数。

大 O 符号

大 O 符号用于表示时间复杂度的上界（最坏情况），即算法的阶不会超过大 O 符号中的函数 $f(n)$ 。

$$0 \leq f(n) \leq cg(n) \quad (1.1)$$

$$\begin{aligned} f(n) &= n^2 + n \\ &= O(n^2) \\ &= O(n^3) \end{aligned}$$

大 Ω 符号

大 Ω 符号用于表示时间复杂度的下界（最好情况），即算法的阶不会低于大 Ω 符号中的函数 $f(n)$ 。

$$0 \leq cg(n) \leq f(n) \quad (1.2)$$

$$\begin{aligned} f(n) &= n^2 + n \\ &= \Omega(n^2) \\ &= \Omega(100n) \end{aligned}$$

小 o 符号

小 o 符号用于表示时间复杂度的上界，即算法的阶一定低于小 o 符号中的函数 $f(n)$ 。

$$0 \leq f(n) < cg(n) \quad (1.3)$$

$$\begin{aligned} f(n) &= n^2 + n \\ &= o(n^3) \end{aligned}$$

小 ω 符号

小 ω 符号用于表示时间复杂度的下界，即算法的阶一定高于小 ω 符号中的函数 $f(n)$ 。

$$0 \leq cg(n) < f(n) \quad (1.4)$$

$$\begin{aligned} f(n) &= n^2 + n \\ &= \omega(n^2) \end{aligned}$$

Θ 符号

若 $f(n) = O(g(n))$ 且 $f(n) = \Omega(g(n))$ ，则称 $f(n)$ 的阶与 $g(n)$ 的阶相等：

$$f(n) = \Theta(g(n)) \quad (1.5)$$

$$f(n) = n^2 + n$$

$$g(n) = 100n^2$$

$$= \Theta(n^2)$$

百钱买百鸡

公鸡 5 文钱 1 只，母鸡 3 文钱 1 只，小鸡 1 文钱 3 只，如果用 100 文钱买 100 只鸡，那么公鸡、母鸡和小鸡各应该买多少只？

```

1 void buy(int n, int money) {
2     for (int x = 0; x <= n / 5; x++) {
3         for (int y = 0; y <= n / 3; y++) {
4             int z = n - x - y;
5             if (z > 0 && z % 3 == 0 && 5*x + 3*y + z/3 == money) {
6                 printf("x = %d, y = %d, z = %d\n", x, y, z);
7             }
8         }
9     }
10 }
```

1.1.3 空间复杂度 (Space Complexity)

算法占用的内存空间自然是越小越好，空间复杂度是对一个算法在运行过程中临时占用存储空间大小的量度，它同样使用了大 O 表示法。

正所谓鱼和熊掌不可兼得，很多时候不得不在时间复杂度和空间复杂度之间进行取舍。绝大多数时候，时间复杂度更为重要，宁可多分配一些内存空间，也要提升程序的执行速度。

1.1.4 均摊时间复杂度 (Amortized Time Complexity)

均摊时间复杂度用于分析一组操作中，虽然某个操作的时间复杂度很高，但是经过若干次操作后，这组操作的平均时间复杂度较低的情况。

动态数组

```
1 public void add(T element) {  
2     if (size == capacity) {  
3         capacity *= 2;  
4         T[] newArr = (T[]) new Object[capacity];  
5         System.arraycopy(arr, 0, newArr, 0, size);  
6         arr = newArr;  
7     }  
8     arr[size++] = element;  
9 }
```

这段代码实现了往数组中添加数据的功能。在数组没有满的情况下，直接将数据添加到数组末尾，时间复杂度为 $O(1)$ 。

但是当数组已满时，需要对数组创建一个更大的数组，将原数组中的数据复制到新数组中，然后再将新数据添加到数组的末尾，这个过程的时间复杂度为 $O(n)$ 。

假设数组的容量为 n ，每执行 n 次 `add()` 操作，才会进行一次扩容。那么，可以将这一次的 $O(n)$ 的时间均摊到前面的 n 次操作中，这样 `add()` 操作的均摊时间复杂度就是 $O(1)$ 。

1.2 递推方程

1.2.1 递推 (Recurrence)

递归算法无法直接根据语句的执行次数计算出时间复杂度，但是可以通过递推方程迭代展开进行求解。

求和

```
1 int sum(int n) {  
2     if (n <= 0) {  
3         return 0;  
4     }  
5     return sum(n - 1) + n;  
6 }
```

$$T(n) = \begin{cases} 1 & n = 0 \\ T(n-1) + 1 & n \geq 1 \end{cases} \quad (1.6)$$

$$\begin{aligned} T(n) &= T(n-1) + 1 \\ &= [T(n-2) + 1] + 1 \\ &= [[T(n-3) + 1] + 1] + 1 \\ &= \dots \\ &= T(n-k) + k \end{aligned}$$

$$n - k = 0 \Rightarrow n = k$$

$$\begin{aligned} T(n) &= T(n-n) + n \\ &= T(0) + n \\ &= 1 + n \\ &= O(n) \end{aligned}$$

插入排序

```
1 void insert_sort(int arr[], int n) {  
2     if (n <= 1) {  
3         return;  
4     }  
5     insert_sort(arr, n-1);  
6     int last = arr[n-1];  
7     int j = n - 2;  
8     while (j >= 0 && arr[j] > last) {  
9         arr[j+1] = arr[j];  
10        j--;  
11    }  
12    arr[j+1] = last;  
13 }
```

$$T(n) = \begin{cases} 1 & n \leq 1 \\ T(n-1) + n & n \geq 1 \end{cases} \quad (1.7)$$

$$\begin{aligned} T(n) &= T(n-1) + n \\ &= [T(n-2) + n] + n \\ &= [[T(n-3) + n] + n] + n \\ &= \dots \\ &= T(n-k) + nk \end{aligned}$$

$$n - k = 0 \Rightarrow n = k$$

$$\begin{aligned} T(n) &= T(n-n) + n^2 \\ &= T(0) + n^2 \\ &= 1 + n^2 \\ &= O(n^2) \end{aligned}$$

汉诺塔

有三根柱子 A、B、C，A 柱子上从下到上套有 n 个圆盘，要求将 A 柱子上的圆盘移动到 C 柱子上。每次只能移动一个圆盘，且大圆盘始终不能叠在小圆盘上面。

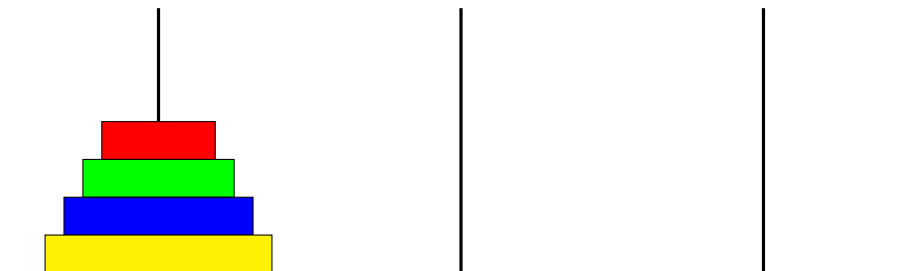


图 1.1: 汉诺塔

递归算法求解汉诺塔问题：

1. 将前 $n - 1$ 个圆盘从 A 柱借助于 C 柱搬到 B 柱。
2. 将最后一个圆盘直接从 A 柱搬到 C 柱。
3. 将 $n - 1$ 个圆盘从 B 柱借助于 A 柱搬到 C 柱。

```
1 def hanoi(n, A, B, C):  
2     if n == 1  
3         move(1, A, C)  
4     else  
5         hanoi(n-1, A, C, B)  
6         move(n, A, C)  
7         hanoi(n-1, B, A, C)
```

$$T(n) = \begin{cases} 1 & n = 1 \\ 2T(n-1) + 1 & n > 1 \end{cases} \quad (1.8)$$

$$\begin{aligned}
T(n) &= 2T(n-1) + 1 \\
&= 2[2T(n-2) + 1] + 1 \\
&= 2[2[2T(n-3) + 1] + 1] + 1 \\
&= \dots \\
&= 2^k * T(n-k) + 2^k - 1
\end{aligned}$$

$$n - k = 1 \Rightarrow n = k + 1$$

$$\begin{aligned}
T(n) &= 2^{n-1} * T(1) + 2^{n-1} - 1 \\
&= 2^{n-1} + 2^{n-1} - 1 \\
&= 2^n - 1 \\
&= O(2^n)
\end{aligned}$$

假设每次移动花费 1 秒，解决一个 64 层的汉诺塔问题大约需要 5800 亿年。

吓得我抱起了

抱着抱着抱着我的小鲤鱼的我的我的我



1.3 P=NP?

1.3.1 旅行商问题 (Traveling Salesman Problem)

小灰最近在工作中遇到了一个棘手的问题。公司正在开发一个物流项目，其中一个需求是为快递员自动规划送快递的路线。

有一个快递员，要分别给三家顾客送快递，他自己到达每个顾客家的路程各不相同，每个顾客之间的路程也各不相同。那么想要把快递依次送达这三家，并最终回到起点，哪一条路线所走的总距离是最短的呢？

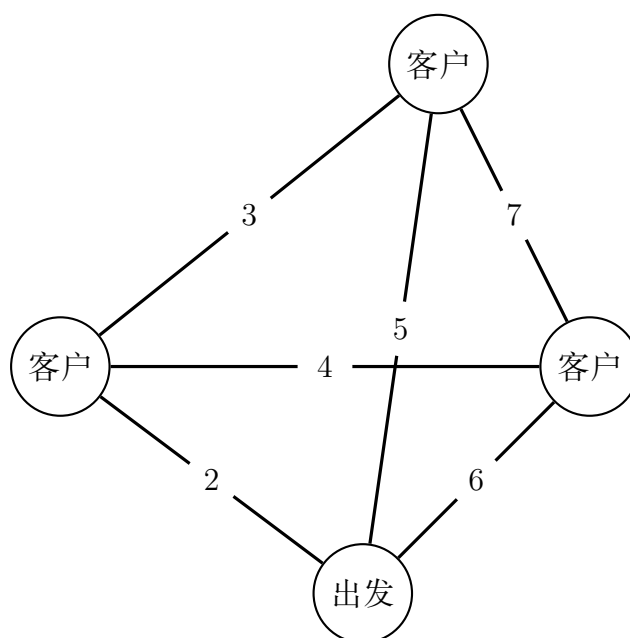


图 1.2: 快递客户路线

为了寻求最优路线，小灰研究了好久，可惜还是没有找到一个高效率的解决方案。不只是小灰，当前的计算机科学家们也没有找到一个行之有效的优化方案，这是典型的旅行商问题。

有一个商品推销员，要去若干个城市推销商品。该推销员从一个城市出发，需要经过所有城市后，回到出发地。每个城市之间都有道路连通，且距离各不相同，推销员应该如何选择路线，使得总行程最短呢？

这个问题看起来很简单，却很难找到一个真正高效的求解算法。其中最容易想到

的，是使用穷举法把所有可能的路线穷举出来，计算出每一条路线的总行程。通过排列组合，从所有路线中找出总行程最短的路线。显然，这个方法的时间复杂度是 $O(n!)$ ，随着城市数量的增长，花费的运算时间简直不可想象！

后来，人们想出了许多相对优化的解决方案，比如动态规划法和分枝定界法等。但是，这些算法的时间复杂度仍然是指数级的，并没有让性能问题得到根本的解决。

像这样的问题有很多，旅行商问题仅仅是其中的一例。对于这类问题统称为 NP 问题。

1.3.2 P=NP?

算法的设计与分析在计算机科学领域有着重要的应用背景。1966 ~ 2005 年期间，Turing 奖获奖 50 人，其中 10 人以算法设计，7 人以计算理论、自动机和复杂性研究领域的杰出贡献获奖。计算复杂性理论中的 P=NP? 问题是世界七大数学难题之一。

一些常见的算法的时间复杂度，例如二分查找法 $O(\log n)$ 、归并排序 $O(n \log n)$ 、Floyd 最短路径 $O(n^3)$ 等，都可以用 $O(n^k)$ 表示。这些算法都是多项式时间算法，即能在多项式时间内解决问题。这类问题被称为 P 问题 (Polynomial)。

人们常说，能用钱解决的问题都不是问题。在计算机科学家眼中，能用多项式时间解决的问题都不是问题。

然而，世间还存在许多变态的问题，是无法（至少是暂时无法）在多项式时间内解决的，比如一些算法的时间复杂度是 $O(2^n)$ ，甚至 $O(n!)$ 。随着问题规模 n 的增长，计算量的增长速度是非常恐怖的。这类问题被称为 NP 问题 (Non-deterministic Polynomial)，意思是“不确定是否能在多项式时间内解决”。

有些科学家认为，所有的 NP 问题终究都可以在多项式时间内解决，只是我们暂时还没有找到方法；也有些科学家认为，某些 NP 问题永远无法在多项式时间内

解决。这个业界争论用 $P=NP?$ 这个公式来表达。

这还不算完，在所有的 NP 问题当中，还存在着一些大 BOSS，被称为 NPC 问题。

1.3.3 NPC

这里所说的 NPC 问题可不是游戏当中的 NPC。要想理解 NPC 问题，需要先了解归约的概念。

归约 (reduction) 可以简单理解成问题之间的转化。例如问题是一个一元一次方程的求解问题 $Q : 3x + 6 = 12$ ，这个问题可以转化成一个一元二次方程 $Q' : 0x^2 + 3x + 6 = 12$ 。

问题 Q 并不比问题 Q' 难解决，只要有办法解决 Q' ，就一定能够解决 Q 。对于这种情况，可以说问题 Q 归约于问题 Q' 。

同时，这种归约可以逐级传递，比如问题 A 归约于问题 B，问题 B 归约于问题 C，问题 C 归约于问题 D，那么可以说问题 A 归约于问题 D。

在 NP 问题之间，也可以存在归约关系。把众多的 NP 问题层层归约，必定会得到一个或多个终极问题，这些归约的终点就是所谓的 NPC 问题 (NP-Complete)。旅行商问题被科学家证明属于 NPC 问题。

俗话说擒贼先擒王，只要有朝一日，我们能够找到 NPC 问题的多项式时间算法，就能够解决掉所有的 NP 问题！但遗憾的是，至今还没有人能够找到可行的方法，很多人认为这个问题是无解的。

回到最初的快递路线规划问题，既然是工程问题，与其钻牛角尖寻求最优解，不如用小得多的代价寻求次优解。最简单的办法是使用贪心算法，先选择距离起点最近的地点 A，再选择距离 A 最近的地点 B，以此类推，每一步都保证局部最优。这样规划出的路线未必是全局最优，但平均情况下也不会比最优方案差多少。

Chapter 2 数组

2.1 数组

2.1.1 数组 (Array)

数组是数据结构中最简单的结构，很多编程语言都内置数组。数组是有限个相同类型的变量所组成的集合，数组中的每一个变量被称为元素。

创建数组时会在内存中划分出一块连续的内存，将数据按顺序进行存储，数组中的每一个元素都有着自己的下标 (index)，下标从 0 开始一直到数组长度-1。因为数组在存储数据时是按顺序存储的，存储数据的内存也是连续的。

对于数组来说，读取元素是最简单的操作。由于数组在内存中顺序存储，所以只要给出数组的下标，就可以读取到对应位置的元素。像这种根据下标读取元素的方式叫作随机读取。但是需要注意的是，数组的下标范围必须在 0 到数组长度-1 之内，否则会出现数组越界。数组读取元素的时间复杂度是 $O(1)$ 。

数组拥有非常高效的随机访问能力，只要给出下标，就可以用常量时间找到对应元素。有一种高效查找元素的算法叫作二分查找，就是利用了数组的这个优势。

数组的劣势体现在插入和删除元素方面。由于数组元素连续紧密地存储在内存中，插入、删除元素都会导致大量元素被迫移动，影响效率。总的来说，数组所适合的是读操作多、写操作少的场景。

更新数组元素

```
1 arr = [3, 1, 2, 5, 4, 9, 7, 2]
2 arr[5] = 10
3 print(arr[5])
```

2.2 查找算法

2.2.1 顺序查找 (Sequence Search)

顺序查找也称线性查找，是一种按照序列原有顺序进行遍历比较查询的基本查找算法。

对于任意一个序列以及一个需要查找的元素（关键字），将关键字与序列中元素依次比较，直到找出与给定关键字相同的元素，或者将序列中的元素与其都比较完为止。若某个元素的值与关键字相等，则查找成功；如果直到最后一个元素，元素的值和关键字比较都不等时，则查找不成功。

最好的情况就是在第一个位置就找到，算法时间复杂度为 $O(1)$ 。

最坏情况是关键字不存在，需要进行 n 次比较，时间复杂度为 $O(n)$ 。

平均查找次数为 $(n + 1)/2$ ，平均时间复杂度为 $O(n)$ 。

顺序查找

```
1 int sequence_search(int *arr, int n, int key) {  
2     for (int i = 0; i < n; i++) {  
3         if (arr[i] == key) {  
4             return i;  
5         }  
6     }  
7     return -1;  
8 }
```

2.2.2 二分查找 (Binary Search)

二分查找法也称折半查找，是一种效率较高的查找方法。折半查找要求线性表必须采用顺序存储结构，而且表中元素按关键字有序排列。

算法思想是假设表中元素是按升序排列，将表中间位置的关键字与查找关键字比较，如果两者相等，则查找成功；否则利用中间位置记录将表分成前、后两个子表，如果中间位置的关键字大于查找关键字，则进一步查找前一子表，否则进一步查找后一子表。重复以上过程，直到找到满足条件的记录，使查找成功，或直到子表不存在为止，此时查找不成功。

二分查找法的时间复杂度为 $O(\log n)$ 。

二分查找

```
1 int binary_search(int *arr, int n, int key) {
2     int start = 0;
3     int end = n - 1;
4     while (start <= end) {
5         int mid = (start + end) / 2;
6         if (arr[mid] == key) {
7             return mid;
8         } else if (arr[mid] < key) {
9             start = mid + 1;
10        } else {
11            end = mid - 1;
12        }
13    }
14    return -1;
15 }
```

2.3 数组元素插入与删除

2.3.1 插入元素

在数组中插入元素存在 3 种情况：

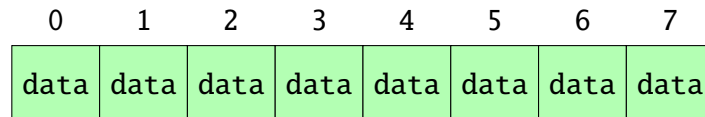


图 2.1: 数组

尾部插入

直接把插入的元素放在数组尾部的空闲位置即可。

中间插入

首先把插入位置及后面的元素向后移动，腾出位置，再把要插入的元素放入该位置上。

插入元素

```
1 int insert(int *arr, int n, int index, int val) {
2     if(index < 0 || index >= n) {
3         return n;
4     }
5     for(int i = n - 1; i >= index; i--) {
6         arr[i+1] = arr[i];
7     }
8     arr[index] = val;
9     n++;
10    return n;
11 }
```

超范围插入

数组的长度在创建时就已经确定了，要实现数组的扩容，只能创建一个新数组，长度是旧数组的 2 倍，再把旧数组中的元素全部复制过去，这样就实现了数组的扩容。

数组插入元素最好情况是尾部插入，无需移动任何元素，时间复杂度为 $O(1)$ 。最坏情况是在第一个位置插入，这样就需要移动后面所有 $n - 1$ 个元素，时间复杂度为 $O(n)$ 。因此，总体的时间复杂度为 $O(n)$ 。

2.3.2 删除元素

数组的删除操作与插入操作过程相反，如果被删除的元素位于数组中间，其后的元素都需要向前挪动一位。

删除元素

```
1 int delete(int *arr, int n, int index) {  
2     if(index < 0 || index >= n) {  
3         return n;  
4     }  
5     for(int i = index + 1; i < n; i++) {  
6         arr[i-1] = arr[i];  
7     }  
8     n--;  
9     return n;  
10 }
```

数组的删除操作，由于只涉及元素的移动，时间复杂度为 $O(n)$ 。

对于删除操作，其实还存在一种取巧的方式，前提是数组元素没有顺序要求。如需要删除数组中某个元素，可直接把最后一个元素复制到被删除元素的位置，然后再删除最后一个元素。这样一来，无须进行大量的元素移动，时间复杂度降低为 $O(1)$ 。当然，这种方式只作参考，并不是删除元素主流的操作方式。

Chapter 3 链表

3.1 链表

3.1.1 单向链表 (Singly Linked List)

为避免元素的移动，采用线性表的另一种存储方式：链式存储结构。链表是一种在物理上非连续、非顺序的数据结构，由若干结点 (node) 所组成。

单向链表的每一个结点又包含两部分，一部分是存放数据的数据域 data，另一部分是指向下一个结点的指针域 next。结点可以在运行时动态生成。

```
1 typedef struct Node {  
2     dataType data;           // 数据域  
3     struct Node *next;       // 指针域  
4 } Node;
```

链表的第一个结点被称为头结点，最后一个节点被称为尾结点，尾结点的 next 指针指向空 NULL。

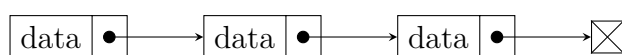


图 3.1: 单向链表

与数组按照下标来随机寻找元素不同，对于链表的其中一个结点 A，只能根据结点 A 的 next 指针来找到该结点的下一个结点 B，再根据结点 B 的 next 指针找到下一个结点 C……

数组在内存中的存储方式是顺序存储，链表在内存中的存储方式则是随机存储。链表采用了见缝插针的方式，每一个结点分布在内存的不同位置，依靠 next 指针关联起来。这样可以灵活有效地利用零散的碎片空间。

3.1.2 双向链表 (Doubly Linked List)

那么，通过链表的一个结点，如何能快速找到它的前一个结点呢？要想让每个结点都能回溯到它的前置结点，可以使用双向链表。

双向链表比单向链表稍微复杂一点，它的每一个结点除了拥有 data 和 next 指针，还拥有指向前置结点的 prev 指针。

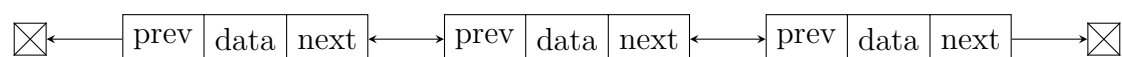


图 3.2: 双向链表

单向链表只能从头到尾遍历，只能找到后继，无法找到前驱，因此遍历的时候不会死循环。而双向链表需要多分配一个指针的存储空间，每个结点有两个指针，分别指向直接前驱和直接后继。

3.1.3 循环链表 (Circular Linked List)

除了单向链表和双向链表以外，还有循环链表。对于单向循环链表，尾结点的 next 指针指向头结点。对于双向循环链表，尾结点的 next 指针指向头结点，并且头结点的 prev 指针指向尾结点。

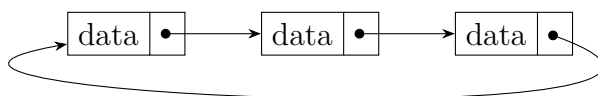


图 3.3: 单向循环链表

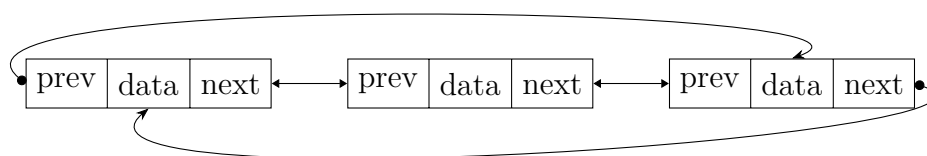


图 3.4: 双向循环链表

3.2 链表的增删改查

3.2.1 查找结点

在查找元素时，链表不像数组那样可以通过下标快速进行定位，只能从头结点开始向后一个一个结点逐一查找。

链表中的数据只能按顺序进行访问，最坏的时间复杂度是 $O(n)$ 。

查找结点

```
1 Node* search(List *head, dataType val) {
2     // 查找元素位置
3     Node *temp = head;
4     while(temp) {
5         if(temp->data == val) {
6             return temp;
7         }
8         temp = temp->next;
9     }
10    return NULL;        // 未找到
11 }
```

3.2.2 更新结点

如果不考虑查找结点的过程，链表的更新过程会像数组那样简单，直接把旧数据替换成新数据即可。

更新结点

```
1 void replace(List *head, int pos, dataType val) {
2     // 找到元素位置
3     Node *temp = head;
```

```

4   for(int i = 0; i < pos; i++) {
5       temp = temp->next;
6   }
7   temp->data = val;
8 }

```

3.2.3 插入结点

链表插入结点，分为 3 种情况：

尾部插入

把最后一个结点的 next 指针指向新插入的结点。

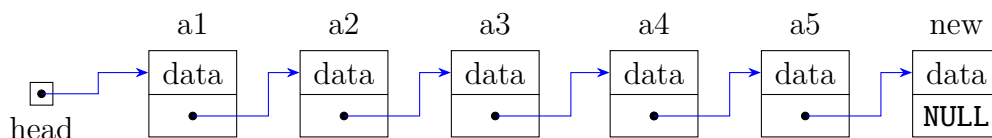


图 3.5: 尾部插入

头部插入

先把新结点的 next 指针指向原先的头结点，再把新结点设置为链表的头结点。

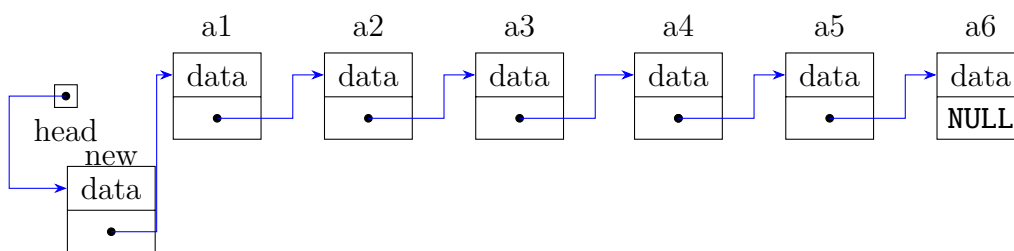


图 3.6: 头部插入

中间插入

先把新结点的 next 指针指向插入位置的结点，再将插入位置的前置结点的 next 指针指向新结点。

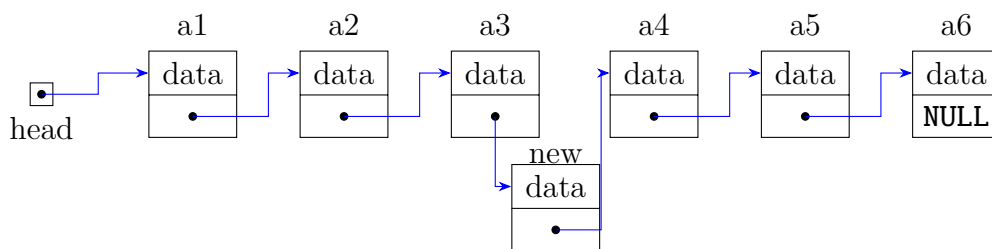


图 3.7: 中间插入

只要内存空间允许，能够插入链表的元素是无穷无尽的，不需要像数组考虑扩容的问题。如果不考虑插入之前的查找元素的过程，只考虑纯粹的插入操作，时间复杂度是 $O(1)$ 。

3.2.4 删除结点

链表的删除操作也分 3 种情况：

尾部删除

把倒数第二个结点的 next 指针指向空。

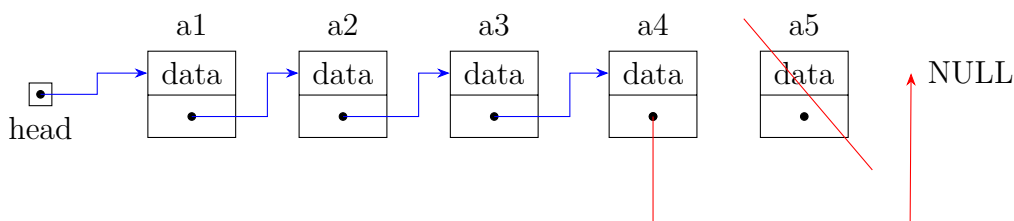


图 3.8: 尾部删除

头部删除

把链表的头结点设置为原先头结点的 next 指针。

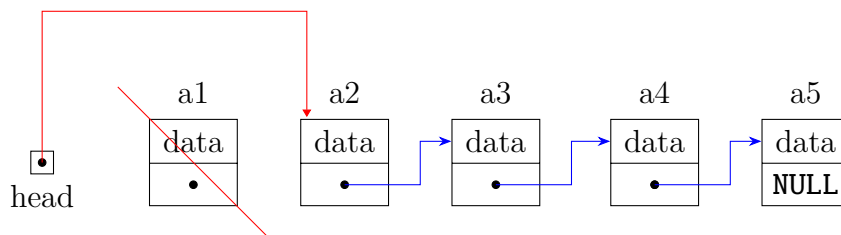


图 3.9: 头部删除

中间删除

把要删除的结点的前置结点的 next 指针，指向要删除结点的下一个结点。

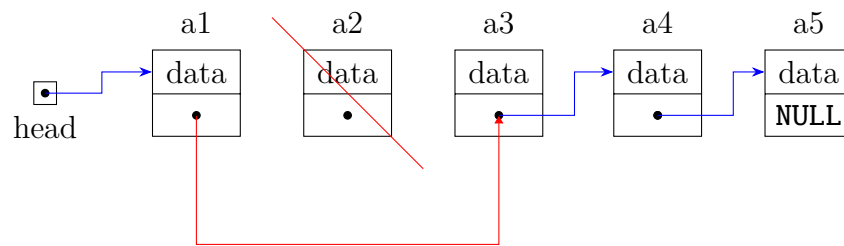


图 3.10: 中间删除

许多高级语言，如 Java，拥有自动化的垃圾回收机制，所以不用刻意去释放被删除的结点，只要没有外部引用指向它们，被删除的结点会被自动回收。

如果不考虑删除操作之前的查找的过程，只考虑纯粹的删除操作，时间复杂度是 $O(1)$ 。

3.3 带头结点的链表

3.3.1 带头结点的链表

为了方便链表的插入、删除操作，在链表加上头结点之后，无论链表是否为空，头指针始终指向头结点。因此对于空表和非空表的处理也统一了，方便了链表的操作，也减少了程序的复杂性和出现 bug 的机会。

插入结点

```
1 void insert(List *head, int pos, dataType val) {
2     Node *newNode = (Node *)malloc(sizeof(Node));
3     newNode->data = val;
4     newNode->next = NULL;
5
6     // 找到插入位置
7     Node *temp = head;
8     for(int i = 0; i < pos; i++) {
9         temp = temp->next;
10    }
11    newNode->next = temp->next;
12    temp->next = newNode;
13 }
```

删除结点

```
1 void delete(List *head, int pos) {
2     Node *temp = head;
3     for(int i = 0; i < pos; i++) {
4         temp = temp->next;
5     }
6     Node *del = temp->next;
7     temp->next = del->next;
8     free(del);
9     del = NULL;
```

3.3.2 数组 VS 链表

数据结构没有绝对的好与坏，数组和链表各有千秋。

比较内容	数组	链表
基本	一组固定数量的数据项	可变数量的数据项
大小	声明期间指定	无需指定，执行期间增长或收缩
存储分配	元素位置在编译期间分配	元素位置在运行时分配
元素顺序	连续存储	随机存储
访问元素	直接访问：索引、下标	顺序访问：指针遍历
插入/删除	速度慢	快速、高效
查找	线性查找、二分查找	线性查找
内存利用率	低效	高效

表 3.1: 数组 VS 链表

数组的优势在于能够快速定位元素，对于读操作多、写操作少的场景来说，用数组更合适一些。

相反，链表的优势在于能够灵活地进行插入和删除操作，如果需要频繁地插入、删除元素，用链表更合适一些。

3.4 倒数第 k 个结点

3.4.1 倒数第 k 个结点

输入一个链表，输出该链表中倒数第 k 个结点。例如一个链表有 6 个结点 [0, 1, 2, 3, 4, 5]，这个链表的倒数第 3 个结点是值为 3 的结点。

算法的思路是设置两个指针 p1 和 p2，它们都从头开始出发，p2 指针先出发 k 个结点，然后 p1 指针再进行出发，当 p2 指针到达链表的尾结点时，则 p1 指针的位置就是链表的倒数第 k 个结点。

倒数第 k 个结点

```
1 public static Node findLastKth(LinkedList list, int k) {
2     Node p1 = list.getHead();
3     Node p2 = list.getHead();
4
5     int i = 0;
6     while(p2 != null && i < k) {
7         p2 = p2.next;
8         i++;
9     }
10    while(p2 != null) {
11        p1 = p1.next;
12        p2 = p2.next;
13    }
14    return p1;
15 }
```

3.5 环形链表

3.5.1 环形链表

一个单向链表中有可能出现环，不允许修改链表结构，如何在时间复杂度 $O(n)$ 、空间复杂度 $O(1)$ 内判断出这个链表是有环链表？如果带环，环的长度是多少？环的入口结点是哪个？

暴力算法首先从头结点开始，依次遍历单链表的每一个结点。对于每个结点都从头重新遍历之前的所有结点。如果发现当前结点与之前结点存在相同 ID，则说明该结点被遍历过两次，链表有环。但是这种方法的时间复杂度太高。

另一种方法就是利用快慢指针，首先创建两个指针 p1 和 p2，同时指向头结点，然后让 p1 每次向后移动一个位置，让 p2 每次向后移动两个位置。在环中，快指针一定会反复遇到慢指针。比如在一个环形跑道上，两个运动员在同一地点起跑，一个运动员速度快，一个运动员速度慢。当两人跑了一段时间，速度快的运动员必然会从速度慢的运动员身后再次追上并超过。

环的长度可以通过从快慢指针相遇的结点开始再走一圈，下一次回到该点的时的移动步数，即环的长度 n。

环的入口可以利用类似获取链表倒数第 k 个结点的方法，准备两个指针 p1 和 p2，让 p2 先走 n 步，然后 p1 和 p2 一块走。当两者相遇时，即为环的入口处。

环形链表

```
1 public static Node cycleNode(LinkedList list) {
2     Node p1 = list.getHead();
3     Node p2 = list.getHead();
4
5     while(p1 != null && p2 != null) {
6         if(p2.next == null || p2.next.next == null) {
7             return null;
```



```

8         }
9         p1 = p1.next;
10        p2 = p2.next.next;
11        if(p1 == p2) {
12            return p1;
13        }
14    }
15    return null;
16 }
17
18 public static int cycleLength(LinkedList list) {
19     Node node = cycleNode(list);
20     if(node == null) {
21         return 0;
22     }
23     int len = 1;
24     Node cur = node.next;
25     while(cur != node) {
26         cur = cur.next;
27         len++;
28     }
29     return len;
30 }
31
32 public static Node cycleEntrance(LinkedList list) {
33     int n = cycleLength(list);
34     if(n == 0) {
35         return null;
36     }
37
38     Node p1 = list.getHead();
39     Node p2 = list.getHead();
40     for(int i = 0; i < n; i++) {
41         p2 = p2.next;
42     }
43
44     while(p1 != p2) {

```

```
45     p1 = p1.next;
46     p2 = p2.next;
47 }
48 return p1;
49 }
```

3.6 反转链表

3.6.1 逆序输出链表

输入一个单链表，从尾到头打印链表每个结点的值。由于单链表的遍历只能从头到尾，所以可以通过递归达到链表尾部，然后在回溯时输出每个结点的值。

逆序输出链表

```
1 public static void inverse(Node head) {  
2     if(head != null) {  
3         inverse(head.next);  
4         System.out.print(head.data + " ");  
5     }  
6 }
```

3.6.2 反转链表

反转一个链表需要调整链表中指针的方向。

递归反转法的实现思想是从链表的尾结点开始，依次向前遍历，遍历过程依次改变各结点的指向，即另其指向前一个结点。

而迭代反转法的实现思想非常直接，就是从当前链表的首结点开始，一直遍历至链表尾部，期间会逐个改变所遍历到的结点的指针域，另其指向前一个结点。

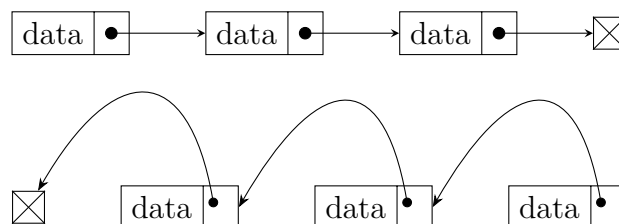


图 3.11: 反转链表

反转链表（递归）

```
1 public static Node reverseList(Node head) {  
2     if(head == null || head.next == null) {  
3         return head;  
4     }  
5     Node next = head.next;  
6     head.next = null;  
7     Node newHead = reverseList(next);  
8     next.next = head;  
9     return newHead;  
10 }
```

反转链表（迭代）

```
1 public static Node reverseListIterative(LinkedList list) {  
2     Node newHead = new Node(-1);  
3     Node head = list.getHead();  
4     while(head != null) {  
5         Node next = head.next;  
6         head.next = newHead.next;  
7         newHead.next = head;  
8         head = next;  
9     }  
10    return newHead.next;  
11 }
```

Chapter 4 栈

4.1 栈

4.1.1 栈 (Stack)

栈，又名堆栈，是一种运算受限的线性数据结构，栈只能在表尾进行插入和删除操作。

栈中的元素只能先进后出 (FILO, First In Last Out)。最早进入栈的元素所存放的位置叫作栈底 (bottom)，最后进入栈的元素存放的位置叫作栈顶 (top)。

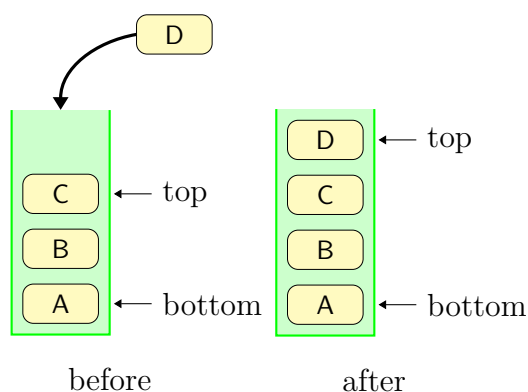


图 4.1: 栈

栈这种数据结构既可以用数组来实现，也可以用链表来实现。

4.1.2 顺序栈

使用数组方式实现的栈称为静态栈。可以根据下标来表示栈顶在数组中的位置，对于空栈，栈顶为-1。

进行入栈操作时，栈顶指针 +1；出栈时，栈顶指针-1。

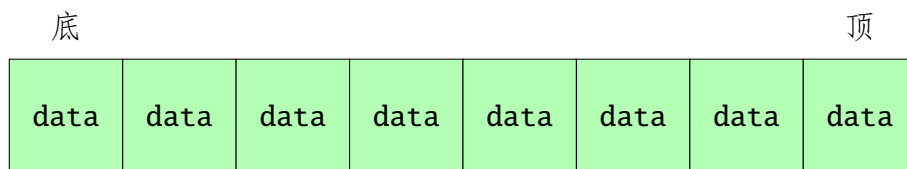


图 4.2: 顺序栈

对满栈进行入栈和对空栈进行出栈操作都会产生数组的越界并引起程序崩溃，称为上溢和下溢。因此使用顺序栈需要提前声明一个数组的大小，如果数组大小不够则可能发生数组越界，如果数组太大则会浪费一定的空间。

使用数组实现的栈的执行效率会比用链表来实现的高，入栈和出栈不需要移动大量元素，只需要移动栈顶指针即可。

4.1.3 链式栈

使用链表方式实现的栈称为动态栈。通过在表头插入一个元素来实现入栈，通过删除表尾元素来实现出栈。

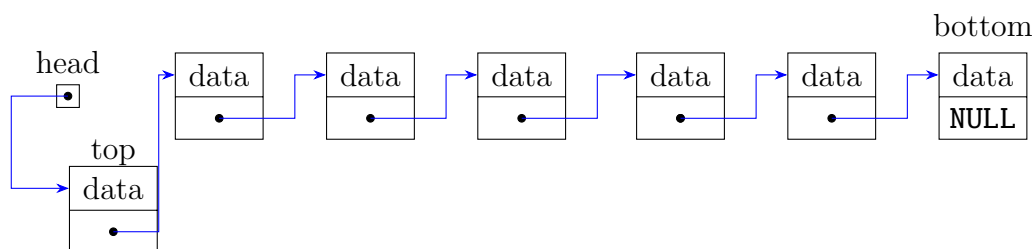


图 4.3: 链式栈

动态栈有链表的部分特性，元素与元素之间在物理存储上可以不连续，但是功能有些受限制，动态栈只能在栈顶处进行插入和删除操作，不能在栈尾或栈中间进行插入和删除操作。

动态栈的元素内存是动态分配的，避免了静态栈可能会浪费空间的问题，但是对申请和释放空间的调用开销会比较大。

4.1.4 栈的应用

栈的输出顺序和输入顺序相反，所以栈常用于对历史的回溯。例如实现递归的逻辑，就可以用栈回溯调用链。

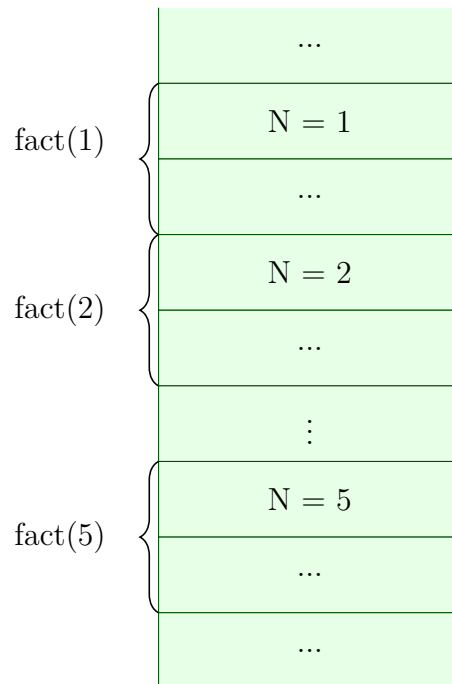


图 4.4: 函数调用栈

栈还有一个著名的应用场景就是面包屑导航，使用户在流浪页面时可以轻松地回溯到上一级更更上一级页面。



图 4.5: 面包屑导航

4.2 入栈与出栈

4.2.1 入栈 (Push)

入栈操作就是把新元素放入栈中，只允许从栈顶一侧放入元素，新元素的位置将会成为新的栈顶。最初，栈为空，栈顶的初始值为-1。每当向栈中添加元素时，栈顶指针 +1。

入栈只影响最后一个元素，不涉及元素的整体移动，所以无论是以数组还是链表实现，时间复杂度都是 $O(1)$ 。

入栈

```
1 void push(Stack *stack, dataType val) {  
2     stack->data[++stack->top] = val;  
3 }
```

4.2.2 出栈 (Pop)

出栈操作就是把新元素从栈中弹出，只有栈顶元素才允许出栈，出栈元素的前一个元素将会成为新的栈顶。从栈中移出元素，栈顶指针-1。数组中元素的删除并非真正意义上把元素从内存中清除，出栈只需对栈顶-1 即可，后期向栈中添加元素时，新元素会将旧元素覆盖。

出栈只影响最后一个元素，不涉及元素的整体移动，所以无论是以数组还是链表实现，时间复杂度都是 $O(1)$ 。

出栈

```
1 dataType pop(Stack *stack) {  
2     return stack->data[stack->top--];  
3 }
```


4.3 最小栈

4.3.1 最小栈

设计一个支持 `push()`、`pop()`、`peek()` 和 `getMin()` 操作的栈，并能在常数时间内检索到最小元素。

对于栈来说，如果一个元素 `a` 在入栈时，栈里有其它的元素 `b`、`c`、`d`，那么无论这个栈在之后经历了什么操作，只要 `a` 在栈中，`b`、`c`、`d` 就一定在栈中。因此，在操作过程中的任意一个时刻，只要栈顶的元素是 `a`，那么就可以确定栈里面现在的元素一定是 `a`、`b`、`c`、`d`。

那么可以在每个元素 `a` 入栈时把当前栈的最小值 `m` 存储起来。在这之后无论何时，如果栈顶元素是 `a`，就可以直接返回存储的最小值 `m`。

当一个元素要入栈时，取辅助栈的栈顶存储的最小值，与当前元素比较得出最小值，将这个最小值插入辅助栈中。当一个元素要出栈时，把辅助栈的栈顶元素也一并弹出。这样在任意一个时刻，栈内元素的最小值就存储在辅助栈的栈顶元素中。

最小栈

```
1 class MinStack:
2     def __init__(self):
3         self.stack = []
4         self.min_stack = [math.inf]
5
6     def push(self, data):
7         self.stack.append(data)
8         self.min_stack.append(min(data, self.min_stack[-1]))
9
10    def pop(self):
11        self.stack.pop()
12        self.min_stack.pop()
```

```
13
14     def peek(self):
15         return self.stack[-1]
16
17     def get_min(self):
18         return self.min_stack[-1]
```

4.4 括号匹配

4.4.1 括号匹配

给定一个只包括"("、")"、 "["、"]"、 "{" 和"}" 的字符串，判断字符串是否有效。有效字符串需满足左括号必须以正确的顺序用相同类型的右括号闭合。

判断括号的有效性可以使用栈来解决。通过遍历字符串，当遇到左括号时，会期望在后续的遍历中，有一个相同类型的右括号将其闭合。由于后遇到的左括号要先闭合，因此将这个左括号放入栈顶。当遇到右括号时，需要将一个相同类型的左括号闭合。此时可以取出栈顶的左括号并判断它们是否是相同类型的括号。如果不是，或者栈中没有左括号，那么字符串无效。在遍历结束后，如果为空栈，说明字符串中的所有左括号闭合。

注意有效字符串的长度一定为偶数，因此如果字符串的长度为奇数，可以直接返回判断出字符串无效，省去后续的遍历判断过程。

括号匹配

```
1 def valid_parentheses(s):
2     if len(s) % 2 == 1:
3         return False
4
5     pairs = {")": "(", "]": "[", "}": "{"}
6     stack = list()
7     for paran in s:
8         if paran in pairs:
9             if not stack or stack[-1] != pairs[paran]:
10                 return False
11             stack.pop()
12         else:
13             stack.append(paran)
14
15     return not stack
```

4.5 表达式求值

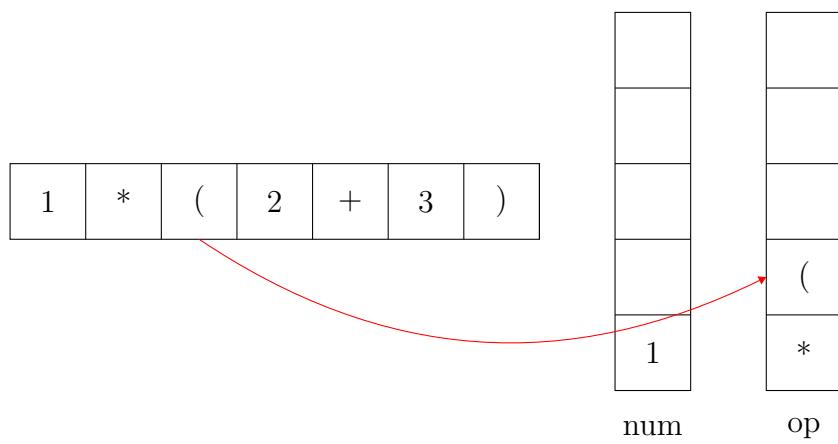
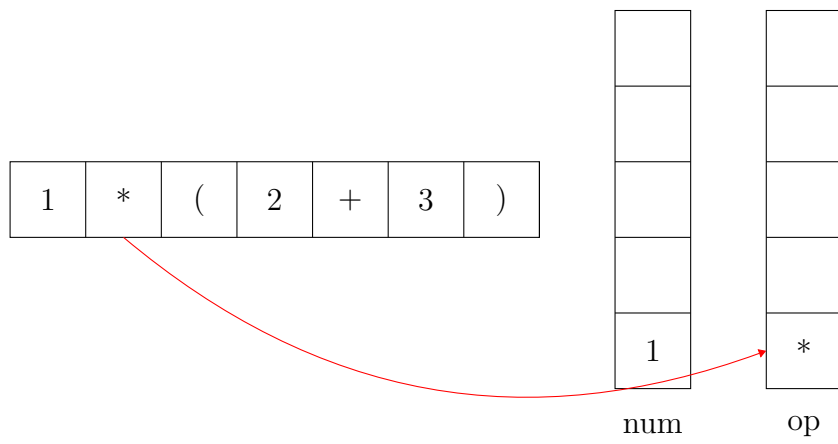
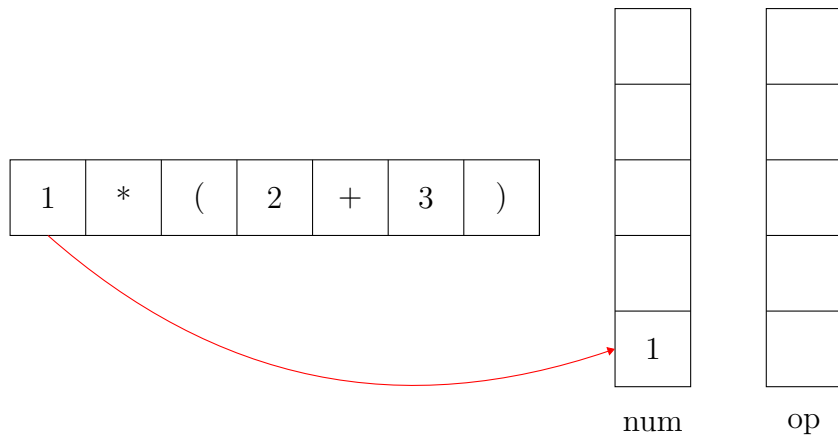
4.5.1 表达式求值

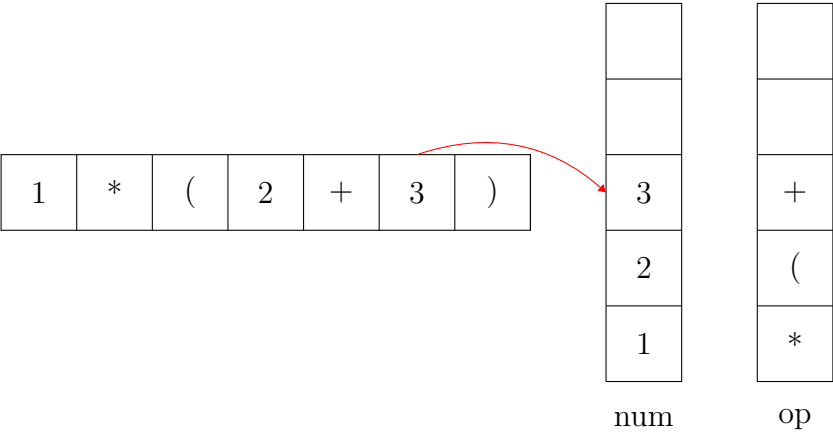
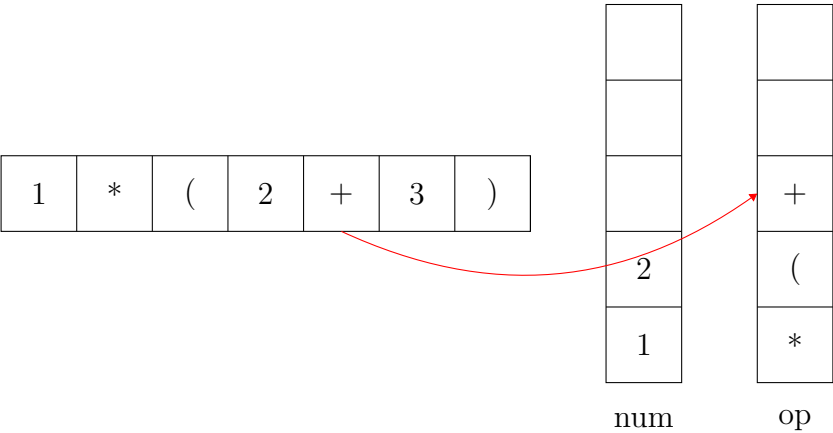
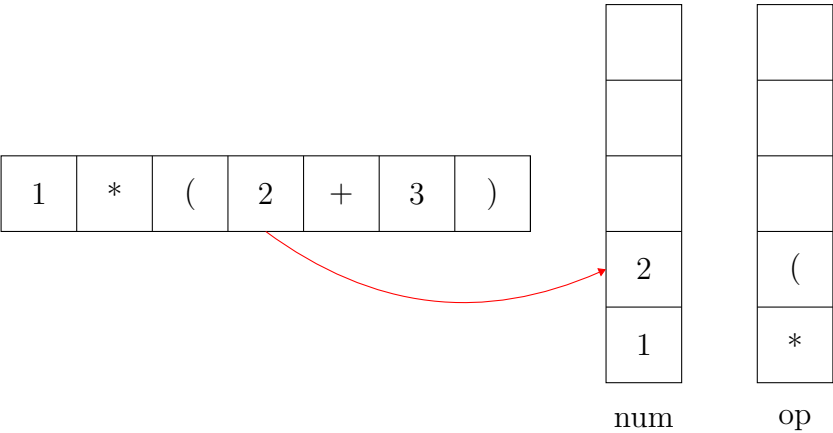
逆波兰表达式是一种后缀表达式，所谓后缀就是指运算符写在运算数的后面。平常使用的算式则是一种中缀表达式，如 $(1 + 2) * (3 + 4)$ ，该算式的逆波兰表达式写法为 $1\ 2\ +\ 3\ 4\ +\ *$ 。

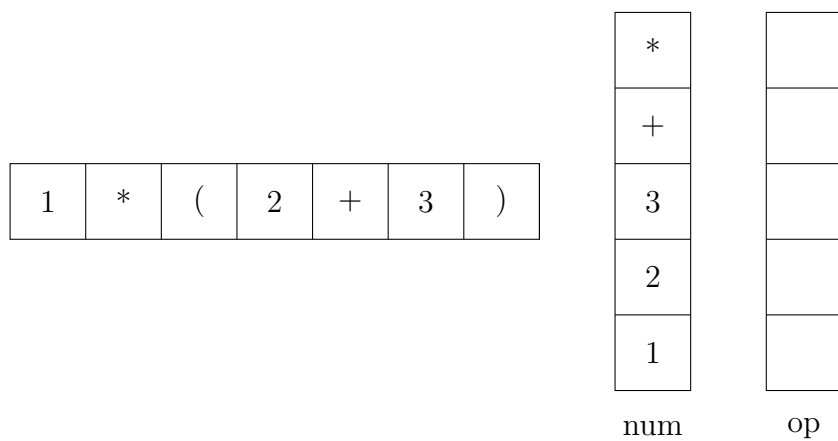
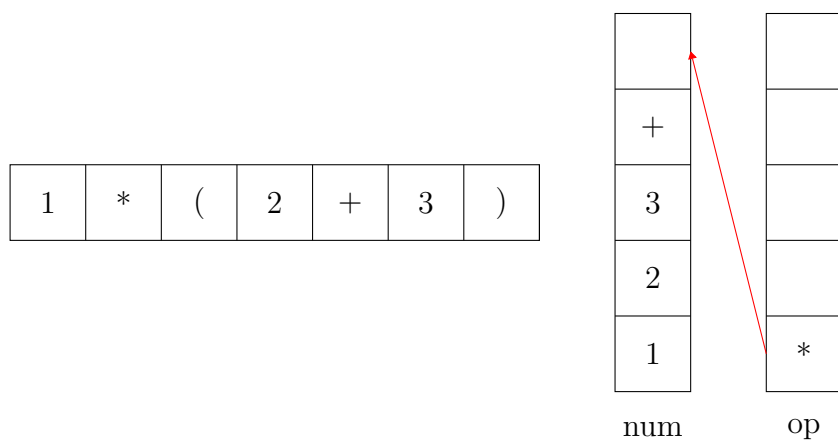
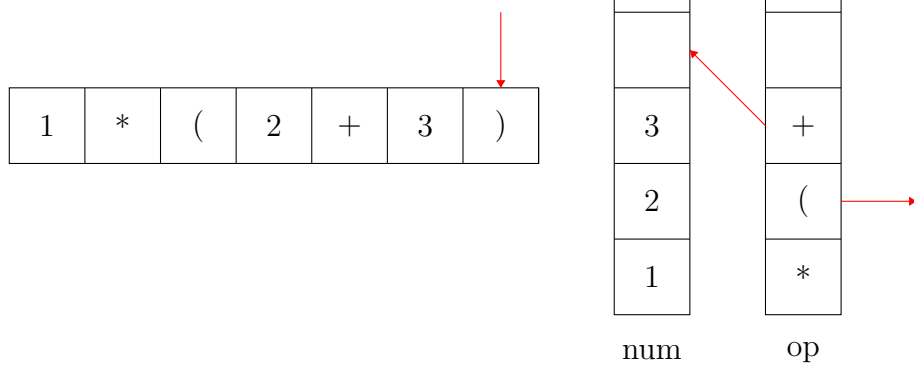
逆波兰表达式的优点在于去掉了中缀表达式中的括号后表达式无歧义，因此适合用栈操作运算。遇到数字则入栈，遇到算符则取出栈顶两个数字进行计算，并将结果压入栈中。

在对中缀表达式求值时，一般都会将其转换为后缀表达式的形式，转换过程同样需要用到栈，规则如下：

1. 如果遇到操作数，就直接将其输出。
2. 如果遇到左括号，将其放入栈中。
3. 如果遇到右括号，则一直出栈并输出，直到遇到左括号为止。注意，左括号只出栈并不输出。
4. 如果遇到任何其它的运算符，如果为栈为空，则直接入栈。否则从栈中出栈元素并输出，直到遇到优先级更低的元素（或者栈为空）位置。在出栈完这些元素后，再将当前遇到的运算符入栈。有一点需要注意，只有在遇到右括号的情况下才将左括号出栈，其它情况都不会出栈左括号。
5. 如果读取到了表达式的末尾，则将栈中所有元素依次出栈输出。







表达式求值

```

1 def priority(op):
2     """
3     运算符的优先级
4     乘除法优先级高于加减法

```

```

5         Args:
6             op (str): 运算符
7         Returns:
8             (int): 优先级
9         """
10        if op == "*" or op == "/":
11            return 2
12        elif op == "+" or op == "-":
13            return 1
14        else:
15            return 0
16
17    def infix_to_postfix(exp):
18        """
19            中缀表达式转换后缀表达式
20            转换后的后缀表达式操作数之前带空格
21        Args:
22            exp (str): 中缀表达式
23        Returns:
24            (str): 后缀表达式
25        """
26        postfix = ""    # 保存生成的后缀表达式
27        s = stack.Stack()
28
29        number = ""
30        for ch in exp:
31            # 如果是数字，保存每一位数字
32            if ch.isdigit():
33                number += ch
34                continue
35
36            # 如果读取一个完整数字，直接输出
37            if len(number) > 0:
38                postfix += number + " "
39                number = ""
40
41            # 空格忽略

```



```

42     if ch == " ":
43         continue
44
45     # 如果是运算符，并且空栈，则直接入栈
46     if s.is_empty():
47         s.push(ch)
48     # 如果遇到左括号，将其放入栈中
49     elif ch == "(":
50         s.push(ch)
51     # 如果遇到右括号，则一直出栈并输出，直到遇到左括号为止
52     # 注意，左括号只出栈并不输出
53     elif ch == ")":
54         while s.peek() != "(":
55             postfix += s.pop() + " "
56             s.pop()
57     # 如果遇到任何其它的运算符，如果为栈为空，则直接入栈
58     # 否则从栈中出栈元素并输出，直到遇到优先级更低的元素（或为空）
59     # 在出栈完这些元素后，再将当前遇到的运算符入栈
60     # 只有遇到右括号的情况下才将左括号出栈
61     else:
62         while not s.is_empty()
63             and priority(ch) <= priority(s.peek()):
64             postfix += s.pop() + " "
65             s.push(ch)
66
67     # 如果读取一个完整数字，直接输出
68     if len(number) > 0:
69         postfix += number + " "
70         number = ""
71
72     while not s.is_empty():
73         postfix += s.pop() + " "
74
75     return postfix.rstrip()
76
77 def calculate(postfix):
78     """

```

```

79     表达式求值
80     Args:
81         postfix (str): 后缀表达式
82     Returns:
83         (int): 表达式结果
84     """
85     s = stack.Stack()
86
87     tokens = postfix.split()
88     for token in tokens:
89         # 数字则入栈
90         try:
91             s.push(int(token))
92         # 运算符则出栈2次，将计算结果入栈
93         except ValueError:
94             num2 = s.pop()
95             num1 = s.pop()
96             if token == '+':
97                 s.push(num1 + num2)
98             elif token == '-':
99                 s.push(num1 - num2)
100             elif token == '*':
101                 s.push(num1 * num2)
102             elif token == '/':
103                 s.push(int(num1 / num2))
104     return s.pop()

```

Chapter 5 队列

5.1 队列

5.1.1 队列 (Queue)

队列是一种运算受限的线性数据结构，不同于栈的先进后出 (FILO)，队列中的元素只能先进先出 (FIFO, First In First Out)。

队列的出口端叫作队头 (front)，队列的入口端叫作队尾 (rear)。队列只允许在队尾进行入队 (enqueue)，在队头进行出队 (dequeue)。

与栈类似，队列既可以用数组来实现，也可以用链表来实现。其中用数组实现时，为了入队操作的方便，把队尾位置规定为最后入队元素的下一个位置。

5.1.2 入队 (enqueue)

入队就是把新元素放入队列中，只允许在队尾的位置放入元素，新元素的下一个位置将会成为新的队尾。入队操作的时间复杂度是 $O(1)$ 。

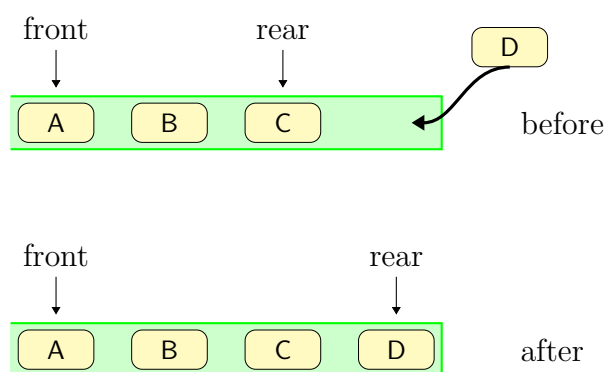


图 5.1: 入队

5.1.3 出队 (dequeue)

出队就是把元素移出队列，只允许在队头一侧移出元素，出队元素的后一个元素将成为新的队头。出队操作的时间复杂度是 $O(1)$ 。

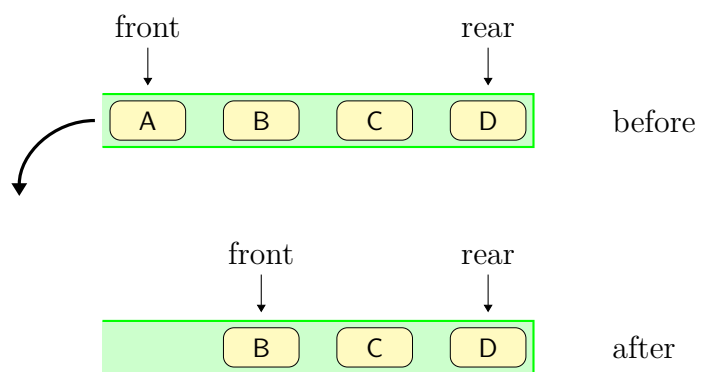


图 5.2: 出队

5.2 循环队列

5.2.1 循环队列 (Circular Queue)

如果不断出队，队头左边的空间就失去了作用，那队列的容量就会变得越来越小。

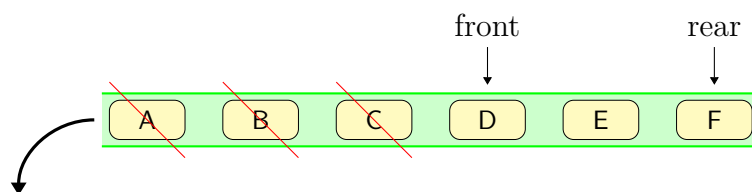


图 5.3: 队列存在的问题

用数组实现的队列可以采用循环队列的方式来维持队列容量的恒定。为充分利用空间，克服假溢出现象，在数组不做扩容的情况下，将队列想象为一个首尾相接的圆环，可以利用已出队元素留下的空间，让队尾指针重新指回数组的首位。这样一来整个队列的元素就循环起来了。

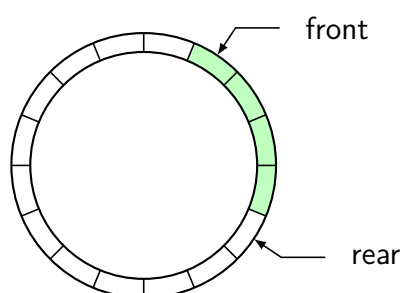


图 5.4: 循环队列

在物理存储上，队尾的位置也可以在队头之前。当再有元素入队时，将其放入数组的首位，队尾指针继续后移即可。队头和队尾互相追赶，这个追赶的过程就是入队的出队的过程。

如果队尾追上队头说明队列满了，如果队头追上队尾说明队列为空。循环队列并非真正地把数组弯曲，利用求余操作就能使队头和队尾指针不会跑出数组的范围，逻辑上实现了弯曲的效果。

假设数组的最大容量为 MAX:

- 入队时队尾指针后移: $(rear + 1) \% MAX$
- 出队时队头指针后移: $(front + 1) \% MAX$
- 判断队满: $(rear + 1) \% MAX == front$
- 判断队空: $front == rear$

需要注意的是，队尾指针指向的位置永远空出一位，所以队列最大容量比数组长度小 1。

入队

```
1 void enqueue(Queue *queue, dataType val) {  
2     queue->data[queue->rear] = val;  
3     queue->rear = (queue->rear + 1) % queue->max;  
4 }
```

出队

```
1 dataType dequeue(Queue *queue) {  
2     dataType ret = queue->data[queue->front];  
3     queue->front = (queue->front + 1) % queue->max;  
4     return ret;  
5 }
```

5.3 栈实现队列

5.3.1 栈实现队列

栈的特性是 FILO，而队列是 FIFO，因此可以使用两个栈来实现队列的效果。

可以将一个栈当作输入栈，用于 `push()` 数据，另一个栈当作输出栈，用于 `pop()` 和 `peek()` 数据。每次 `pop()` 或 `peek()` 时，若输出栈为空则将输入栈的全部数据依次弹出并压入输出栈，这样输出栈从栈顶往栈底的顺序就是队列从队首往队尾的顺序。

用两个栈来实现队列的情况在生活中也经常出现。去医院挂号等待，等待的时候把病历给护士，护士面前放了两堆病历。等着无聊就看护士是怎么管理病历的。发现一个堆是倒着放的，一个堆是正着放的。新过来的病人把病历给她时，她就病历倒着放到第一堆，有病人看病结束后，从第二堆的开头翻一个新的病历，然后叫病人。如果第二堆没了话，就直接把第一堆翻过来放到第二堆上面。

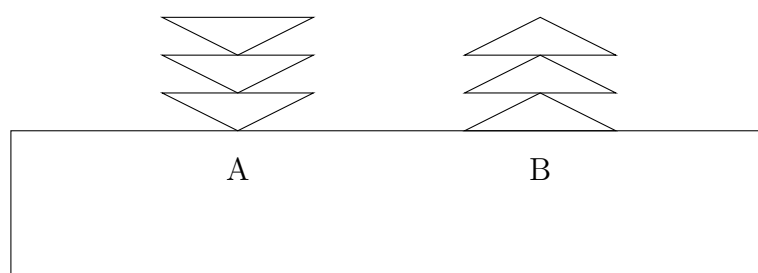


图 5.5: 双栈实现队列

双栈实现队列

```
1 class Queue:
2     def __init__(self):
3         self.in_stack = list()
4         self.out_stack = list()
5
6     def is_empty(self):
7         return not self.in_stack and not self.out_stack
8
```

```
9     def enqueue(self, data):
10         self.in_stack.append(data)
11
12     def dequeue(self):
13         if not self.out_stack:
14             while self.in_stack:
15                 self.out_stack.append(self.in_stack.pop())
16         return self.out_stack.pop()
```

用双栈实现的队列 `push()` 的时间复杂度为 $O(1)$, `pop()` 和 `peek()` 为均摊 $O(1)$, 因为对于每个元素, 至多入栈和出栈各两次。

5.4 队列实现栈

5.4.1 双队列实现栈

为了满足栈的特性，在使用队列实现栈时，应满足队列前端的元素是最后入栈的元素。可以使用两个队列实现栈的操作，其中 queue1 用于存储栈内的元素，queue2 作为入栈操作的辅助队列。

入栈操作时，首先将元素入队到 queue2，然后将 queue1 的全部元素依次出队并入队到 queue2，此时 queue2 的前端的元素即为新入栈的元素，再将 queue1 和 queue2 互换，则 queue1 的元素即为栈内的元素，queue1 的前端和后端分别对应栈顶和栈底。

由于 queue1 用于存储栈内的元素，判断栈是否为空时，只需要判断 queue1 是否为空即可。

5.4.2 单队列实现栈

在两个队列实现栈的方法中，其中一个队列的作用相当于临时变量。因此只使用一个队列就能实现栈了。

入栈操作时，首先获得入栈前的元素个数 n，然后将元素入队到队列，再将队列中的前 n 个元素（即除了新入栈的元素之外的全部元素）依次出队并入队到队列，此时队列的前端的元素即为新入栈的元素，且队列的前端和后端分别对应栈顶和栈底。

单队列实现栈

```
1 import collections
2
3 class Stack:
4     def __init__(self):
```

```
5         self.queue = collections.deque()
6
7     def is_empty(self):
8         return not self.queue
9
10    def push(self, data):
11        n = len(self.queue)
12        self.queue.append(data)
13        for _ in range(n):
14            self.queue.append(self.queue.popleft())
15
16    def pop(self):
17        return self.queue.popleft()
18
19    def peek(self):
20        return self.queue[0]
```

5.5 双端队列

5.5.1 双端队列 (Deque, Double Ended Queue)

双端队列是一种同时具有队列和栈的性质的数据结构，双端队列可以从其两端插入和删除元素。

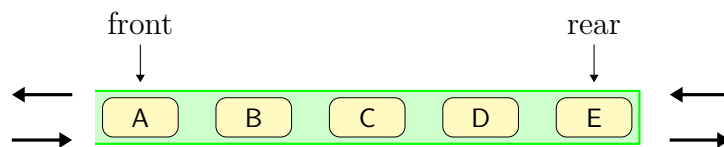


图 5.6: 双端队列

双端队列

```
1 class Deque:
2     def __init__(self):
3         self.data = []
4
5     def is_empty(self):
6         return not self.data
7
8     def add_front(self, val):
9         self.data.insert(0, val)
10
11     def add_rear(self, val):
12         self.data.append(val)
13
14     def remove_front(self):
15         if not self.is_empty():
16             return self.data.pop(0)
17
18     def remove_rear(self):
19         if not self.is_empty():
20             return self.data.pop()
21
22     def get_front(self):
```

```
23         if not self.is_empty():
24             return self.data[0]
25
26     def get_rear(self):
27         if not self.is_empty():
28             return self.data[len(self.data)-1]
```

Chapter 6 哈希表

6.1 哈希表

6.1.1 哈希表 (Hash Table)

例如开发一个学生管理系统，需要有通过输入学号快速查出对应学生的姓名的功能。这里不必每次都去查询数据库，而可以在内存中建立一个缓存表，这样做可以提高查询效率。

学号	姓名
001121	张三
002123	李四
002931	王五
003278	赵六

表 6.1: 学生名单

再例如需要统计一本英文书里某些单词出现的频率，就需要遍历整本书的内容，把这些单词出现的次数记录在内存中。

单词	出现次数
this	108
and	56
are	79
by	46

表 6.2: 词频统计

因为这些需要，一个重要的数据结构诞生了，这个数据结构就是哈希表。哈希表也称散列表，哈希表提供了键（key）和值（value）的映射关系，只要给出一个 key，就可以高效地查找到它所匹配的 value。

哈希表的时间复杂度几乎是常量 $O(1)$ ，即查找时间与问题规模无关。

哈希表的两项基本工作：

1. 计算位置：构造哈希函数确定关键字的存储位置。
2. 解决冲突：应用某种策略解决多个关键字位置相同的问题。

6.2 哈希函数

6.2.1 哈希函数 (Hash Function)

哈希的基本思想是将键 key 通过一个确定的函数，计算出对应的函数值 value 作为数据对象的存储地址，这个函数就是哈希函数。

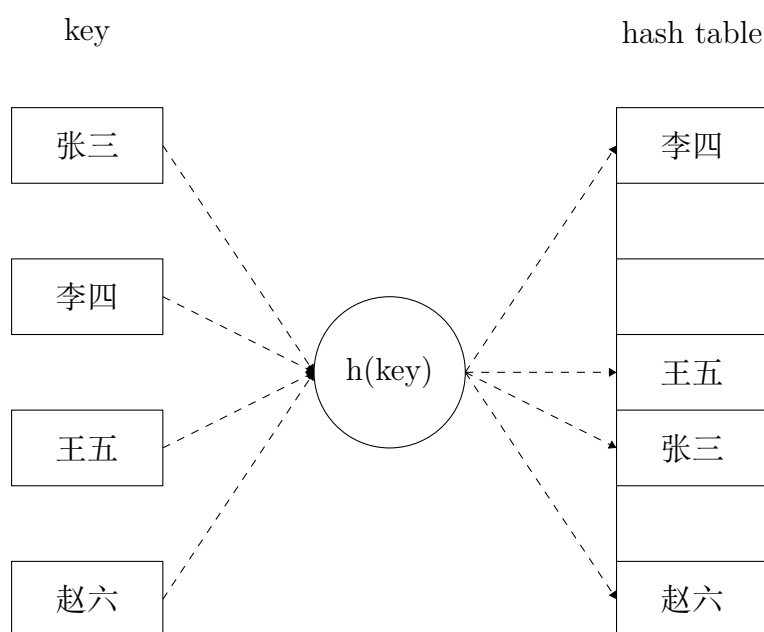


图 6.1: 哈希函数

哈希表本质上也是一个数组，可是数组只能根据下标来访问，而哈希表的 key 则是以字符串类型为主的。

在不同的语言中，哈希函数的实现方式是不一样的。假设需要存储整型变量，转化为数组的下标就不难实现了。最简单的转化方式就是按照数组长度进行取模运算。

一个好的哈希函数应该考虑两个因素：

1. 计算简单，以便提高转换速度。
2. 关键字对应的地址空间分布均匀，以尽量减少冲突。

6.2.2 数字关键字的哈希函数构造方法

对于数字类型的关键字，哈希函数有以下几种常用的构造方法：

直接定址法

取关键字的某个线性函数值为散列地址。

$$h(key) = a * key + b$$

例如根据出生年份计算人口数量 $h(key) = key - 1990$ ：

地址	出生年份	人数
0	1990	1285 万
1	1991	1281 万
2	1992	1280 万
...
10	2000	1250 万
...
21	2011	1180 万

表 6.3: 直接定址法

除留余数法

哈希函数为 $h(key) = key \% p$ ， p 一般取素数。

例如 $h(key) = key \% 17$ ：

地址	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
关键字	34	18	2	20			23	7	42		27	11		30		15	

表 6.4: 除留余数法

数字分析法

分析数字关键字在各位上的变化情况，取比较随机的位作为散列地址。

例如取 11 位手机号码的后 4 位作为地址 $h(\text{key}) = \text{int}(\text{key} + 7)$ 。

再例如取 18 位身份证号码中变化较为随机的位数：

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
3	3	0	1	0	6	1	9	9	0	1	0	0	8	0	4	1	9
省		市		区		年				月		日		辖		校验	

表 6.5: 数字分析法

折叠法

把关键字分割成位数相同的几个部分，然后叠加。

例如将整数 56793542 每三位进行分割：

$$\begin{array}{r}
 542 \\
 793 \\
 + \quad 056 \\
 \hline
 = \quad 1319
 \end{array}$$

$$h(56793542) = 319$$

平方取中法

计算关键字的平方，取中间几位。

例如整数 56793542：

$$\begin{array}{r}
 56793542 \\
 * \quad 56793542 \\
 \hline
 = \quad 3225506412905764
 \end{array}$$

$$h(56793542) = 641$$

6.2.3 字符串关键字的哈希函数构造方法

对于字符串类型的关键字，哈希函数有以下几种常用的构造方法：

ASCII 码加和法

$$h(key) = \left(\sum key[i] \right) \bmod N$$

但是对于某些字符串会导致严重冲突，例如：a3、b2、c1 或 eat、tea 等。

移位法

取前 3 个字符移位。

$$h(key) = (key[0] \times 27^2 + key[1] \times 27 + key[2]) \bmod N$$

对于一些字符串仍然会冲突，例如 string、strong、street、structure 等。

一个有效的改进是涉及关键字中所有 n 个字符：

$$h(key) = \left(\sum_{i=0}^{n-1} key[n-i-1] \times 32^i \right) \bmod N$$

哈希函数

快速计算 $h(abcde) = a * 32^4 + b * 32^3 + c * 32^2 + d * 32 + e$

```
1 int hash(char *key, int tableSize) {
2     int h = 0;           // hash value
3     int i = 0;
4     while(key[i] != '\0') {
5         h = (h << 5) + key[i];
6         i++;
7     }
8     return h % tableSize;
9 }
```

凯撒加密

```
1  /**
2  * @brief 凯撒加密
3  * @note 加密算法: ciphertext[i] = (plaintext[i] + Key) % 128
4  * @param plaintext: 明文
5  * @retval 密文
6  */
7  char* encrypt(char *plaintext) {
8      int n = strlen(plaintext);
9      char *ciphertext = (char *)malloc((n + 1) * sizeof(char));
10     for(int i = 0; i < n; i++) {
11         ciphertext[i] = (plaintext[i] + KEY) % 128;
12     }
13     ciphertext[n] = '\0';
14     return ciphertext;
15 }
16
17 /**
18 * @brief 凯撒解密
19 * @note 解密算法: plaintext[i] = (ciphertext[i] - key + 128) % 128
20 * @param ciphertext: 密文
21 * @retval 明文
22 */
23 char* decrypt(char *ciphertext) {
24     int n = strlen(ciphertext);
25     char *plaintext = (char *)malloc((n + 1) * sizeof(char));
26     for(int i = 0; i < n; i++) {
27         plaintext[i] = (ciphertext[i] - KEY + 128) % 128;
28     }
29     plaintext[n] = '\0';
30     return plaintext;
31 }
```

6.3 冲突处理

6.3.1 装填因子 (Load Factor)

假设哈希表空间大小为 m ，填入表中元素个数是 n ，则称 $\alpha = n/m$ 为哈希表的装填因子。

当哈希表元素太多，即装填因子 α 太大时，查找效率会下降。实用最大装填因子一般取 $0.5 \leq \alpha \leq 0.85$ 。当装填因子过大时，解决的方法是加倍扩大哈希表，这个过程叫作再散列 (rehashing)。

再散列的过程需要遍历原哈希表，把所有的关键字重新散列到新数组中。为什么需要重新散列呢？因为长度扩大以后，散列的规则也随之改变。经过扩容，原本拥挤的哈希表重新变得稀疏，原有的关键字也重新得到了尽可能均匀的分配。

装填因子也是影响产生哈希冲突的因素之一。当不同的关键字可能会映射到同一个散列地址上，就导致了哈希冲突 (collision)，即 $h(key_i) = h(key_j)$, $key_i \neq key_j$ ，因此需要某种冲突解决策略。

例如有 11 个数据对象的集合 $\{18, 23, 11, 20, 2, 7, 27, 30, 42, 15, 34, 35\}$ ，哈希表的大小为 17，哈希函数选择 $h(key) = key \% size$ 。

地址	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
关键字	34	18	2	20			23	7	42		27	11		30		15	

在插入最后一个关键字 35 之前，都没有产生任何冲突。但是 $h(35) = 1$ ，位置已有对象，就导致了冲突。

常用的处理冲突的思路有两种：

1. 开放地址法 (open addressing)：一旦产生了冲突，就按某种规则去寻找另一空地址。开放地址法主要有线性探测法、平方探测法（二次探测法）和双散列法。
2. 分离链接法：将相应位置上有冲突的所有关键字存储在同一个单链表中。

6.3.2 线性探测法 (Linear Probing)

当产生冲突时，以增量序列 1, 2, 3, ..., n - 1 循环试探下一个存储地址。

例如序列 {47, 7, 29, 11, 9, 84, 54, 20, 30}，哈希表表长为 13，哈希函数 $h(\text{key}) = \text{key} \% 11$ ，用线性探测法处理冲突。

key	47	7	29	11	9	84	54	20	30
$h(\text{key})$	3	7	7	0	9	7	10	9	8
冲突次数	0	0	1	0	0	3	1	3	6

地址	0	1	2	3	4	5	6	7	8	9	10	11	12	Δ
插入 47				47										0
插入 7				47				7						0
插入 29				47				7	29					1
插入 11	11			47				7	29					0
插入 9	11			47				7	29	9				0
插入 84	11			47				7	29	9	84			3
插入 54	11			47				7	29	9	84	54		1
插入 20	11			47				7	29	9	84	54	20	3
插入 30	11	30		47				7	29	9	84	54	20	6

表 6.6: 线性探测法

线性探测法的缺陷在于容易出现聚集现象。

6.3.3 平方探测法 (Quadratic Probing)

平方探测法也称为二次探测法，以增量序列 $1^2, -1^2, 2^2, -2^2, \dots, q^2, -q^2$ ($q \leq \lfloor N/2 \rfloor$) 循环试探下一个存储地址。

例如序列 {47, 7, 29, 11, 9, 84, 54, 20, 30}，哈希表表长为 11，哈希函数 $h(\text{key}) = \text{key} \% 11$ ，用平方探测法处理冲突。

key	47	7	29	11	9	84	54	20	30
h(key)	3	7	7	0	9	7	10	9	8
冲突次数	0	0	1	0	0	2	0	3	3

地址	0	1	2	3	4	5	6	7	8	9	10	Δ
插入 47				47								0
插入 7				47				7				0
插入 29				47				7	29			1
插入 11	11			47				7	29			0
插入 9	11			47				7	29	9		0
插入 84	11			47			84	7	29	9		-1
插入 54	11			47			84	7	29	9	54	0
插入 20	11		20	47			84	7	29	9	54	4
插入 30	11	30	20	47			84	7	29	9	54	4

表 6.7: 平方探测法

但是只要还有空间，平方探测法就一定能找到空闲位置吗？

例如对于以下哈希表，插入关键字 11，哈希函数 $h(\text{key}) = \text{key} \% 5$ ，用平方探测法处理冲突。

下标	0	1	2	3	4
key	5	6	7		

表 6.8: 平方探测法存在的问题

对关键字 11 进行平方探测的结果一直在下标 0 和 2 之间波动，永远无法达到其它空的位置。

但是有定理证明，如果哈希表长度是满足 $4k + 3$ ($k \in \mathbb{Z}^+$) 形式的素数时，平方探测法就可以探查整个哈希表空间。

6.3.4 双散列探测法 (Double Hashing)

设定另一个哈希函数 $h_2(key)$ ，探测序列为 $h_2(key), 2h_2(key), 3h_2(key), \dots$ 。

探测序列应该保证所有的散列存储单元都应该能够被探测到，选择以下形式有良好的效果：

$$h_2(key) = p - (key \% p) \quad (p < N \wedge p, N \in \text{素数})$$

6.3.5 分离链接法

分离链接法也称拉链法、链地址法，将相应位置上有冲突的所有关键字存储在同一个单链表中。

例如关键字序列为 $\{47, 7, 29, 11, 16, 92, 22, 8, 3, 50, 37, 89, 94, 21\}$ ，哈希函数 $h(key) = key \% 11$ ，用分离链接法处理冲突。

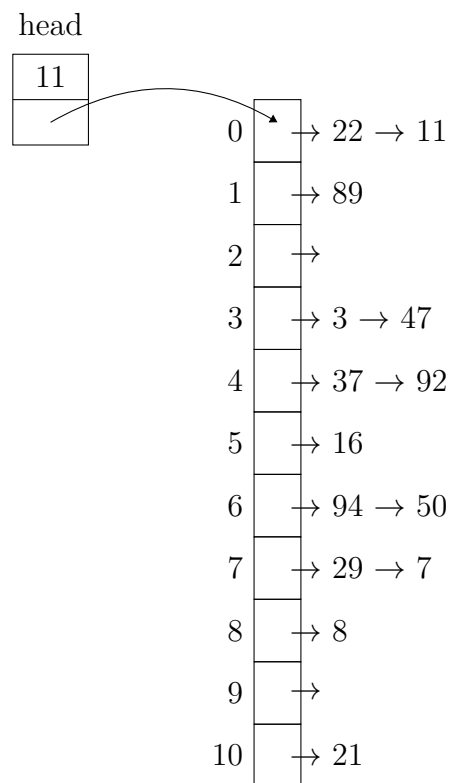


图 6.2: 分离链接法

6.4 性能分析

6.4.1 性能分析

哈希表的平均查找长度 (ASL, Average Search Length) 用来度量哈希表查找效率。关键字的比较次数, 取决于产生冲突的多少。影响产生冲突多少有三个因素:

1. 哈希函数是否均匀
2. 处理冲突的方法
3. 哈希表的装填因子 α

合理的最大装填因子 α 应该不超过 0.85, 选择合适的哈希函数可以使哈希表的查找效率期望是常数 $O(1)$, 它几乎与关键字的空间大小 n 无关。这是以较小的 α 为前提, 因此哈希表是一个以空间换时间的结构。

哈希表的存储对关键字是随机的, 因此哈希表不便于顺序查找、范围查找、最大值/最小值查找等操作。

Chapter 7 分治法

7.1 分治法

7.1.1 分治法 (Divide and Conquer)

分治策略是将原问题分解为 k 个子问题，并对 k 个子问题分别求解。如果子问题的规模仍然不够小，则再划分为 k 个子问题，如此递归的进行下去，直到问题规模足够小，很容易求出其解为止。

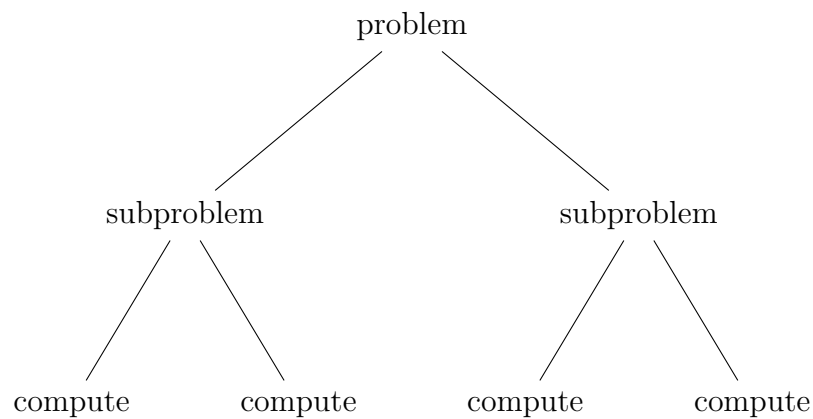


图 7.1: 分治法

将求出的小规模的问题的解合并为一个更大规模的问题的解，自底向上逐步求出原来问题的解。

分治法的适用条件有以下四点：

1. 该问题的规模缩小到一定的程度就可以容易地解决。
2. 该问题可以分解为若干个规模较小的相同问题，即该问题具有最优子结构性质。
3. 利用该问题分解出的子问题的解可以合并为该问题的解。

4. 该问题所分解出的各个子问题是相互独立的，即子问题之间不包含公共的子问题（如果各子问题是不独立的，则分治法要做许多不必要的工作，重复地解公共的子问题，此时虽然也可用分治法，但一般用动态规划较好）。

人们在大量事件中发现，在用分治法设计算法时，最好使子问题的规模大致相同。即将一个问题分成大小相等的 k 个子问题的处理方法是行之有效的。这种使子问题规模大致相等的做法是出自一种平衡子问题的思想，它几乎总是比子问题规模不等的做法要好。

7.1.2 二分搜索

一个装有 16 个硬币的袋子，16 个硬币中有一个是伪造的，并且那个伪造的硬币比真的硬币要轻一些，要求找出这个伪造的硬币。只提供一台可用来比较两组硬币重量的仪器，利用这台仪器，可以知道两组硬币的重量是否相同。

算法思想是将 16 个硬币等分成 A、B 两份，将 A 放仪器的一边，B 放另一边，如果 A 袋轻，则表明伪币在 A，解子问题 A 即可，否则解子问题 B。

二分搜索每执行一次循环，待搜索数组的大小减少一半，在最坏情况下，循环被执行了 $O(\log n)$ 次，循环体内运算需要 $O(1)$ 时间。因此，整个算法在最坏情况下的计算时间复杂性为 $O(\log n)$ 。

7.2 大整数加法

7.2.1 大整数加法

如果有两个很大的整数，如何求出它们的和？

这还不简单？直接用 long 类型存储，在程序里相加不就行了？

C/C++ 中的 int 类型能表示的范围是 $-2^{31} \sim 2^{31} - 1$ ，unsigned 类型能表示的范围是 $0 \sim 2^{32} - 1$ ，所以 int 和 unsigned 类型变量都不能保存超过 10 位的整数。

有时需要参与运算的数可能会远远不止 10 位，例如计算 100! 的精确值。即便使用能表示很大数值范围的 double 变量，但是由于 double 变量只有 64 位，精度也不足以表示一个超过 100 位的整数。我们称这种基本数据类型无法表示的整数为大整数。

在小学的时候，老师教我们用列竖式的方式计算两个整数的和。

$$\begin{array}{r} 426709752318 \\ + 95481253129 \\ \hline = 522191005447 \end{array}$$

不仅仅是人脑，对于计算机来说同样如此。对于大整数，我们无法一步到位直接算出结果，所以不得不把计算拆解成一个一个子步骤。

可是，既然大整数已经超出了 long 类型的范围，我们如何来存储这样的整数呢？

存放大整数最简单的方法就是使用数组，可以用数组的每一个元素存储整数的每一个数位。如果给定大整数的最长位数是 n ，那么按位计算的时间复杂度是 $O(n)$ 。

大整数加法

```
1 def big_int_add(num1, num2):
```

```

2     # 其中一个数为0，直接返回另一个数
3     if num1 == "0":
4         return num2
5     elif num2 == "0":
6         return num1
7
8     # 计算两个数中较长的整数位数
9     max_len = max(len(num1), len(num2))
10    # 让位数较短的整数前面补0对齐
11    num1 = '0' * (max_len - len(num1)) + num1
12    num2 = '0' * (max_len - len(num2)) + num2
13
14    result = ""          # 结果
15    carry = 0           # 保存进位
16    # 从右往左逐位相加
17    for i in range(max_len - 1, -1, -1):
18        s = int(num1[i]) + int(num2[i]) + carry
19        result = str(s % 10) + result
20        carry = s // 10
21
22    # 判断最高位是否有进位
23    if carry > 0:
24        result = str(carry) + result
25
26    # 去除结果前面多余的0
27    i = 0
28    while result[i] == '0':
29        i += 1
30    return result[i:]

```

这种思路其实还存在一个可优化的地方。我们之前是把大整数按照每一个十进制数位来拆分，比如较大整数的长度有 50 位，那么需要创建一个 51 位的数组，数组的每个元素存储其中一位。

真的有必要把原整数拆分得那么细吗？显然不需要，只需要拆分到可以被直接计算的度就够了。int 类型的取值范围是 $-2147483648 \sim 2147483647$ ，最多有 10

位整数。为了防止溢出，可以把大整数的每 9 位作为数组的一个元素，进行加法运算。如此一来，占用空间和运算次数，都被压缩了 9 倍。

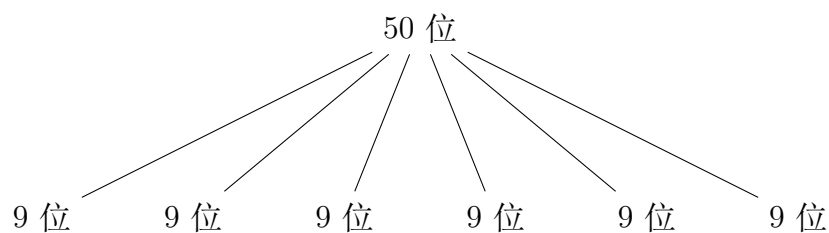


图 7.2: 大整数加法优化

在 Java 中，工具类 `BigInteger` 和 `BigDecimal` 的底层实现也是把大整数拆分成数组进行运算，与此处的思路大体类似。

Chapter 8 排序算法

8.1 排序算法

8.1.1 排序算法

应用到排序的常见比比皆是，例如当开发一个学生管理系统时需要按照学号从小到大进行排序，当开发一个电商平台时需要把同类商品按价格从低到高进行排序，当开发一款游戏时需要按照游戏得分从多到少进行排序。

根据时间复杂度的不同，主流的排序算法可以分为三类：

1. $O(n^2)$ ：冒泡排序、选择排序、插入排序
2. $O(n\log n)$ ：归并排序、快速排序、堆排序
3. $O(n)$ ：计数排序、桶排序、基数排序

在算法界还存在着更多五花八门的排序，它们有些基于传统排序变形而来，有些则是脑洞大开，如鸡尾酒排序、猴子排序、睡眠排序等。

例如睡眠排序，对于待排序数组中的每一个元素，都开启一个线程，元素值是多少，就让线程睡多少毫秒。当这些线程陆续醒来的时候，睡得少的线程线性来，睡得多的线程后醒来。睡眠排序虽然挺有意思，但是没有任何实际价值。启动大量线程的资源消耗姑且不说，数值接近的元素也未必能按顺序输出，而且一旦遇到很大的元素，线程睡眠时间可能超过一个月。

8.1.2 稳定性

排序算法还可以根据其稳定性，划分为稳定排序和不稳定排序：

- 稳定排序：值相同的元素在排序后仍然保持着排序前的顺序。
- 不稳定排序：值相同的元素在排序后打乱了排序前的顺序。

	0	1	2	3	4
原始数列	5	8	6	6	3
不稳定排序	3	5	6	6	8
稳定排序	3	5	6	6	8

图 8.1: 排序稳定性

8.2 冒泡排序

8.2.1 冒泡排序 (Bubble Sort)

冒泡排序是最基础的交换排序。冒泡排序之所以叫冒泡排序，正是因为这种排序算法的每一个元素都可以像小气泡一样，根据自身大小，一点一点向着数组的一侧移动。

按照冒泡排序的思想，要把相邻的元素两两比较，当一个元素大于右侧相邻元素时，交换它们的位置；当一个元素小于或等于右侧相邻元素时，位置不变。

例如一个有 8 个数字组成的无序序列，进行升序排序。

0	1	2	3	4	5	6	7
5	8	6	3	9	2	1	7
5	8	6	3	9	2	1	7
5	6	8	3	9	2	1	7
5	6	3	8	9	2	1	7
5	6	3	8	9	2	1	7
5	6	3	8	2	9	1	7
5	6	3	8	2	1	9	7
5	6	3	8	2	1	7	9

图 8.2: 冒泡排序第 1 轮

这样一来，元素 9 作为数列中最大的元素，就像是汽水里的小气泡一样，浮到了最右侧。这时，冒泡排序的第 1 轮就结束了。数列最右侧元素 9 的位置可以认为是一个有序区域，有序区域目前只有 1 个元素。

接着进行第 2 轮排序：

0	1	2	3	4	5	6	7
5	6	3	8	2	1	7	9
5	6	3	8	2	1	7	9
5	3	6	8	2	1	7	9
5	3	6	8	2	1	7	9
5	3	6	2	8	1	7	9
5	3	6	2	1	8	7	9
5	3	6	2	1	7	8	9

图 8.3: 冒泡排序第 2 轮

第 2 轮排序结束后，数列右侧的有序区有了 2 个元素。

根据相同的方法，完成剩下的排序：

	0	1	2	3	4	5	6	7
第 3 轮	3	5	2	1	6	7	8	9
第 4 轮	3	2	1	5	6	7	8	9
第 5 轮	2	1	3	5	6	7	8	9
第 6 轮	1	2	3	5	6	7	8	9
第 7 轮	1	2	3	5	6	7	8	9

图 8.4: 冒泡排序第 3 ~ 7 轮

8.2.2 算法分析

冒泡排序是一种稳定排序，值相等的元素并不会打乱原本的顺序。由于该排序算法的每一轮都要遍历所有元素，总共遍历 $n - 1$ 轮。

时间复杂度	空间复杂度	稳定性	设计思想
$O(n^2)$	$O(1)$	稳定	贪心法

表 8.1: 冒泡排序算法分析

冒泡排序

```

1 void bubbleSort(int *arr, int n) {
2     for(int i = 0; i < n; i++) {
3         for(int j = 0; j < n-i-1; j++) {
4             if(arr[j] > arr[j+1]) {
5                 swap(&arr[j], &arr[j+1]);
6             }
7         }
8     }
9 }

```

逆序对

假设数组有 n 个元素，如果 $A[i] > A[j]$, $i < j$ ，那么 $A[i]$ 和 $A[j]$ 就被称为逆序对 (inversion)。

```

1 int countInversions(int *arr, int n) {
2     int cnt = 0;          // 逆序对数
3     for(int i = 0; i < n-1; i++) {
4         for(int j = i+1; j < n; j++) {
5             if(arr[i] > arr[j]) {
6                 cnt++;
7             }
8         }
9     }
10    return cnt;
11 }

```

8.2.3 冒泡排序第一次优化

常规的冒泡排序需要进行 $n - 1$ 轮循环，即使在中途数组已经有序，但是还是会继续剩下的循环。例如当数组是 $\{2, 1, 3, 4, 5\}$ 时，在经过一轮排序后已经变为有序状态，再进行多余的循环就会浪费时间。

为了解决这个问题，可以在每一轮循环中设置一个标志。如果该轮循环中有元素发生过交换，那么就有必要进行下一轮循环。如果没有发生过交换，说明当前数组已经完成排序。

冒泡排序第一次优化

```
1 void bubbleSortOptimize1(int *arr, int n) {
2     for(int i = 0; i < n - 1; i++) {
3         bool hasSwapped = false; // 标记是否发生交换
4         for(int j = 0; j < n - i - 1; j++) {
5             if(arr[j] > arr[j+1]) {
6                 swap(&arr[j], &arr[j+1]);
7                 hasSwapped = true; // 发生交换
8             }
9         }
10        // 该轮未发生交换，已经有序
11        if(!hasSwapped) {
12            return;
13        }
14    }
15 }
```

8.2.4 冒泡排序第二次优化

在经过一次优化后，算法还存在一个问题，例如数组 $\{2, 3, 1, 4, 5, 6\}$ 在经过一轮交换后变为 $\{2, 1, 3, 4, 5, 6\}$ ，但是在下一轮时后面有很多次比较都是多余的，因为并没有产生交换操作。

为了解决这个问题，可以再设置一个标志位，用于记录当前轮所交换的最后一个元素的下标。在下一轮排序中，只需比较到该标志位即可，因此之后的元素在上一轮中没有交换过，在这一轮中也不可能交换了。

冒泡排序第二次优化

```
1 void bubbleSortOptimize2(int *arr, int n) {
2     int len = n - 1;          // 内层循环执行次数
3     for(int i = 0; i < n - 1; i++) {
4         bool hasSwapped = false; // 标记是否发生交换
5         int last = 0;          // 标记最后一次发生交换的位置
6         for(int j = 0; j < len; j++) {
7             if(arr[j] > arr[j+1]) {
8                 swap(&arr[j], &arr[j+1]);
9                 hasSwapped = true; // 发生交换
10                last = j;
11            }
12        }
13        // 该轮未发生交换，已经有序
14        if(!hasSwapped) {
15            return;
16        }
17        len = last;          // 最后一次发生交换的位置
18    }
19 }
```

8.3 选择排序

8.3.1 选择排序 (Selection Sort)

有了冒泡排序为什么还要发明选择排序？冒泡排序有个很大的弊端，就是元素交换次数太多了。

想象一个场景，假设你是一名体育老师，正在指挥一群小学生按照个头从矮到高的顺序排队。采用冒泡排序的方法需要频繁交换相邻学生的位置，同学们心里恐怕会想：“这体育老师是不是有毛病啊？”

在程序运行的世界里，虽然计算机并不会产生什么负面情绪，但是频繁的数组元素交换意味着更多的内存读写操作，严重影响了代码运行效率。

有一个简单的办法，就是每一次找到个子最矮的学生，直接交换到队伍的前面。

例如一个有 8 个数字组成的无序序列，进行升序排序。

	0	1	2	3	4	5	6	7
原数组	5	8	6	3	9	2	1	7
第 1 轮	1	8	6	3	9	2	5	7
第 2 轮	1	2	6	3	9	8	5	7
第 3 轮	1	2	3	6	9	8	5	7
第 4 轮	1	2	3	5	9	8	6	7
第 5 轮	1	2	3	5	6	8	9	7
第 6 轮	1	2	3	5	6	7	9	8
第 7 轮	1	2	3	5	6	7	8	9

图 8.5: 选择排序

8.3.2 算法分析

算法每一轮选出最小值，再交换到左侧的时间复杂度是 $O(n)$ ，一共迭代 $n - 1$ 轮，总的时间复杂度是 $O(n^2)$ 。

由于算法所做的是原地排序，并没有利用额外的数据结构，所以空间复杂度是 $O(1)$ 。

时间复杂度	空间复杂度	稳定性	设计思想
$O(n^2)$	$O(1)$	不稳定	减治法

表 8.2: 选择排序算法分析

选择排序

```
1 void selectionSort(int *arr, int n) {
2     for(int i = 0; i < n-1; i++) {
3         int minIndex = i;
4         for(int j = i+1; j < n; j++) {
5             if(arr[j] < arr[minIndex]) {
6                 minIndex = j;
7             }
8         }
9         if(i != minIndex) {
10             swap(&arr[i], &arr[minIndex]);
11         }
12     }
13 }
```

8.3.3 选择排序优化

选择排序的整体思想是在一个序列当中选出一个最小的元素，和第一个元素交换，然后在剩下的找最小的，和第二个元素交换。这样最终就可以得到一个有序序列。但是为了更加高效，可以每次选择出一个最小值和一个最大值，分别放在

序列的最左和最右边。

选择排序优化

```
1 void selectionSortOptimize(int *arr, int n) {
2     int left = 0;
3     int right = n - 1;
4     while(left < right) {
5         int min = left;
6         int max = right;
7         for(int i = left; i <= right; i++) {
8             if(arr[i] < arr[min]) {
9                 min = i;
10            }
11            if(arr[i] > arr[max]) {
12                max = i;
13            }
14        }
15        swap(&arr[max], &arr[right]);
16        // 考虑特殊情况，最小值在最右位置
17        if(min == right) {
18            min = max;
19        }
20        swap(&arr[min], &arr[left]);
21        left++;
22        right--;
23    }
24 }
```

8.4 插入排序

8.4.1 插入排序 (Insertion Sort)

如何对扑克牌进行排序呢？例如现在手上有红桃 6, 7, 9, 10 这四张牌，已经处于升序排序状态。这时候抓到了一张红桃 8，如何让手上的五张牌重新变成升序呢？

使用冒泡排序？选择排序？恐怕正常人打牌的时候都不会那么做。最自然最简单的方式，是在已经有序的四张牌中找到红桃 8 应该插入的位置，也就是 7 和 9 之间，把红桃 8 插入进去。

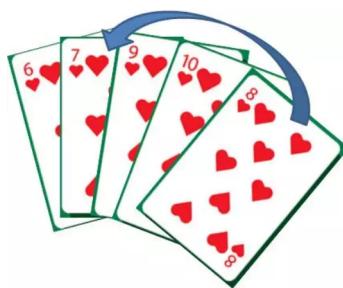


图 8.6: 理牌

例如一个有 8 个数字组成的无序序列，进行升序排序。

	0	1	2	3	4	5	6	7
原数组	5	8	6	3	9	2	1	7
第 1 轮	5	8	6	3	9	2	1	7
第 2 轮	5	6	8	3	9	2	1	7
第 3 轮	3	5	6	8	9	2	1	7
第 4 轮	3	5	6	8	9	2	1	7
第 5 轮	2	3	5	6	8	9	1	7
第 6 轮	1	2	3	5	6	8	9	7
第 7 轮	1	2	3	5	6	7	8	9

图 8.7: 插入排序

8.4.2 算法分析

插入排序要进行 $n - 1$ 轮，每一轮在最坏情况下的比较复制次数分别是 1 次、2 次、3 次、4 次... 一直到 $n - 1$ 次，所以最坏时间复杂度是 $O(n^2)$ 。

至于空间复杂度，由于插入排序是在原地进行排序，并没有引入额外的数据结构，所以空间复杂度是 $O(1)$ 。

时间复杂度	空间复杂度	稳定性	设计思想
$O(n^2)$	$O(1)$	稳定	减治法

表 8.3: 插入排序算法分析

插入排序

```
1 void insertionSort(int *arr, int n) {
2     for(int i = 1; i < n; i++) {
3         int temp = arr[i];
4         int j = i - 1;
5         while(j >= 0 && temp < arr[j]) {
6             arr[j+1] = arr[j];
7             j--;
8         }
9         arr[j+1] = temp;
10    }
11 }
```

8.4.3 折半插入排序 (Binary Insertion Sort)

折半插入排序是对插入排序的改进，其过程就是不断依次将元素插入前面已经排好序的序列中，在寻找插入点时采用了折半查找。

折半插入排序

```
1 void binaryInsertionSort(int *arr, int n) {
2     for(int i = 1; i < n; i++) {
3         int temp = arr[i];
4         int start = 0;
5         int end = i - 1;
6         while(start <= end) {
7             int mid = start + (end - start) / 2;
8             if(arr[mid] > temp) {
9                 end = mid - 1;
10            } else {
11                start = mid + 1;
12            }
13        }
14        int j;
15        for(j = i - 1; j > end; j--) {
16            arr[j+1] = arr[j];
17        }
18        arr[j+1] = temp;
19    }
20 }
```

8.5 鸡尾酒排序

8.5.1 鸡尾酒排序 (Cocktail Sort)

对一个无序数组 {2, 3, 4, 5, 6, 7, 8, 1} 进行升序排序，如果按照冒泡排序，步骤如下：

	0	1	2	3	4	5	6	7
原数组	2	3	4	5	6	7	8	1
第 1 轮	2	3	4	5	6	7	1	8
第 2 轮	2	3	4	5	6	1	7	8
第 3 轮	2	3	4	5	1	6	7	8
第 4 轮	2	3	4	1	5	6	7	8
第 5 轮	2	3	1	4	5	6	7	8
第 6 轮	2	1	3	4	5	6	7	8
第 7 轮	1	2	3	4	5	6	7	8

图 8.8: 冒泡排序

从 2 到 8 已经是有序了，只有元素 1 的位置不对，却还要进行 7 轮排序，这也太憋屈了吧！

鸡尾酒排序正是用于解决这种问题的。鸡尾酒排序又叫快乐小时排序，它基于冒泡排序做了一点小小的优化。

鸡尾酒排序的第一轮与冒泡排序相同，从左向右比较和交换。

0	1	2	3	4	5	6	7
2	3	4	5	6	7	1	8

图 8.9: 鸡尾酒排序第 1 轮

第二轮开始不一样了，反过来从右向左比较和交换。

0	1	2	3	4	5	6	7
2	3	4	5	6	7	1	8
2	3	4	5	6	1	7	8
2	3	4	5	1	6	7	8
2	3	4	1	5	6	7	8
2	3	1	4	5	6	7	8
2	1	3	4	5	6	7	8
1	2	3	4	5	6	7	8

图 8.10: 鸡尾酒排序第 2 轮

第二轮结束后，此时虽然已经有序，但是算法并没有结束。

第三轮再次从左向右比较和交换，在此过程中，没有元素发生交换，证明已经有序，排序结束。

鸡尾酒排序的过程就像钟摆一样，左右来回比较和交换。本来冒泡要用 8 轮排序的场景，鸡尾酒用 3 轮就解决了。

鸡尾酒排序的优点是能够在大部分元素已经有序的情况下，减少排序的回合数；而缺点也很明显，就是代码量几乎扩大了一倍。

鸡尾酒排序

```
1 void cocktailSort1(int *arr, int n) {  
2     for(int i = 0; i < n / 2; i++) {  
3         // 从左向右  
4         bool isSorted = true;    // 标记当前轮是否有序  
5         for(int j = i; j < n - i - 1; j++) {
```

```

6         if(arr[j] > arr[j+1]) {
7             swap(&arr[j], &arr[j+1]);
8             isSorted = false;    // 发生交换
9         }
10    }
11    if(isSorted) {
12        break;
13    }
14
15    // 从右向左
16    isSorted = true;
17    for(int j = n - i - 1; j > i; j--) {
18        if(arr[j] < arr[j-1]) {
19            swap(&arr[j], &arr[j-1]);
20            isSorted = false;
21        }
22    }
23    if(isSorted) {
24        break;
25    }
26 }
27 }

```

8.5.2 鸡尾酒排序优化

类似冒泡排序的优化方法，鸡尾酒排序也可以记录每一轮排序后最后一次元素交换的位置。但是对于双向的鸡尾酒排序而言，需要设置两个边界值。

鸡尾酒排序优化

```

1 void cocktailSort2(int *arr, int n) {
2     int lastLeft = 0;        // 左侧最后一次交换位置
3     int lastRight = 0;       // 右侧最后一次交换位置
4     int leftBorder = 0;      // 无序区左边界
5     int rightBorder = n - 1;  // 无序区右边界

```

```

6
7   for(int i = 0; i < n / 2; i++) {
8       // 从左向右
9       bool isSorted = true;    // 标记当前轮是否有序
10      for(int j = leftBorder; j < rightBorder; j++) {
11          if(arr[j] > arr[j+1]) {
12              swap(&arr[j], &arr[j+1]);
13              isSorted = false;  // 发生交换
14              lastRight = j;
15          }
16      }
17      if(isSorted) {
18          break;
19      }
20      rightBorder = lastRight;
21
22      // 从右向左
23      isSorted = true;
24      for(int j = n - i - 1; j > i; j--) {
25          if(arr[j] < arr[j-1]) {
26              swap(&arr[j], &arr[j-1]);
27              isSorted = false;
28              lastLeft = j;
29          }
30      }
31      if(isSorted) {
32          break;
33      }
34      leftBorder = lastLeft;
35  }
36 }

```

8.6 归并排序

8.6.1 归并排序 (Merge Sort)

归并排序算法采用分治法：

1. 分解：将序列每次折半划分。
2. 合并：将划分后的序列两两按序合并。

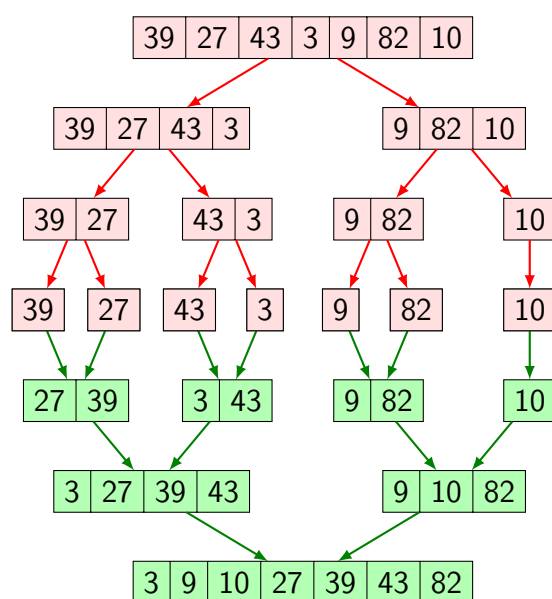


图 8.11: 归并排序

8.6.2 算法分析

归并排序每次将数组折半对分，一共分了 $\log n$ 次，每一层进行合并操作的运算量是 n ，所以时间复杂度为 $O(n \log n)$ 。归并排序的速度仅次于快速排序。

时间复杂度	空间复杂度	稳定性	设计思想
$O(n \log n)$	$O(n)$	稳定	分治法

表 8.4: 归并排序算法分析

```

1 void merge(int *arr, int start, int mid, int end, int *temp) {
2     int i = start;
3     int j = mid + 1;
4     int k = 0;
5
6     while(i <= mid && j <= end) {
7         if(arr[i] <= arr[j]) {
8             temp[k++] = arr[i++];
9         } else {
10            temp[k++] = arr[j++];
11        }
12    }
13
14    while(i <= mid) {
15        temp[k++] = arr[i++];
16    }
17    while(j <= end) {
18        temp[k++] = arr[j++];
19    }
20
21    for(int i = 0; i < k; i++) {
22        arr[start+i] = temp[i];
23    }
24 }
25
26 void mergeSort(int *arr, int start, int end, int *temp) {
27     if(start < end) {
28         int mid = start + (end - start) / 2;
29         mergeSort(arr, start, mid, temp);
30         mergeSort(arr, mid+1, end, temp);
31         merge(arr, start, mid, end, temp);
32     }
33 }

```


8.7 快速排序

8.7.1 快速排序 (Quick Sort)

快速排序是很重要的算法，与傅里叶变换等算法并称二十世纪十大算法。

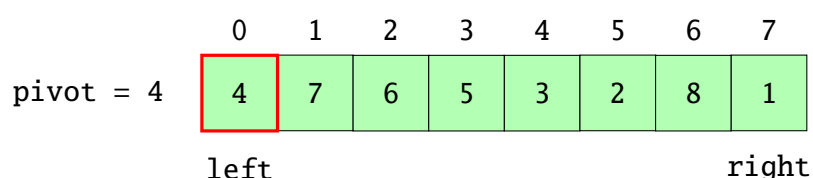
快速排序之所以快，是因为它使用了分治法。快速排序在每一轮挑选一个基准 (pivot) 元素，并让其它比它小的元素移动到数列一边，比它大的元素移动到数列的另一边，从而把数列拆解成了两个部分。

选择基准元素最简单的方式是选择数列的第一个元素。这种选择在绝大多数情况下是没有问题的，但是如果对一个原本逆序的数列进行升序排序，整个数列并没有被分成一半，每一轮仅仅确定了基准元素的位置。这种情况下数列第一个元素要么是最小值，要么是最大值，根本无法发挥分治法的优势。在这种极端情况下，快速排序需要进行 n 轮，时间复杂度退化成了 $O(n^2)$ 。

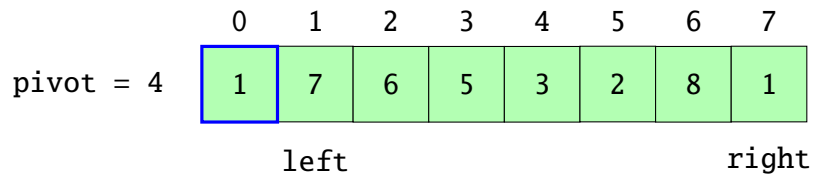
如何避免这种极端情况呢？可以不选择数列的第一个元素，而是随机选择一个元素作为基准元素。这样一来，即使是在数列完全逆序的情况下，也可以有效地将数列分成两部分。当然，即使是随机选择，每一次也有极小的几率选到数列的最大值或最小值，同样会对分治造成一定影响。

确定了基准值后，如何实现将小于基准的元素都移动到基准值一边，大于基准值的都移动到另一边呢？

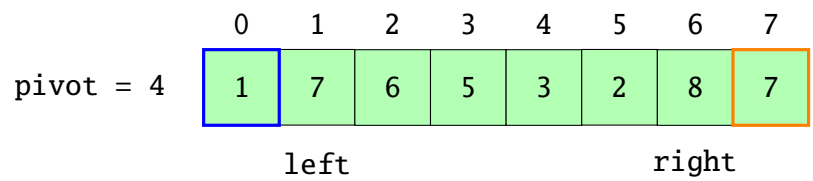
例如一个有 8 个数字组成的无序序列，进行升序排序。选定基准元素 pivot，设置两个指针 left 和 right，指向数列的最左和最右两个元素。



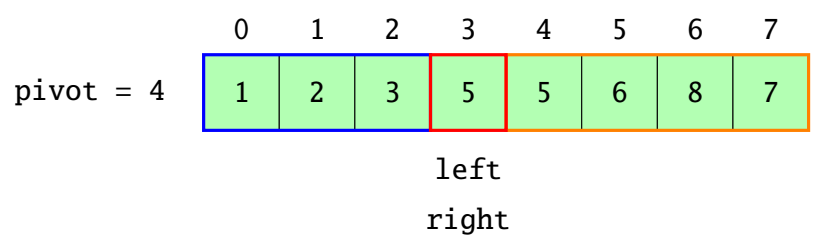
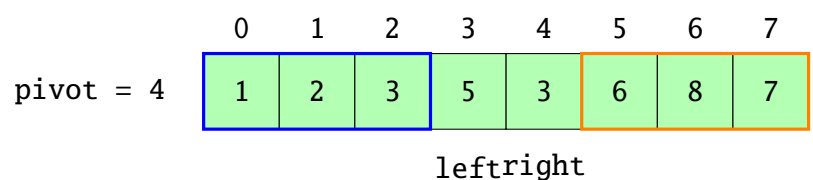
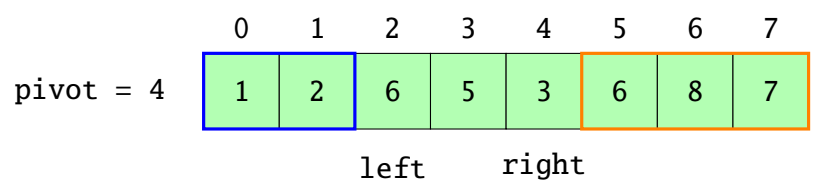
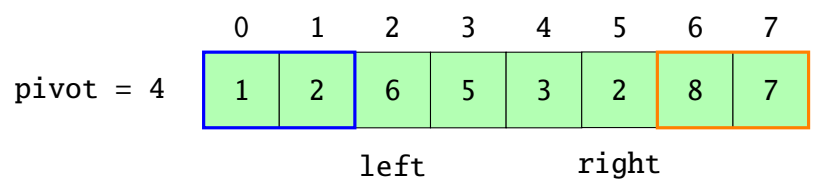
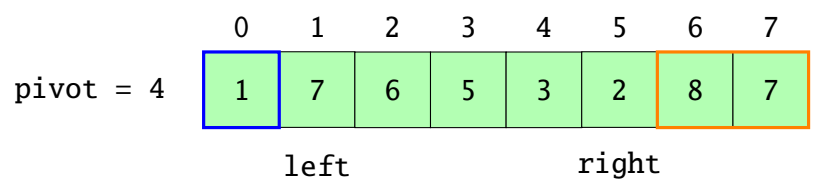
从 right 指针开始，把指针所指向的元素和基准元素做比较。如果比 pivot 大，则 right 指针向左移动；如果比 pivot 小，则把 right 所指向的元素填入 left 指针所指向的位置，同时 left 向右移动一位。



接着，切换到 left 指针进行比较，把指针所指向的元素和基准元素做比较。如果小于 pivot，则 left 指针向右移动；如果大于 pivot，则把 left 所指向的元素填入 right 指针所指向的位置，同时 right 向左移动一位。



重复之前的步骤继续排序：



当 left 和 right 指针重合在同一位置的时候，把之前的 pivot 元素的值填入该重合的位置。此时数列左边的元素都小于基准元素，数列右边的元素都大于基准元素。

8.7.2 算法分析

分治法的思想下，原数列在每一轮被拆分成两部分，每一部分在下一轮又被拆分成两部分，直到不可再分为止。这样平均情况下需要 $\log n$ 轮，因此快速排序算法的平均时间复杂度是 $O(n \log n)$ 。

时间复杂度	空间复杂度	稳定性	设计思想
$O(n \log n) \sim O(n^2)$	$O(\log n) \sim O(n)$	不稳定	分治法

表 8.5: 快速排序算法分析

快速排序

```
1 void quickSort(int *arr, int start, int end) {
2     if(start < end) {
3         int i = start;
4         int j = end;
5         int pivot = arr[start];
6
7         while(i < j) {
8             while(i < j && arr[j] > pivot) {
9                 j--;
10            }
11            if(i < j) {
12                arr[i] = arr[j];
13                i++;
14            }
15            while(i < j && arr[i] < pivot) {
16                i++;
17            }
18            if(i < j) {
19                arr[j] = arr[i];
```

```
20         j--;  
21     }  
22 }  
23 arr[i] = pivot;  
24 quickSort(arr, start, i-1);  
25 quickSort(arr, i+1, end);  
26 }  
27 }
```

8.8 计数排序

8.8.1 计数排序 (Counting Sort)

基于比较的排序算法的最优下界为 $\Omega(n \log n)$ 。计数排序是一种不基于比较的排序算法，而是利用数组下标来确定元素的正确位置。

遍历数列，将每一个整数按照其值对号入座，对应数组下标的元素加 1。数组的每一个下标位置的值，代表了数列中对应整数出现的次数。有了这个统计结果，直接遍历数组，输出数组元素的下标值，元素的值是多少就输出多少次。

从功能角度，这个算法可以实现整数的排序，但是也存在一些问题。如果只以最大值来决定统计数组的长度并不严谨，例如数列 {95, 94, 91, 98, 99, 90, 99, 93, 91, 92}，这个数列的最大值是 99，但最小值是 90。如果创建长度为 100 的数组，前面的从 0 到 89 的空间位置都浪费了。

因此，不应再以数列的 $\max + 1$ 作为统计数组的长度，而是以数列 $\max - \min + 1$ 作为统计数组的长度。同时，数列的最小值作为一个偏移量，用于统计数组的对号入座。

计数排序适用于一定范围的整数排序，在取值范围不是很大的情况下，它的性能甚至快过那些 $O(n \log n)$ 的排序算法。

计数排序

```
1 void countingSort(int *arr, int n) {  
2     int max = arr[0];  
3     int min = arr[0];  
4     for(int i = 1; i < n; i++) {  
5         if(arr[i] > max) {  
6             max = arr[i];  
7         }  
8         if(arr[i] < min) {
```

```
9         min = arr[i];
10     }
11 }
12
13 int range = max - min + 1;
14 int table[range];
15 memset(table, 0, sizeof(table));
16
17 for(int i = 0; i < n; i++) {
18     table[arr[i] - min]++;
19 }
20
21 int cnt = 0;
22 for(int i = 0; i < range; i++) {
23     while(table[i]--) {
24         arr[cnt++] = i + min;
25     }
26 }
27 }
```

8.9 桶排序

8.9.1 桶排序 (Bucket Sort)

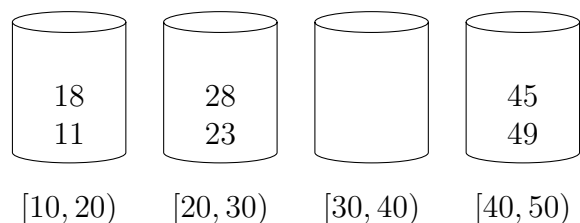
桶排序是计数排序的扩展版本。计数排序可以看成每个桶只存储相同元素，而桶排序每个桶存储一定范围的元素。

每一个桶代表一个区间范围，里面可以承载一个或多个元素。通过划分多个范围相同的区间，将每个子区间自排序，最后合并。桶排序需要尽量保证元素分散均匀，否则当所有数据集中在同一个桶中时，桶排序失效。

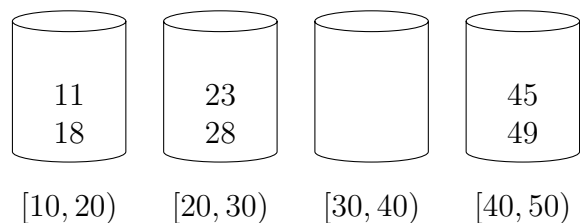
例如一个待排序的序列：

0	1	2	3	4	5
18	11	28	45	23	49

确定桶的个数与每个桶的取值范围，遍历待排序序列，将元素放入对应的桶中：



分别对每个桶中的元素进行排序：



将桶中的元素按顺序赋值到原始数组中：

0	1	2	3	4	5
11	18	23	28	45	49

创建桶的数量取决于数据的区间范围，一般创建桶的数量等于待排序的元素数量，每个桶的区间跨度为：

$$\frac{\max - \min}{\text{buckets} - 1}$$

桶排序

```
1 public static void bucketSort(int[] arr) {
2     int max = arr[0];
3     int min = arr[0];
4     for(int i = 1; i < arr.length; i++) {
5         if(arr[i] > max) {
6             max = arr[i];
7         }
8         if(arr[i] < min) {
9             min = arr[i];
10        }
11    }
12    int diff = max - min;
13
14    // 创建桶
15    int bucketNum = arr.length;
16    List<LinkedList<Integer>> buckets
17        = new ArrayList<>(bucketNum);
18    for(int i = 0; i < bucketNum; i++) {
19        buckets.add(new LinkedList<>());
20    }
21
22    // 遍历数组，把元素放入对应桶中
23    for(int i = 0; i < arr.length; i++) {
24        // 计算当前元素所放置的桶号
25        int num = (arr[i]-min) / (diff / (bucketNum-1));
26        buckets.get(num).add(arr[i]);
27    }
28
29    // 桶内排序
```



```
30     for(int i = 0; i < bucketNum; i++) {
31         Collections.sort(buckets.get(i));
32     }
33
34     // 取出元素
35     int cnt = 0;
36     for(LinkedList<Integer> list : buckets) {
37         for(int data : list) {
38             arr[cnt++] = data;
39         }
40     }
41 }
```

8.10 基数排序

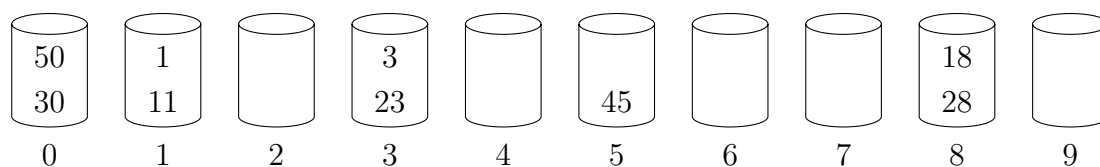
8.10.1 基数排序 (Radix Sort)

基数排序可以看作是桶排序的扩展，主要思想是将整数按位划分。基数排序需要准备 10 个桶，分别代表 0~9，根据整数个位数字的数值将元素放入对应的桶中，之后按照输入赋值到原序列中，再依次对十位、百位等进行同样的操作。

例如一个待排序的序列：

0	1	2	3	4	5	6	7	8
3	1	18	11	28	45	23	50	30

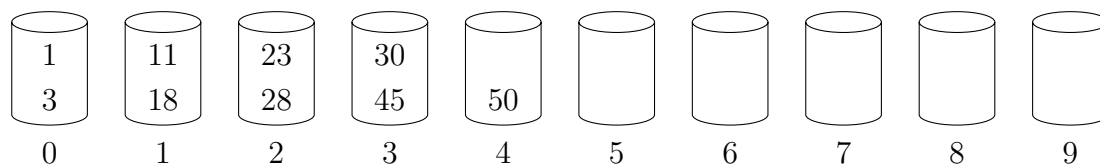
根据个位数值放入对应的桶中：



将桶中的元素按顺序赋值到原始数组中：

0	1	2	3	4	5	6	7	8
50	30	1	11	3	23	45	18	28

根据十位数值放入对应的桶中：



将桶中的元素按顺序赋值到原始数组中：

0	1	2	3	4	5	6	7	8
1	3	11	18	23	28	30	45	50

```

1 public static void radixSort(int[] arr) {
2     int max = arr[0];
3     for(int i = 1; i < arr.length; i++) {
4         if(arr[i] > max) {
5             max = arr[i];
6         }
7     }
8
9     int bucketNum = 10;
10    // 从个位开始
11    for(int exp = 1; max / exp > 0; exp *= 10) {
12        // 创建桶
13        List<LinkedList<Integer>> buckets
14            = new ArrayList<>(bucketNum);
15        for(int i = 0; i < bucketNum; i++) {
16            buckets.add(new LinkedList<>());
17        }
18
19        // 把元素放到对应桶中
20        for(int data : arr) {
21            int num = (data / exp) % 10;
22            buckets.get(num).add(data);
23        }
24
25        // 按顺序取出元素
26        int cnt = 0;
27        for(LinkedList<Integer> bucket : buckets) {
28            for(int data : bucket) {
29                arr[cnt++] = data;
30            }
31        }
32    }
33 }

```

8.11 珠排序

8.11.1 珠排序 (Bead Sort)

珠排序的算法和算盘非常类似。算盘上有许多圆圆的珠子被串在细杆上，如果把算盘竖起来，算盘上的小珠子会在重力的作用下滑到算盘底部。其中有一个很神奇的细节，如果统计每一横排珠子的数量，下落之后每一排珠子数量正好是下落前珠子数量的升序排序。

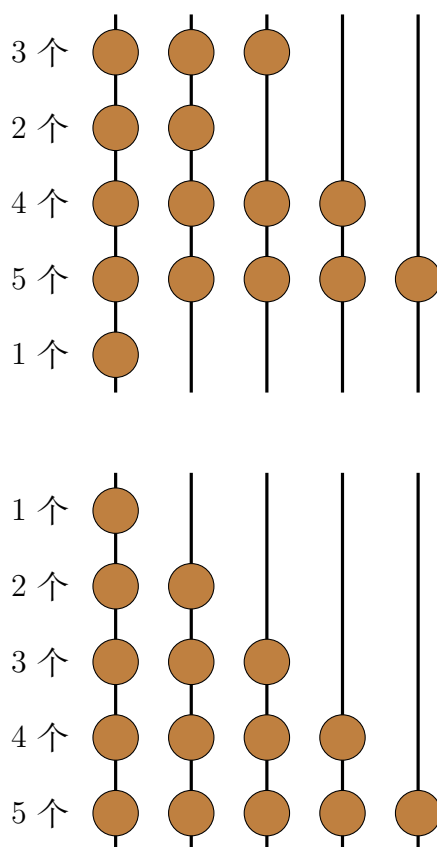


图 8.12: 珠子下落

通过模拟珠子下落的元素可以对一组正整数进行排序。使用二维数组来模拟算盘，有珠子的位置设为 1，没有珠子的位置设为 0。例如一个有 5 个数字组成的无序整型数组 {3, 2, 4, 5, 1} 就可以转化为以下的二维数组：

1	1	1	0	0
1	1	0	0	0
1	1	1	1	0
1	1	1	1	1
1	0	0	0	0

表 8.6: 模拟算盘

接下来模拟算盘珠子掉落的过程，让所有的元素 1 都落到二维数组的最底部。

1	0	0	0	0
1	1	0	0	0
1	1	1	0	0
1	1	1	1	0
1	1	1	1	1

最后把掉落后的算盘转换为一维有序数组 {1, 2, 3, 4, 5}。

8.12 猴子排序

8.12.1 猴子排序 (Bogo Sort)

听说过“猴子和打字机”的理论吗？

无限猴子定理 (Infinite Monkey Theorem) 与薛定谔的猫、电车实验等并居十大思想实验，所谓思想实验即用想象力去进行，而在现实中基本无法去实现的实验。

无限猴子定理讲的是如果让一只猴子在打字机上胡乱打字，只要有无限的时间，总有一天可以恰好打出莎士比亚的著作。如果让无限只猴子在无限的空间、无限的时间里不停地敲打打字机，总有一天可以完整打出一本《哈姆雷特》，甚至是可以打出所有可能的文章。



图 8.13: 无限猴子定理

这个看似不可能的事情，却可以用现有数学原理被推导出来。但在现实中往往被认为是无法实现的，因为人们认为“无限”这个条件通常无法被满足。根据概率论证，即使可观测宇宙中充满了猴子一直不停地打字，能够打出一部《哈姆雷特》的概率仍然小于 $\frac{1}{10^{183900}}$ 。

无限猴子定理同样可以用在排序中。如果给数组随机排列顺序，每一次排列之后验证数组是否有序，只要次数足够多，总有一次数组刚好被随机成有序数组。可是要想真的随机出有序数列，恐怕要等到猴年马月了。

Chapter 9 树

9.1 树

9.1.1 树 (Tree)

许多逻辑关系并不是简单的线性关系，在实际场景中，常常存在着一对多，甚至多对多的情况。树和图就是典型的非线性数据结构。

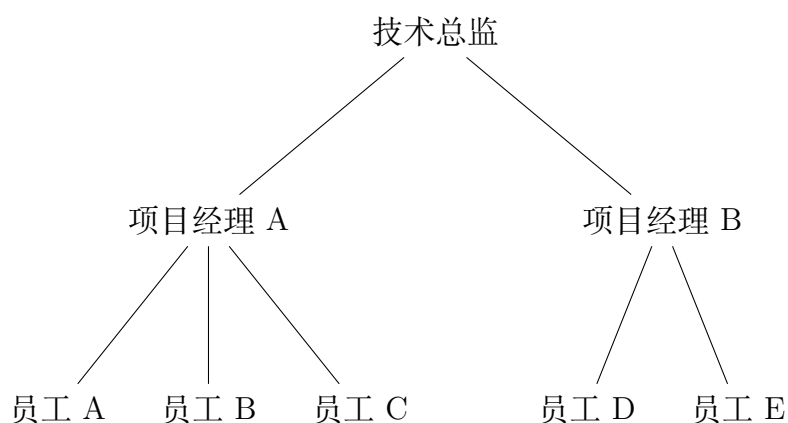


图 9.1: 职级关系

这些结构都像自然界中的树一样，从同一个根衍生出许多枝干，再从每一个枝干衍生出许多更小的枝干，最后衍生出更多的叶子。

树是由 n ($n \geq 0$) 个有限节点组成的一个具有层次关系的集合，当 $n = 0$ 时称为空树。

在任意一个非空树中，有以下特点：

- 有且仅有一个特定的结点称为根 (root)。
- 当 $n > 1$ 时,其余结点可分为 m ($m > 0$) 个互不相交的有限集 T_1, T_2, \dots, T_m , 其中每一个集合本身又是一棵树, 并且称为根的子树 (subtree)。

9.1.2 树的术语

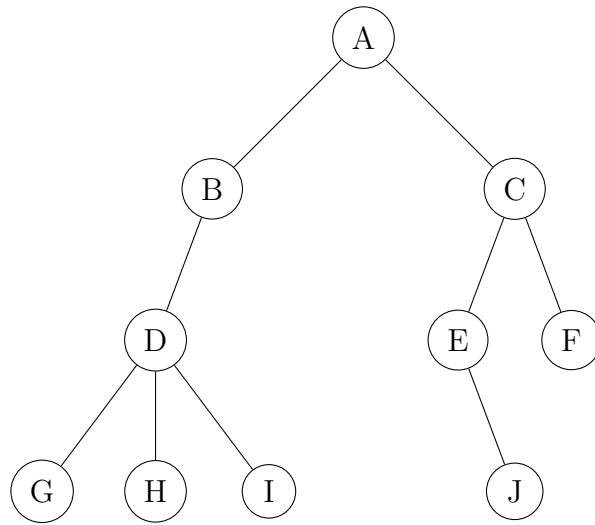


图 9.2: 树

- 根：没有父结点（parent）的结点。
- 内部结点（internal node）：至少有一个子结点（child）的结点。
- 外部结点（external node） / 叶子结点（leaf node）：没有子结点的结点。
- 度（degree）：结点分支的个数。
- 路径（path）：从根结点到树中某结点经过的分支构成了路径。
- 祖先结点（ancestors）：包含父结点、父结点的父结点等。
- 子孙结点（descendants）：包含子结点、子结点的子结点等。
- 深度（depth） / 高度（height）：最大层级数。

9.2 二叉树

9.2.1 二叉树 (Binary Tree)

二叉树是树的一种特殊形式。二叉树的每个结点最多有两个孩子结点，即最多有 2 个，也可能只有 1 个，或者没有孩子结点。

二叉树结点的两个孩子结点，分别被称为左孩子 (left child) 和右孩子 (right child)。这两个孩子结点的顺序是固定的，不能颠倒或混淆。

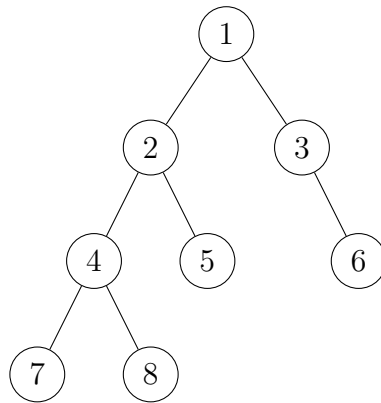


图 9.3: 二叉树

二叉树还有几种特殊的形式：

左斜树 (left skew tree) / 右斜树 (right skew tree)

只有左子树或只有右子树的二叉树。

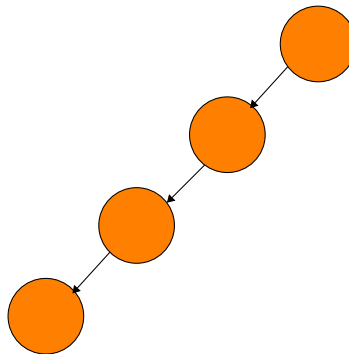


图 9.4: 左斜树

满二叉树 (full binary tree)

所有非叶子结点都存在左右孩子，并且所有叶子结点都在同一层。

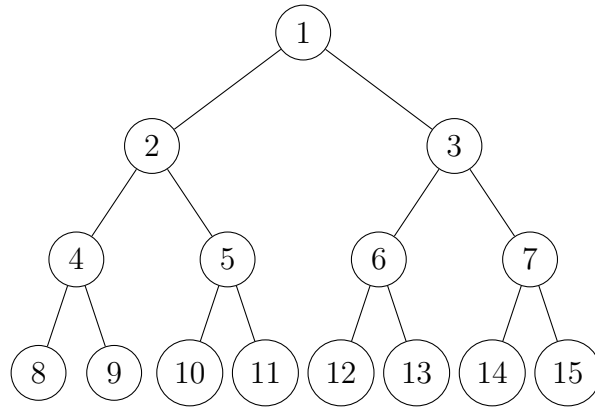


图 9.5: 满二叉树

完全二叉树

对于一个有 n 个结点的二叉树，按层级顺序编号，则所有结点的编号从 1 到 n ，完全二叉树所有结点和同样深度的满二叉树的编号从 1 到 n 的结点位置相同。简单来说，就是除最后一层外，其它各层的结点数都达到最大，并且最后一层从右向左连续缺少若干个结点。

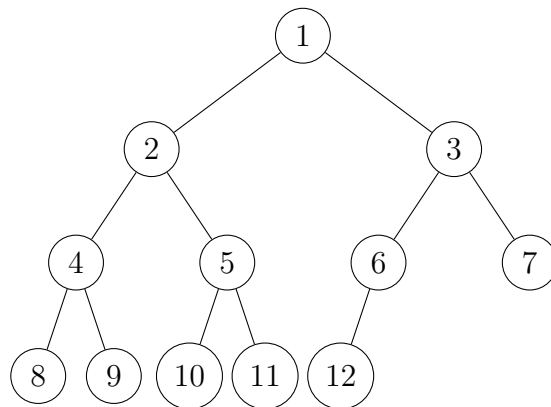


图 9.6: 完全二叉树

9.2.2 二叉树的存储结构

二叉树既可以通过链式存储，也可以使用数组存储：

链式存储结构

一个结点最多可以指向左右两个孩子结点，所以二叉树的每一个结点包含三个部分：

- 存储数据的数据域 data
- 指向左孩子的指针 left
- 指向右孩子的指针 right

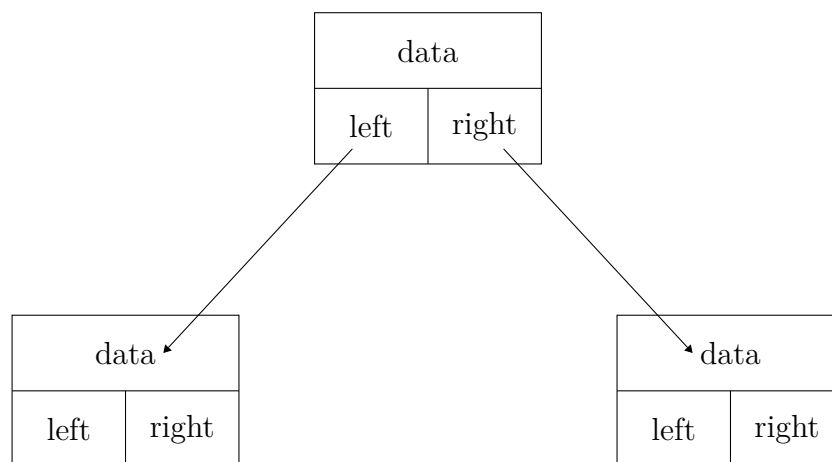


图 9.7: 链式存储结构

数组存储

按照层级顺序把二叉树的结点放到数组中对应的位置上。如果某一结点的左孩子或右孩子空缺，则数组的相应位置也空出来。

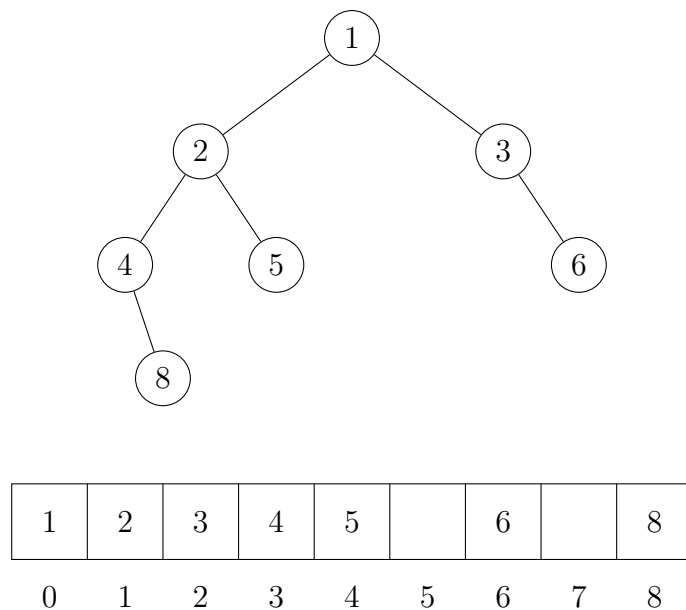


图 9.8: 数组存储

采用数组存储可以更方便地定位二叉树的孩子结点和父结点。假设一个父结点的下标是 parent ，那么它的左孩子结点的下标就是 $2 * \text{parent} + 1$ ，右孩子结点的下标就是 $2 * \text{parent} + 2$ 。反过来，假设一个左孩子结点的下标是 leftChild ，那么它的父结点的下标就是 $(\text{leftChild} - 1) / 2$ 。

但是，对于一个稀疏的二叉树来说，用数组表示法是非常浪费空间的。对于一种特殊的完全二叉树——二叉堆而言，就是使用数组进行存储的。

9.3 二叉树的遍历

9.3.1 二叉树的遍历

在计算机程序中，遍历（traversal）本身是一个线性操作，所以遍历同样具有线性结构的数组或链表是一件轻而易举的事情。

反观二叉树，是典型的非线性数据结构，遍历时需要把非线性关联的结点转化成一个线性的序列，以不同的方式来遍历，遍历出的序列顺序也不同。

二叉树的遍历方式分为 4 种：

1. 前序遍历（pre-order）：访问根结点，遍历左子树，遍历右子树。
2. 中序遍历（in-order）：遍历左子树，访问根结点，遍历右子树。
3. 后序遍历（post-order）：遍历左子树，遍历右子树，访问根结点。
4. 层次遍历（level-order）：按照从根结点到叶子结点的层次关系，一层一层横向遍历。

9.3.2 前序遍历

二叉树的前序遍历，首先访问根结点然后遍历左子树，最后遍历右子树。在遍历左、右子树时，仍然先访问根结点，然后遍历左子树，最后遍历右子树，如果结点为空则返回。

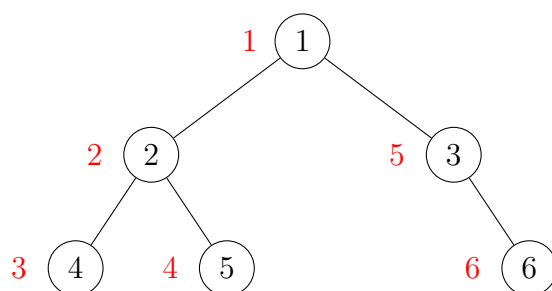


图 9.9: 前序遍历

前序遍历

```
1 void preOrder(BST *root) {  
2     if(!root) {  
3         return;  
4     }  
5     printf("%d ", root->data);  
6     preOrder(root->left);  
7     preOrder(root->right);  
8 }
```

9.3.3 中序遍历

二叉树的中序遍历，首先遍历左子树，然后访问根结点，最后遍历右子树，如果结点为空则返回。

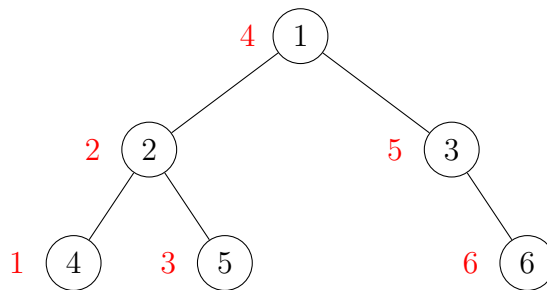


图 9.10: 中序遍历

中序遍历

```
1 void inOrder(BST *root) {  
2     if(!root) {  
3         return;  
4     }  
5     inOrder(root->left);  
6     printf("%d ", root->data);  
7     inOrder(root->right);  
8 }
```

9.3.4 后序遍历

二叉树的后序遍历，首先遍历左子树，然后遍历右子树，最后访问根结点，如果结点为空则返回。

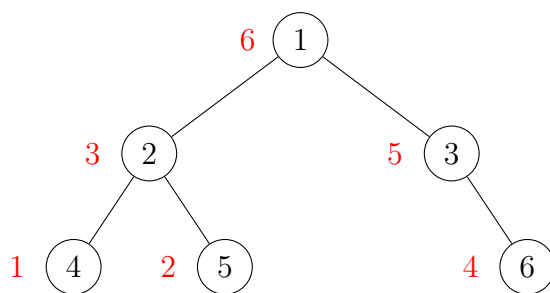


图 9.11: 后序遍历

后序遍历

```
1 void postOrder(BST *root) {  
2     if(!root) {  
3         return;  
4     }  
5     postOrder(root->left);  
6     postOrder(root->right);  
7     printf("%d ", root->data);  
8 }
```

9.3.5 二叉树遍历非递归实现

绝大多数可以用递归解决的问题，其实都可以用栈来解决，因为递归和栈都有回溯的特性。

以二叉树的中序遍历为例。当遇到一个结点时，就把它入栈，并去遍历它的左子树。当左子树遍历结束后，从栈顶弹出这个结点并访问它，然后按其右指针再去

中序遍历该结点的右子树。

中序遍历（非递归）

```
1 public void inOrderNonRecursive(BSTNode node) {
2     Stack s = new Stack();
3     while(node != null || !s.empty()) {
4         // 一直向左并将沿途结点压入堆栈
5         while(node != null) {
6             s.push(node);
7             node = node.left;
8         }
9         if(!s.empty()) {
10            node = s.pop();           // 结点弹出堆栈
11            System.out.println(node.data); // 访问结点
12            node = node.right;        // 转向右子树
13        }
14    }
15 }
```

9.3.6 层次遍历

二叉树同一层次的结点之间是没有直接关联的，需要队列来辅助完成层序遍历。

层次遍历从根结点开始首先将根结点入队，然后开始循环执行以下操作直到队列为空：结点出队、访问该结点、其左右儿子入队。

层次遍历

```
1 public void levelOrder(BSTNode node) {
2     if(node == null) {
3         return;
4     }
5 }
```



```
6   Queue q = new Queue();
7   q.enqueue(node);
8   while(!q.empty()) {
9       node = q.dequeue();
10      System.out.println(node.data);    // 访问结点
11      if(node.left != null) {
12          q.enqueue(node.left);
13      }
14      if(node.right != null) {
15          q.enqueue(node.right);
16      }
17  }
18 }
```

9.4 二叉搜索树

9.4.1 二叉搜索树 (Binary Search Tree)

二叉搜索树，也称二叉查找树或二叉排序树，可以是一棵空树。

如果不为空树，那么二叉搜索树满足以下性质：

1. 非空左子树的所有结点的值小于其根结点的值。
2. 非空右子树的所有结点的值大于其根结点的值。
3. 左、右子树均是二叉搜索树。

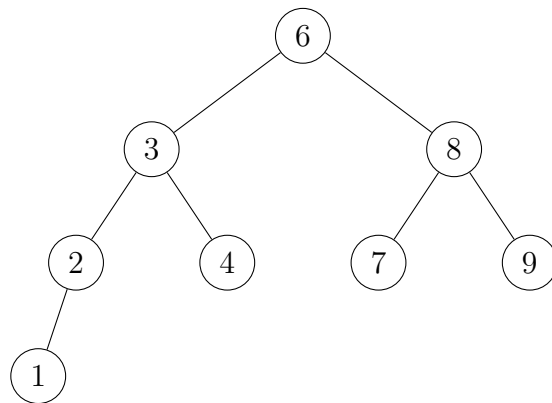


图 9.12: 二叉搜索树

9.4.2 查找结点

在二叉搜索树中查找一个元素从根结点开始，如果树为空，返回 NULL。

如果树不为空，则将根结点的值和被查找的 key 值进行比较：

1. 如果 key 值小于根结点的值，只需在左子树中继续查找。
2. 如果 key 值大于根结点的值，只需在右子树中继续查找。
3. 如果 key 值与根结点的值相等，查找成功。

查找结点（递归）

```
1 Node* search(Node *root, dataType val) {
2     if(!root) {
3         return NULL;
4     }
5     if(val == root->data) {
6         return root;
7     } else if(val < root->data) {
8         return search(root->left, val);
9     } else {
10        return search(root->right, val);
11    }
12 }
```

由于非递归函数的执行效率高，可将尾递归（在函数最后才使用递归返回）的函数改为迭代函数。

查找结点（迭代）

```
1 Node* search(Node *root, dataType val) {
2     if(!root) {
3         return NULL;
4     }
5     while(root) {
6         if(root->data == val) {
7             return root;
8         } else if(val < root->data) {
9             root = root->left;
10        } else {
11            root = root->right;
12        }
13    }
14 }
```

9.4.3 查找最小值和最大值

二叉搜索树中，最小值一定在树的最左分枝的叶子结点上，最大值一定在树的最右分枝的叶子结点上。

查找最小值（递归）

```
1 Node* findMin(Node *root) {  
2     if(!root) {  
3         return NULL;  
4     } else if(!root->left) {  
5         return root;  
6     } else {  
7         return findMin(root->left);    //沿左分枝继续查找  
8     }  
9 }
```

查找最大值（迭代）

```
1 Node* findMax(Node *root) {  
2     if(!root) {  
3         return NULL;  
4     }  
5     while(root->right) {  
6         root = root->right;  
7     }  
8     return root;  
9 }
```

9.4.4 插入结点

在二叉搜索树中插入结点与查找的算法相似，需要找到插入的位置并将新结点插入。

插入结点

```
1 BST* insert(BST *root, dataType val) {  
2     // 空树，插入结点设为树根  
3     if(!root) {  
4         return init(val);  
5     }  
6     if(val < root->data) {  
7         root->left = insert(root->left, val);  
8     } else {  
9         root->right = insert(root->right, val);  
10    }  
11    return root;  
12 }
```

9.5 哈夫曼树

9.5.1 哈夫曼树 (Huffman Tree)

树的每一个结点都可以拥有自己的权值 (weight)，假设二叉树有 n 个叶子结点，每个叶子结点都带有权值 w_k ，从根结点到每个叶子结点的长度为 l_k ，则树的带权路径长度 (WPL, Weighted Path Length) 为：

$$WPL = \sum_{k=1}^n w_k l_k$$

哈夫曼树是由麻省理工学院的哈夫曼博士于 1952 年发明的，哈夫曼树是在叶子结点和权重确定的情况下，带权路径长度最小的二叉树，也被称为最优二叉树。

例如，有五个叶子结点，它们的权值为 $\{1, 2, 3, 4, 5\}$ ，用此权值序列可以构造出形状不同的多个二叉树。

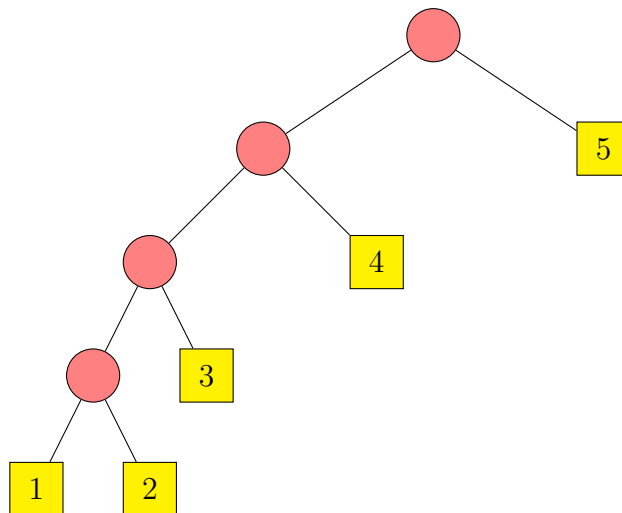


图 9.13: $WPL = 5 * 1 + 4 * 2 + 3 * 3 + 2 * 4 + 1 * 4 = 34$

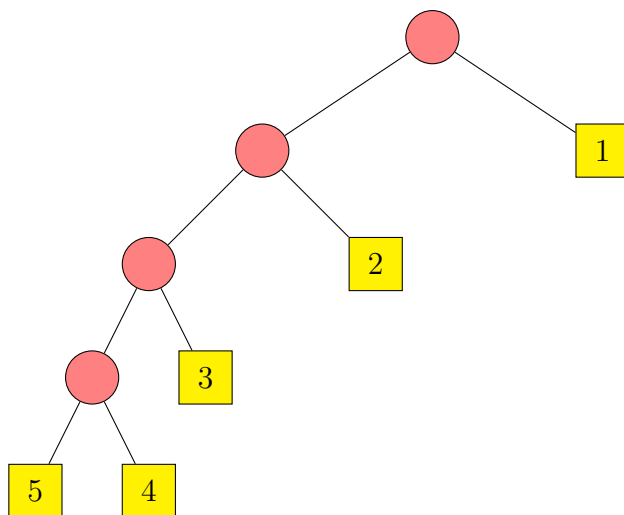


图 9.14: $WPL = 1 * 1 + 2 * 2 + 3 * 3 + 4 * 4 + 5 * 4 = 50$

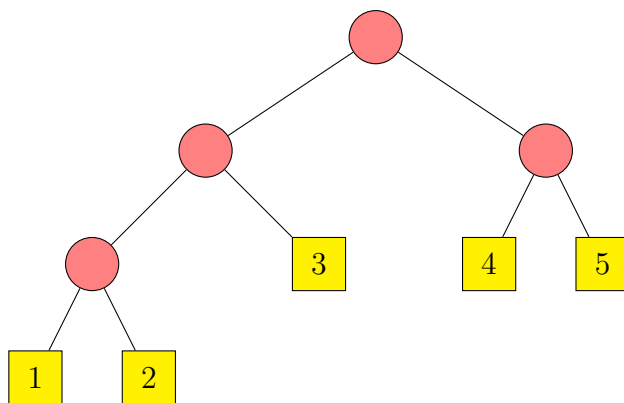


图 9.15: $WPL = 3 * 2 + 4 * 2 + 5 * 2 + 1 * 3 + 2 * 3 = 33$

怎样才能保证构建出的二叉树带权路径长度最小呢？原则上，应该让权重小的叶子结点远离树根，权重大的叶子结点靠近树根。需要注意的是，同样叶子结点所构成的哈夫曼树可能不止一棵。

9.5.2 哈夫曼树的构造

哈夫曼树的构造方法就是每次把权值最小的两棵二叉树合并。

例如有 6 个叶子结点，权重依次是 2、3、7、9、18、25。

第一步：把每一个叶子结点都当成一棵独立的树（只有根结点的树），这样就形成了一个森林。



第二步：从森林中移除权值最小的两个结点，生成父结点，父结点的权值是这两个结点权值之和，把父结点加入森林。重复该步骤，直到森林中只有一棵树为止。

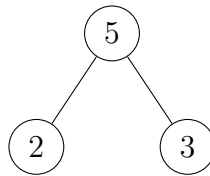


图 9.16: 合并 2 和 3

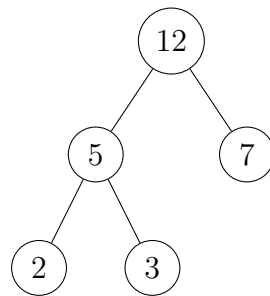


图 9.17: 合并 5 和 7

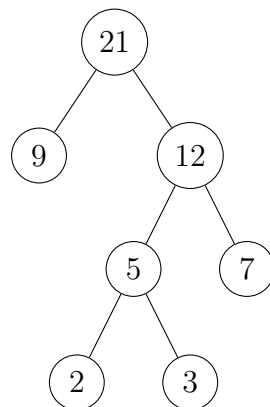


图 9.18: 合并 9 和 12

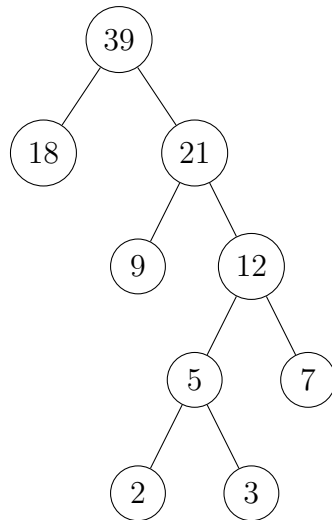


图 9.19: 合并 18 和 21

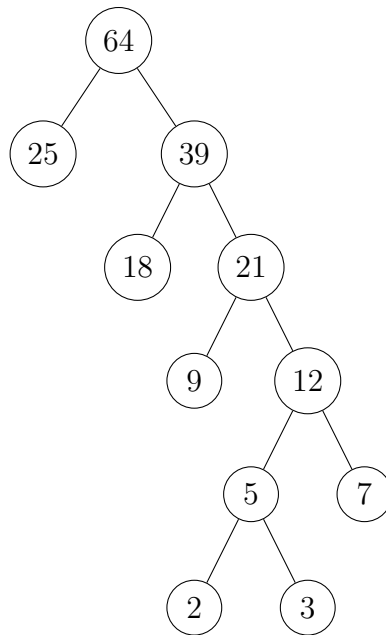


图 9.20: 合并 25 和 39

哈夫曼树有以下几个特点：

1. 没有度为 1 的结点。
2. 哈夫曼树的任意非叶结点的左右子树交换后仍是哈夫曼树。
3. 对同一组权值，可能存在不同构的两棵哈夫曼树。

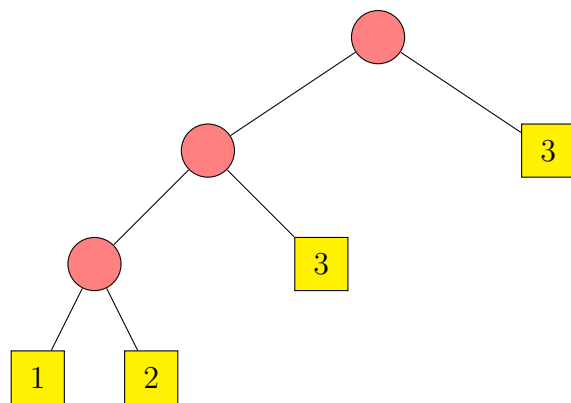


图 9.21: $WPL = 3 * 1 + 3 * 2 + 1 * 3 + 2 * 3 = 18$

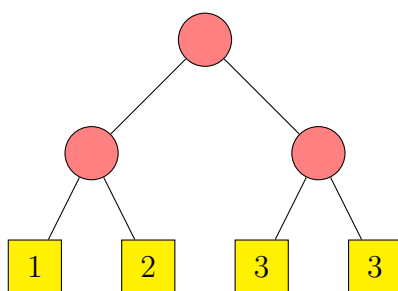


图 9.22: $WPL = 1 * 2 + 2 * 2 + 3 * 2 + 3 * 2 = 18$

9.6 哈夫曼编码

9.6.1 哈夫曼编码 (Huffman Code)

哈夫曼编码是一种高效的编码方式，在信息存储和传输过程中用于对信息进行压缩。要理解哈夫曼编码，需要从信息存储的底层逻辑讲起。

计算机不是人，它不认识中文和英文，更不认识图片和视频，它唯一认识的就是 0（低电平）和 1（高电平）。因此，计算机上一切文字、图象、音频、视频，底层都是用二进制来存储和传输的。

将信息转换成计算机能够识别的二进制形式的过程被称为编码。在 ASCII 码中，每一个字符表示成特定的 8 位二进制数。例如字符串 APPLE 表示成 8 位二进制编码为 01000001 01010000 01010000 01001100 01000101。

显然，ASCII 码是一种等长编码，也就是任何字符的编码长度都相等。等长编码的有点明显，因为每个字符对应的二进制编码长度相等，所以很容易设计，也很方便读写。但是计算机的存储空间以及网络传输的带宽是有限的，等长编码最大的缺点就是编码结果太长，会占用过多资源。

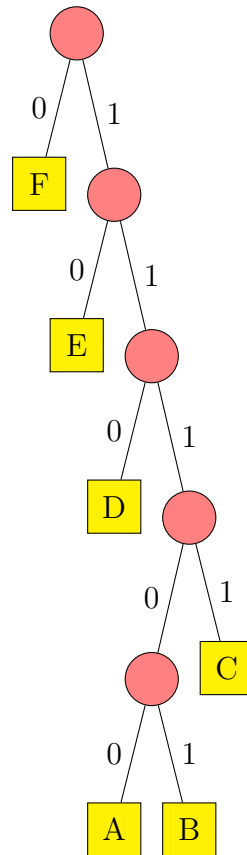
使用不等长编码，让出现频率高的字符用的编码短一些，出现频率低的字符编码长一些，可以使编码的总长度减小。但是不等长编码是不能随意设计的，如果一个字符的编码恰好是另一个字符编码的前缀，就会产生歧义的问题。

哈夫曼编码就是一种不等长的编码，并且任何一个字符的编码都不是另一个字符编码的前缀，因此可以无二义地进行解码，并且信息编码的总长度最小。

哈夫曼编码并非一套固定的编码，而是根据给定信息中各个字符出现的频次，动态生成最优的编码。哈夫曼编码的生成过程就用到了哈夫曼树。

例如一段信息里只有 A、B、C、D、E、F 这 6 个字符，出现的次数分别是 2 次、3 次、7 次、9 次、18 次、25 次。通过把这 6 个字符当成 6 个叶子结点，将出

现次数作为结点的权重，生成一颗哈夫曼树。将哈夫曼树中结点的左分支当做 0、结点的右分支当做 1，从哈夫曼树的根结点到每一个叶子结点的路径，都可以等价为一串二进制编码。



字符	编码
A	11100
B	11101
C	1111
D	110
E	10
F	0

表 9.1: 哈夫曼编码

因为每一个字符对应的都是哈夫曼树的叶子结点，从根结点到这些叶子结点的路径并没有包含关系，最终得到的二进制编码自然也不会是彼此的前缀。

Chapter 10 图

10.1 图

10.1.1 图 (Graph)

你的微信中有若干好友，而你的好友又有若干好友。许许多多的用户组成了一个多对多的关系网，这个关系网就是数据结构中的图。

再例如使用地图导航功能时，导航会根据你的出发地和目的地规划最佳的地铁换乘路线。许许多多的地铁站组成的交通网络也可以认为是图。

图是一种比树更为复杂的数据结构。树的结点之间是一对多的关系，并且存在父与子的层级划分。而图的顶点之间是多对多关系，并且所有顶点都是平等的，无所谓谁是父子。

在图中，最基本的单元是顶点 (vertex)，相当于树中的结点。顶点之间的关联关系被称为边 (edge)。图中包含一组顶点和一组边，通常用 V 表示顶点集合，用 E 表示边集合。边可以看作是顶点对，即 $(v, w) \in E, v, w \in V$ 。

在有些图中，每一条边并不是完全等同的。例如地铁线路，站与站之间的距离都有可能不同。因此图中会涉及边的权重 (weight)，涉及到权重的图被称为带权图 (weighted graph)，也称为网络。

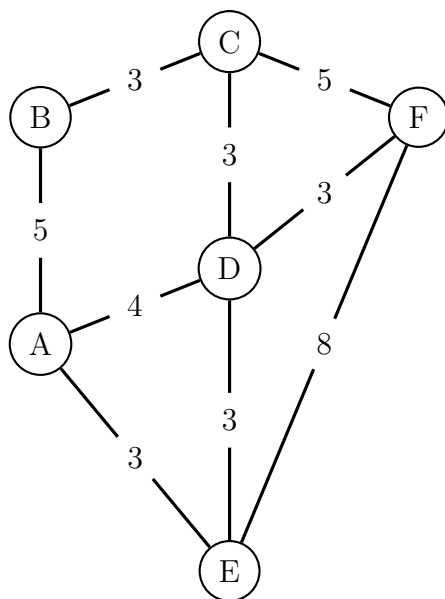


图 10.1: 带权图

还有一种图，顶点之间的关联并不是完全对称的。拿微信举例，你的好友列表里有我，但我的好友列表里未必有你。

这样一来，顶点之间的边就有了方向的区分，这种带有方向的图被称为有向图 (directed graph)。有向边可以使用 $\langle v, w \rangle$ 表示从 v 指向 w 的边。

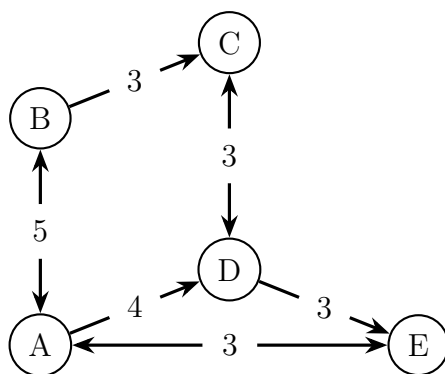


图 10.2: 有向图

相应地，在 QQ 中，只要我把你从好友里删除，你在自己的好友列表里就看不到我了。因此 QQ 的好友关系可以认为是一个没有方向区分的图，这种图被称为无向图 (undirected graph)。

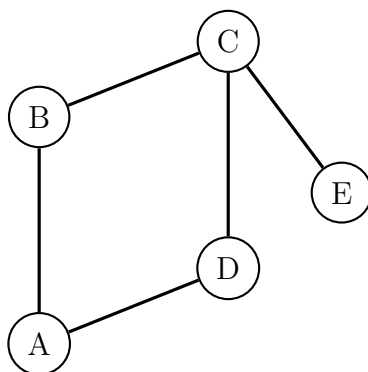
10.2 图的表示

10.2.1 邻接矩阵 (Adjacency Matrix)

拥有 n 个顶点的图，它所包含的边的数量最多是 $n(n-1)$ 条，因此，要表达各个顶点之间的关联关系，最清晰易懂的方式是使用邻接矩阵 $G[N][N]$ 。

对于无向图来说，如果顶点之间有关联，那么邻接矩阵中对应的值为 1；如果顶点之间没有关联，那么邻接矩阵中对应的值为 0。

$$G[i][j] = \begin{cases} 1 & \langle v_i, v_j \rangle \text{ 是 } G \text{ 中的边} \\ 0 & \langle v_i, v_j \rangle \text{ 不是 } G \text{ 中的边} \end{cases}$$

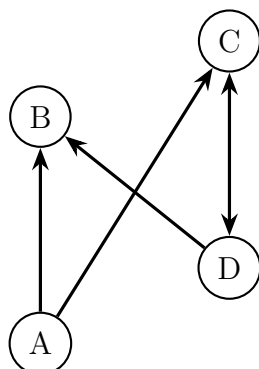


	A	B	C	D	E
A	0	1	0	1	0
B	1	0	1	0	0
C	0	1	0	1	1
D	1	0	1	0	0
E	0	0	1	0	0

表 10.1: 无向图邻接矩阵

需要注意的是，邻接矩阵从左上到右下的一条对角线上的元素值必然是 0，因为任何一个顶点与它自身是没有连接的。同时，无向图对应的邻接矩阵是一个对称矩阵，假如 A 和 B 有关联，那么 B 和 A 也必定有关联。

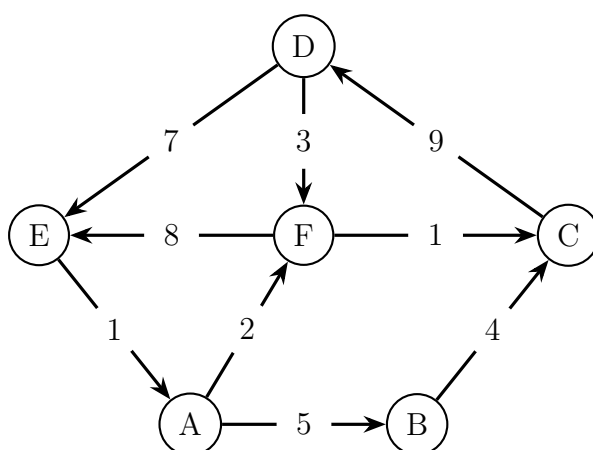
但是对于有向图的邻接矩阵，不一定是一个对称矩阵，假如 A 可以达到 B，从 B 未必能达到 A。



	A	B	C	D
A	0	1	1	0
B	0	0	0	0
C	0	0	0	1
D	0	1	1	0

表 10.2: 有向图邻接矩阵

对于网络，只要把邻接矩阵对应位置的值定义为边 $\langle v_i, v_j \rangle$ 的权重即可。



	A	B	C	D	E	F
A	∞	5	∞	∞	∞	2
B	∞	∞	4	∞	∞	∞
C	∞	∞	∞	9	∞	∞
D	∞	∞	∞	∞	7	3
E	1	∞	∞	∞	∞	∞
F	∞	∞	1	∞	8	∞

表 10.3: 带权图邻接矩阵

对于带权图，如果 v_i 和 v_j 之前没有边应该将权值设为 ∞ 。

邻接矩阵的优点：

1. 简单、直观。
2. 可以快速查到一个顶点和另一顶点之间的关联关系。
3. 方便计算任一顶点的度，对于有向图，从顶点发出的边数为出度，指向顶点的边数为入度。

邻接矩阵的缺点：

1. 浪费空间，对于稀疏图（点很多而边很少）有大量无效元素。但对于稠密图（特别是完全图）还是很合算的。
2. 浪费时间，统计稀疏图中边的个数，也就是计算邻接矩阵中元素 1 的个数。

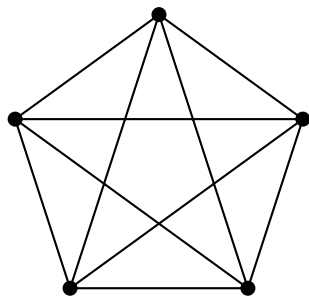


图 10.3: 完全图

10.2.2 邻接表 (Adjacency List)

为了解决邻接矩阵占用空间的问题，人们想到了另一种图的表示方法——邻接表。在邻接表中，图的每一个顶点都是一个链表的头结点，其后连接着该顶点能够直接到达的相邻顶点。对于稀疏图而言，邻接表存储方式占用的空间比邻接矩阵要小得多。

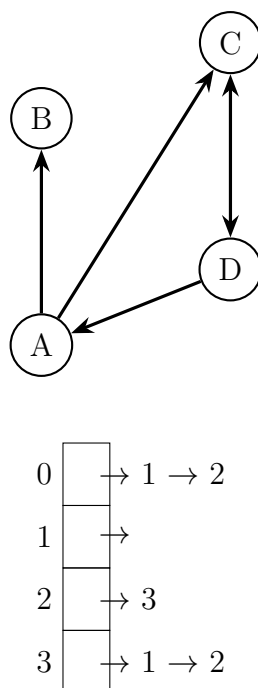


图 10.4: 邻接表

通过遍历邻接表可以查找到所有能够到达的相邻顶点，但是对于逆向查找，即哪些顶点可以达到一个顶点就会很麻烦。

逆邻接表和邻接表是正好相反的，逆邻接表每一个顶点作为链表的头结点，后继结点所存储的是能够直接到达该顶点的相邻顶点。

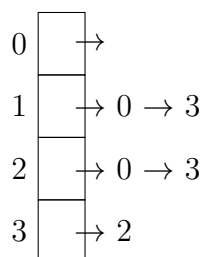


图 10.5: 逆邻接表

可是，一个图要维护正反两个邻接表，也太麻烦了吧？

通过十字链表可以把邻接表和逆邻接表结合在一起。

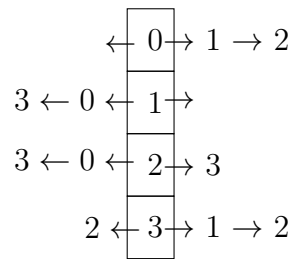


图 10.6: 十字链表

10.3 图的遍历

10.3.1 深度优先搜索 (DFS, Depth First Search)

深度优先搜索是一种一头扎到底的遍历方法，选择一条路，尽可能不断地深入，遇到死路就回退，回退过程中如果遇到没探索的支路，就进入该支路继续深入。

例如一个小镇的每个地点都藏有可以实现愿望的光玉，现在要从 0 号出生点出发去收集所有的光玉。

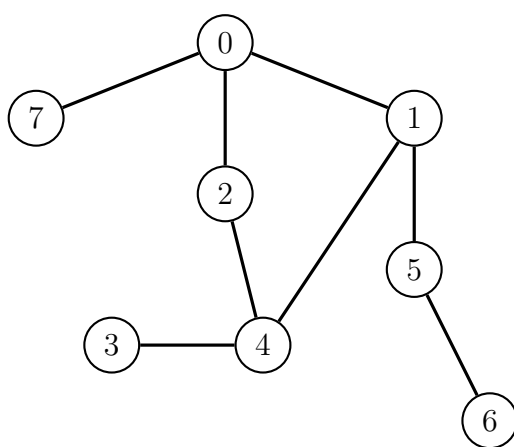


图 10.7: 深度优先搜索

二叉树的先序遍历本质上也可以认为是图的深度优先遍历。要想实现回溯，可以利用栈的先进后出的特性，也可以采用递归的方式，因为递归本身就是基于方法调用栈来实现的。

Algorithm 1 深度优先搜索

```
1: procedure DFS(Vertex V)
2:   isVisited[V] = true
3:   for v in V do
4:     if !isVisited[v] then
5:       DFS(v)
6:     end if
7:   end for
8: end procedure
```

10.3.2 广度优先搜索 (BFS, Breath First Search)

除了深度优先搜索一头扎到底的方法以外，还有一种方法就是首先把从源点相邻的顶点遍历，然后再遍历稍微远一点的顶点，再去遍历更远一点的顶点。

二叉树的层次遍历本质上也可以认为是图的广度优先遍历，需要借助队列来实现重放。

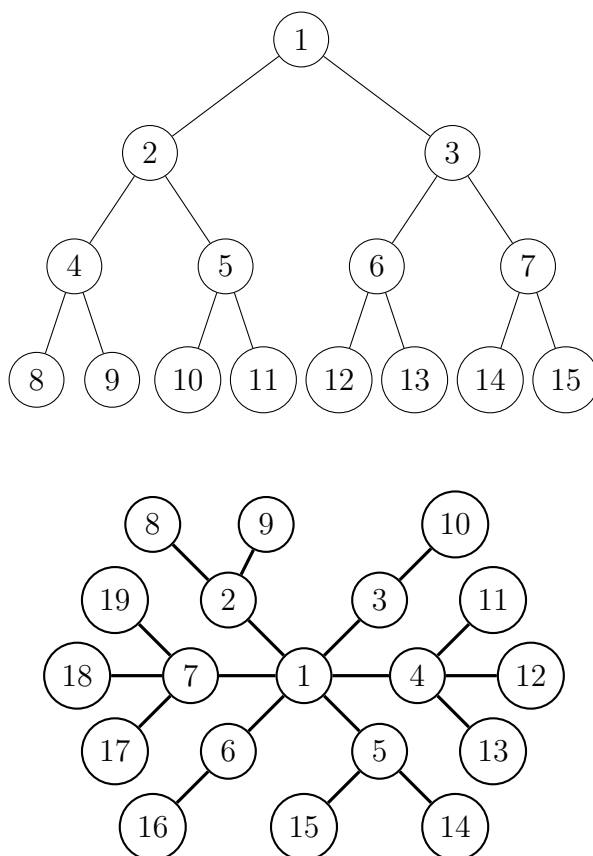


图 10.8: 广度优先搜索

Algorithm 2 广度优先搜索

```
1: procedure BFS(Vertex V)
2:   isVisited[V] = true
3:   enqueue(Q, V)
4:   while !isEmpty(Q) do
5:     V = dequeue(Q)
6:     for v in V do
7:       if !isVisited[v] then
8:         isVisited[v] = true
9:         enqueue(Q, v)
10:      end if
11:    end for
12:  end while
13: end procedure
```

10.4 连通图

10.4.1 连通图

图还有一些有关路径的术语：

- 连通：如果从顶点 V 到 W 存在一条路径，则称 V 和 W 是连通的。
- 路径：顶点 V 到 W 的路径是一系列顶点 $\{V, v_1, v_2, \dots, v_n, W\}$ 的集合，其中任意一对相邻的顶点间都有图中的边。
- 路径长度：路径中边的个数，如果是带权图（网络），则是所有边的权重和。
- 简单路径：顶点 V 到 W 之间的路径中所有顶点都不同。
- 回路：起点等于终点的路径。

如果图中任意两顶点均连通，那么称这个图是一个连通图。

一个图的连通分量指的是图的极大连通子图，极大连通子图需要满足两点要求：

1. 顶点数到达极大，即再加一个顶点就不连通了。
2. 边数达到极大，即包含子图中所有顶点相连的所有边。

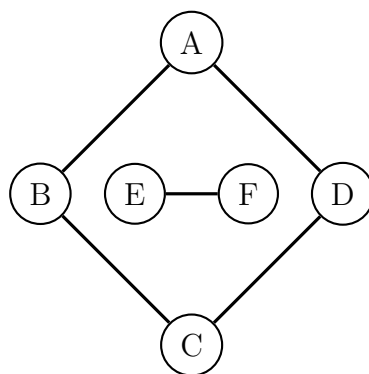


图 10.9: 图 G

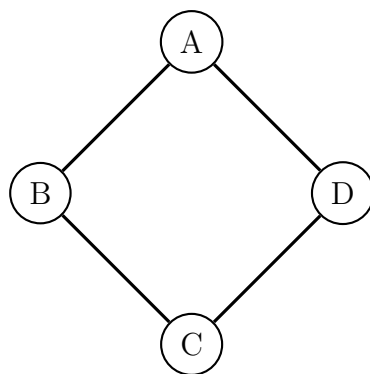


图 10.10: 是图 G 的极大连通子图

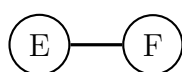


图 10.11: 是图 G 的极大连通子图

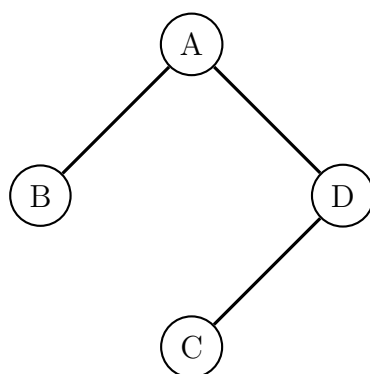


图 10.12: 不是图 G 的极大连通子图

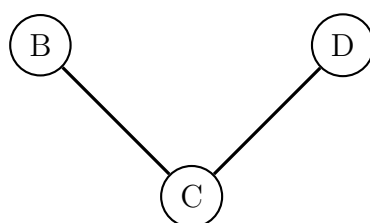


图 10.13: 不是图 G 的极大连通子图

对于有向图而言，如果有向图中任意一对顶点 V 和 W 之间存在双向路径，既可以从 V 走到 W ，也可以从 W 走到 V ，但这两条路径不一定是同一条，则称该图为强连通图。

如果有向图不是强连通图，但将所有的有向边替换为无向边之后可以变为连通图，则称该图为弱连通图。

有向图的极大强连通子图称为强连通分量。

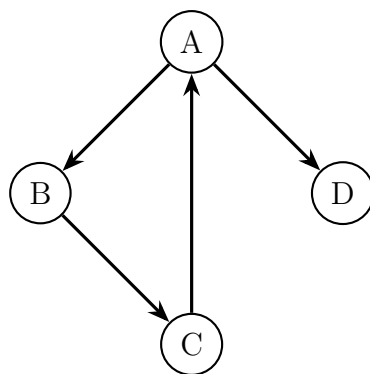


图 10.14: 有向图 G

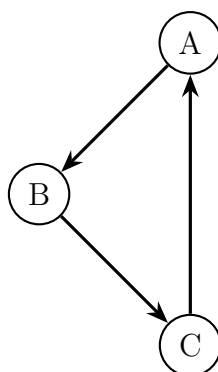


图 10.15: 是有向图 G 的强连通分量



图 10.16: 是有向图 G 的强连通分量

10.4.2 非连通图的遍历

如果一个图不是连通图，那么无论使用深度优先遍历还是广度优先遍历，都会有顶点无法被访问到。解决这个问题的方式是每调用一次 $\text{dfs}(V)$ 或 $\text{bfs}(V)$ ，就把

顶点 V 所在的连通分量遍历一遍。

Algorithm 3 非连通图的深度优先遍历

```
1: procedure DFS(Vertex  $V$ )
2:    $\text{isVisited}[V] = \text{true}$ 
3:   for  $v$  in  $V$  do
4:     if  $\text{!isVisited}[v]$  then
5:       DFS( $v$ )
6:     end if
7:   end for
8: end procedure
9:
10: procedure LISTCOMPONENTS(Graph  $G$ )
11:   for  $V$  in  $G$  do
12:     if  $\text{!isVisited}[V]$  then
13:       DFS( $V$ )
14:     end if
15:   end for
16: end procedure
```

10.5 最短路径

10.5.1 最短路径 (Shortest Path)

在现实中很多需要都运用到了最短路径的算法，例如从一个地铁站到另一个地铁站的最快换乘路线等。地铁线路图中，地铁站可以看作是图的顶点，站与站之间的线路可以看作是边，权重可以是距离、时间、费用等。



图 10.17: 上海地铁线路图

在网络中，求两个不同顶点之间的所有路径中，边的权值之和最小的那一条路径，这条路径就是两点之间的最短路径。其中最短路径的第一个顶点称为源点 (source)，最后一个顶点为终点 (destination)。

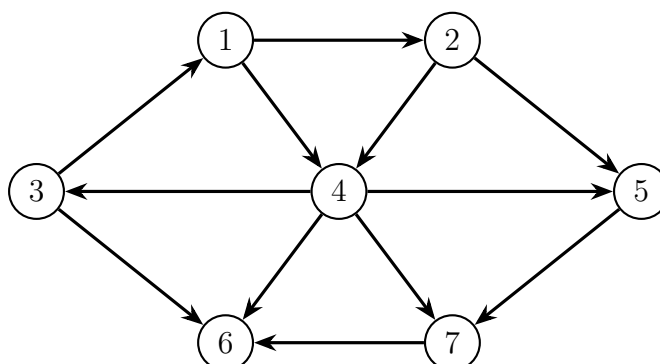
图的最短路径问题分为 2 种类型：

1. 单源最短路径：从某固定源点出发，求到所有其它顶点的最短路径。
2. 多源最短路径：求任意两顶点间的最短路径。

10.5.2 无权图的单源最短路径算法(SSSP, Single-Source Shortest Path)

无权图的单源最短路径算法可以按照递增（非递减）的顺序找出到各个顶点的最短路，算法类似广度优先遍历。

例如在一个无权图中，以顶点 3 作为源点，离源点距离为 1 的顶点有 1 和 6，距离为 2 的顶点有 2 和 4，距离为 3 的顶点有 5 和 7。



无权图的单元最短路径算法中， $\text{dist}[v]$ 存储从源点 S 到 v 的最短路径，初始化源点 $\text{dist}[S]$ 的距离为 0， $\text{path}[v]$ 表示达到顶点路径 v 上一个经过的顶点。

顶点	1	2	3	4	5	6	7
dist	1	2	0	2	3	1	3
path	3	1	-1	1	2	3	4

表 10.4: 最短路径表

Algorithm 4 无权图的单源最短路径

```
1: procedure UNWEIGHTEDSSSP(Vertex S)
2:   enqueue(Q, S)
3:   while !isEmpty(Q) do
4:     V = dequeue(Q)
5:     for v in V do
6:       if dist[v] == -1 then
7:         dist[v] = dist[V] + 1
8:         path[v] = V
9:         enqueue(Q, v)
10:      end if
11:    end for
12:  end while
13: end procedure
```

10.5.3 有权图的单源最短路径算法

有权图的最短路径不一定是经过顶点数最少的路。如果图中存在负值圈(negative-cost cycle)的话会导致算法失效，因为沿着回路走无穷多次，花销是负无穷。

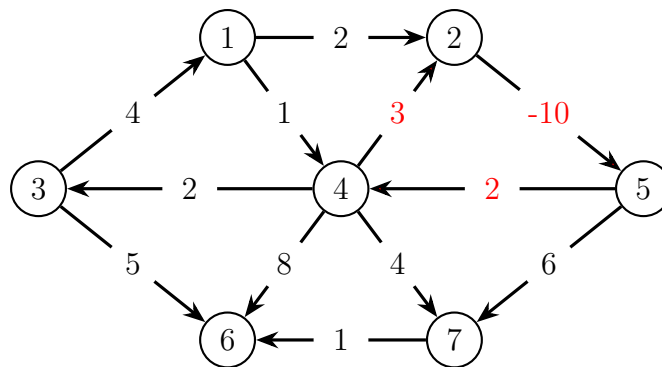


图 10.18: 负值圈

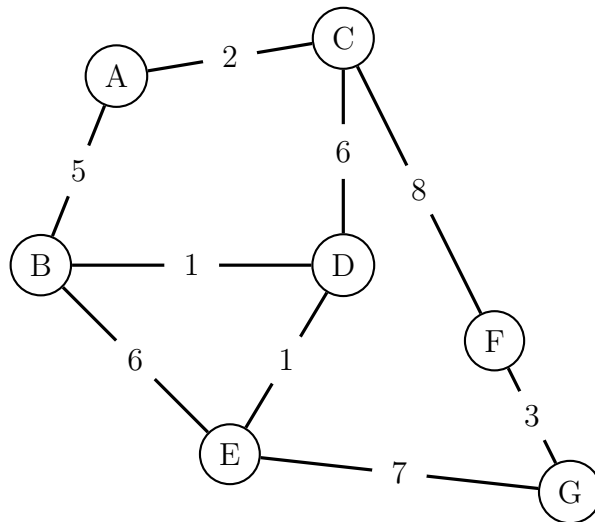
带权图的单源最短路径可以通过迪杰斯特拉 (Dijkstra) 算法解决。

特斯拉？什么鬼？

迪杰斯特拉算法的本质是不断刷新起点与其他各个顶点之间的距离表。Dijkstra 算法采用了贪心的思想，每次都未收录的顶点中选取 dist 值最小的收录。每当收录一个顶点时，可能会影响另外一个顶点的 dist 值。

$$dist[w] = \min\{dist[w], dist[v] + weight_{<v,w>}\}$$

例如计算从源点 A 到其它各顶点的最短路径。



第 1 步：创建距离表。其中表中 key 是顶点名称，value 是源点 A 到对应顶点的已知最短距离。一开始并不知道最短路径是多少，因此 value 都为 ∞ 。

B	C	D	E	F	G
∞	∞	∞	∞	∞	∞

第 2 步：找到源点 A 的邻接点 B 和 C，从 A 到 B 的距离是 5，从 A 到 C 的距离是 2。

B	C	D	E	F	G
5	2	∞	∞	∞	∞

第 3 步：从距离表中找到从 A 出发距离最短的顶点，也就是顶点 C。找到顶点 C 的邻接点 D 和 F（A 已经遍历过不需要考虑）。从 C 到 D 的距离是 6，所以从 A 到 D 的距离是 $2 + 6 = 8$ ；从 C 到 F 的距离是 8，所以从 A 到 F 的距离是 $2 + 8 = 10$ 。

B	C	D	E	F	G
5	2	8	∞	10	∞

第 4 步：从距离表中找到从 A 出发距离最短的顶点（C 已经遍历过不需要考虑），也就是顶点 B。找到顶点 B 的邻接点 D 和 E（A 已经遍历过不需要考虑）。从 B 到 D 的距离是 1，所以从 A 到 D 的距离是 $5 + 1 = 6$ ，小于距离表中的 8；从 B 到 E 的距离是 6，所以从 A 到 E 的距离是 $5 + 6 = 11$ 。

B	C	D	E	F	G
5	2	6	11	10	∞

第 5 步：从距离表中找到从 A 出发距离最短的顶点（B 和 C 不用考虑），也就是顶点 D。找到顶点 D 的邻接点 E 和 F。从 D 到 E 的距离是 1，所以从 A 到 E 的距离是 $6 + 1 = 7$ ，小于距离表中的 11；从 D 到 F 的距离是 2，所以从 A 到 F 的距离是 $6 + 2 = 8$ ，小于距离表中的 10。

B	C	D	E	F	G
5	2	6	7	8	∞

第 6 步：从距离表中找到从 A 出发距离最短的顶点，也就是顶点 E。找到顶点 E 的邻接点 G。从 E 到 G 的距离是 7，所以从 A 到 G 的距离是 $7 + 7 = 14$ 。

B	C	D	E	F	G
5	2	6	7	8	14

第 7 步：从距离表中找到从 A 出发距离最短的顶点，也就是顶点 F。找到顶点 F 的邻接点 G。从 F 到 G 的距离是 3，所以从 A 到 G 的距离是 $8 + 3 = 11$ ，小于距离表中的 14。

B	C	D	E	F	G
5	2	6	7	8	11

最终，距离表中存储的是从源点 A 到所有顶点的最短距离。

10.5.4 多源最短路径算法

如何能够求出一个带权图中所有顶点两两之间的最短距离呢？

对了！刚刚学过了 Dijkstra 算法，可以对每个顶点都使用一次 Dijkstra 算法，这样就求出了所有顶点之间的最短距离。

这个思路确实可以实现，但是 Dijkstra 算法的代码逻辑比较复杂，有没有更简单的方法呢？

弗洛伊德（Floyd-Warshall）算法是专门用于寻找带权图中多源点之间的最短路径算法。Floyd 算法的思想是，若想缩短两点间的距离，仅有一种方式，那就是通过第三顶点绕行。

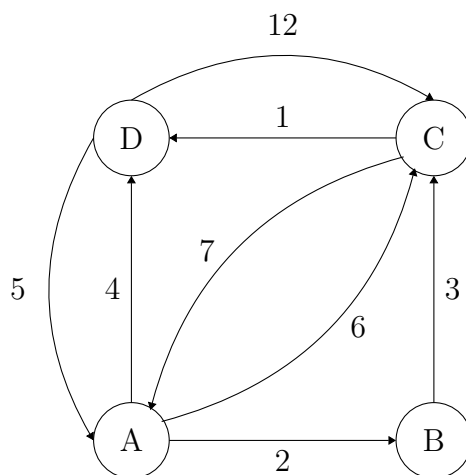
假设 $D^k[i][j]$ 为路径 $\{i \rightarrow \{l \leq k\} \rightarrow j\}$ 的最小长度。当 D^{k-1} 已经完成，递推到 D^k 时：

1. 如果 $k \notin$ 最短路径 $\{i \rightarrow \{l \leq k\} \rightarrow j\}$ ，则 $D^k = D^{k-1}$ 。
2. 如果 $k \in$ 最短路径 $\{i \rightarrow \{l \leq k\} \rightarrow j\}$ ，该路径必定由两段最短路径组成，则 $D^k[i][j] = D^{k-1}[i][k] + D^{k-1}[k][j]$ 。

例如，小哼准备去一些城市旅游，有些城市之间有公路，有些城市之间则没有。为了节省经费以及方便计划旅程，小哼希望在出发之前直到任意两个城市之间的最短路程。

如果要想任意两点之间的路程变短，只能引入第三个点，并通过这个顶点中转才有可能缩短原来的路程。这个中转的顶点甚至有时候不只通过一个顶点，而是经过两个或更多点中转会更短。

当任意两点之间不允许经过第三个点中转时，这些城市之间的最短路径就是邻接矩阵的初始路径。



	A	B	C	D
A	0	2	6	4
B	∞	0	3	∞
C	7	∞	0	1
D	5	∞	12	0

在只允许经过 1 号顶点中转的情况下，任意两点之间的最短路程更新为：

	A	B	C	D
A	0	2	6	4
B	∞	0	3	∞
C	7	9	0	1
D	5	7	11	0

在只允许经过 1 号和 2 号顶点中转的情况下，任意两点之间的最短路程更新为：

	A	B	C	D
A	0	2	5	4
B	∞	0	3	∞
C	7	9	0	1
D	5	7	10	0

在只允许经过 1 号、2 号和 3 号顶点中转的情况下，任意两点之间的最短路程更新为：

	A	B	C	D
A	0	2	5	4
B	10	0	3	4
C	7	9	0	1
D	5	7	10	0

最后允许通过所有顶点作为中转，任意两点之间的最短路程更新为：

	A	B	C	D
A	0	2	5	4
B	9	0	3	4
C	6	8	0	1
D	5	7	10	0

Floyd 最短路径

```

1 void floyd(Graph *g, int dist[MAX][MAX]) {
2     // 最短路径矩阵初始化为图的邻接矩阵
3     for(int i = 0; i < g->vertexNum; i++) {
4         for(int j = 0; j < g->vertexNum; j++) {
5             dist[i][j] = g->weight[i][j];
6         }
7     }
8
9     // Floyd算法
10    for(int k = 0; k < g->vertexNum; k++) {
11        for(int i = 0; i < g->vertexNum; i++) {
12            for(int j = 0; j < g->vertexNum; j++) {
13                if(dist[i][k] + dist[k][j] < dist[i][j]) {
14                    dist[i][j] = dist[i][k] + dist[k][j];
15                }
16            }
17        }
18    }
19 }

```

10.6 最小生成树

10.6.1 最小生成树 (MST, Minimum Spanning Tree)

所谓最小生成树，就是一个图的极小连通子图，它包含原图的所有顶点，并且所有边的权值之和尽可能小。

最小生成树需要满足 3 个条件：

1. 是一棵树：树不能有回路，并且 $|V|$ 个顶点一定有 $|V| - 1$ 条边。
2. 是生成树：包含原图的全部顶点，树的 $|V| - 1$ 条边都必须在图里，并且如果向生成树中任意加一条边都一定构成回路。
3. 边的权重和最小。

如果最小生成树存在，那么图一定连通，反之亦然。

例如一个带权图，蓝色边可以把所有顶点连接起来，又保证边的权值和最小。

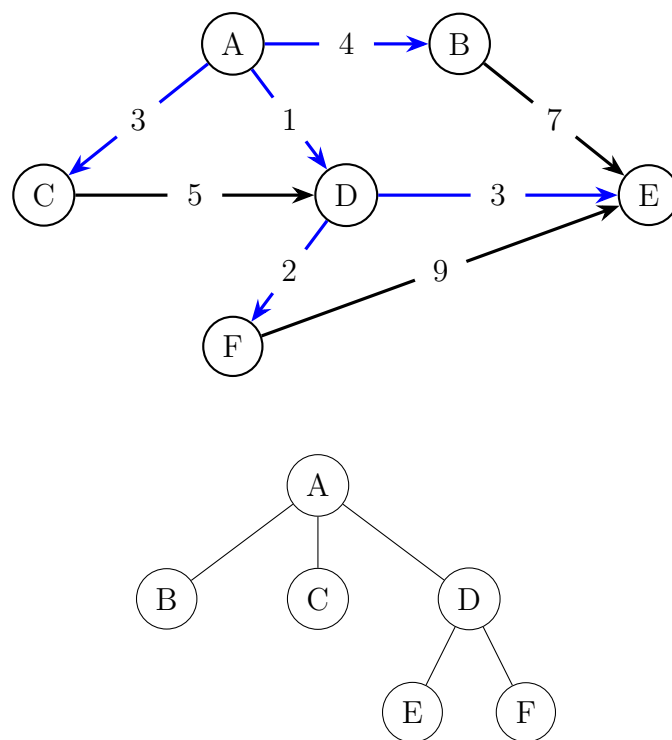


图 10.19: 最小生成树

图的最小生成树不是唯一的，同一个图有可能对应多个最小生成树。

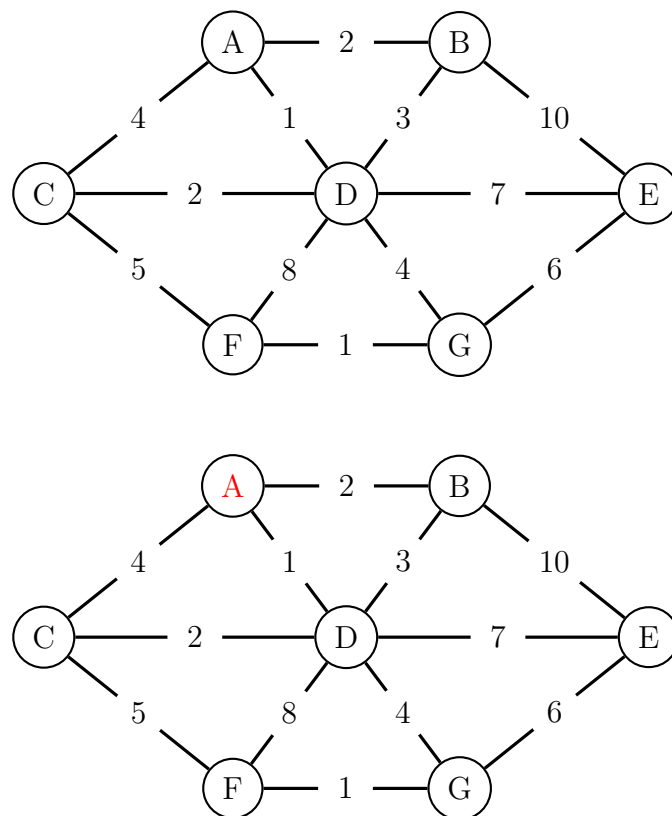
最小生成树在现实中很很多用处。假如要在若干个城市之间铺设铁路，而预算又是有限的，那么就需要寻找成本最低的铺设方式。城市之间的交通网就像一个连通图，其实并不需要在每两个城市之间都直接进行连接，只需要一个最小生成树，保证所有的城市都有铁路可以达到即可。

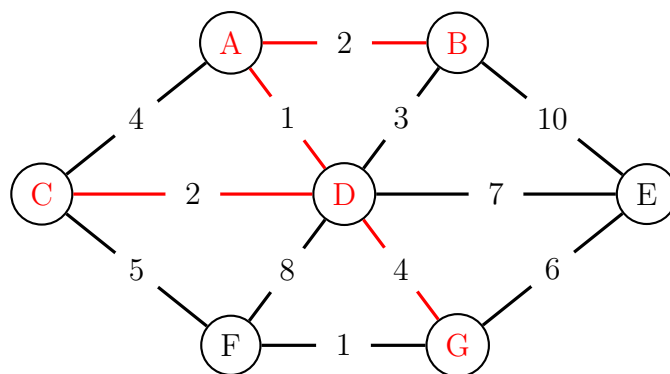
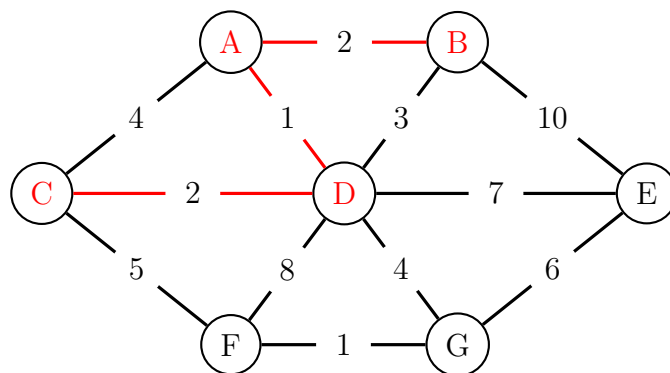
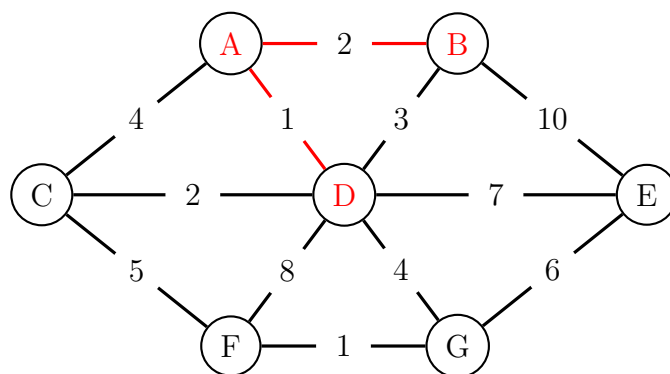
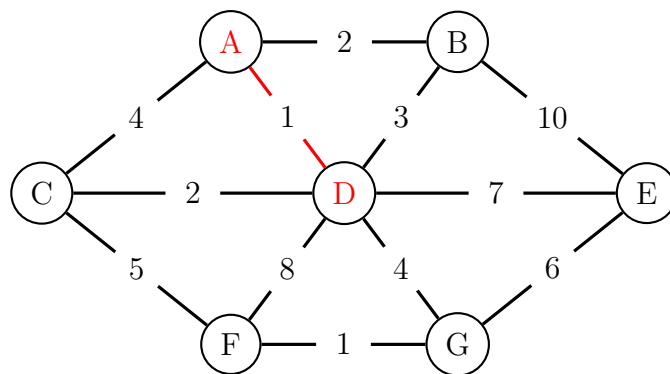
10.6.2 Prim

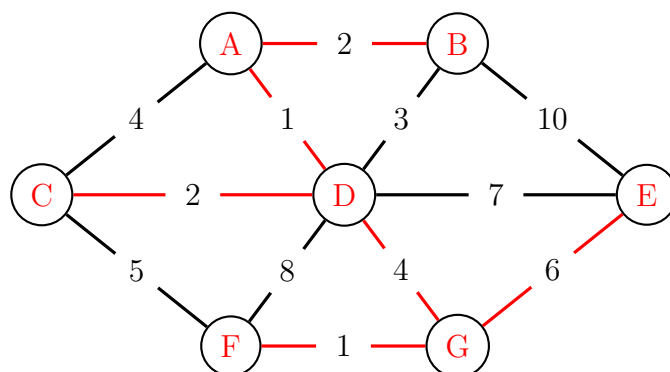
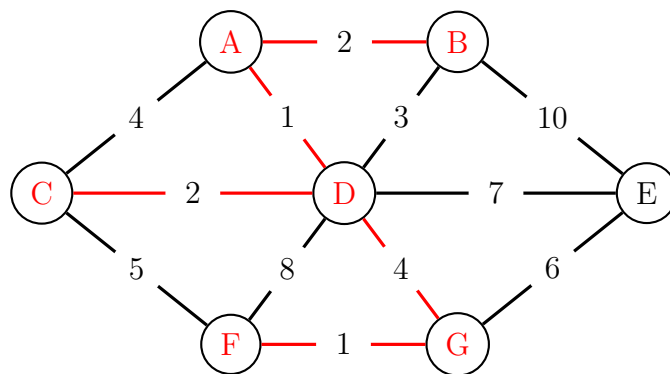
Prim 算法是以图的顶点为基础，从一个初始顶点开始，寻找达到其它顶点权值最小的边，并把该顶点加入到已触达顶点的集合中。当全部顶点都加入到集合时，算法的工作就完成了。Prim 算法的本质是基于贪心算法（greedy algorithm）。

Prim 算法可以理解为让一棵小树长大，每次找能够向外生长的最小边。

例如使用 Prim 算法获得一个带权图的最小生成树：





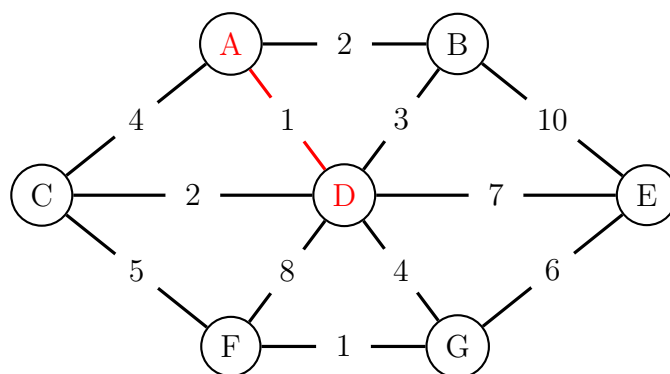


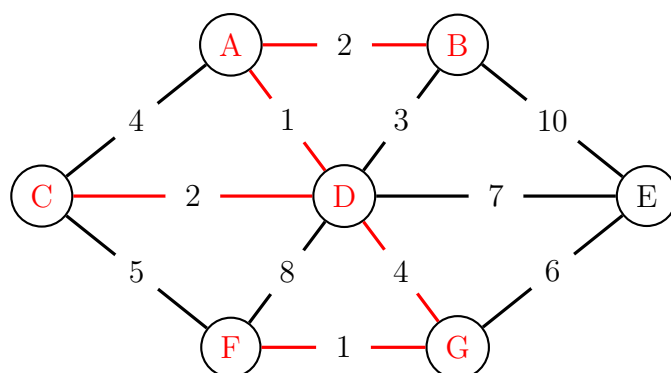
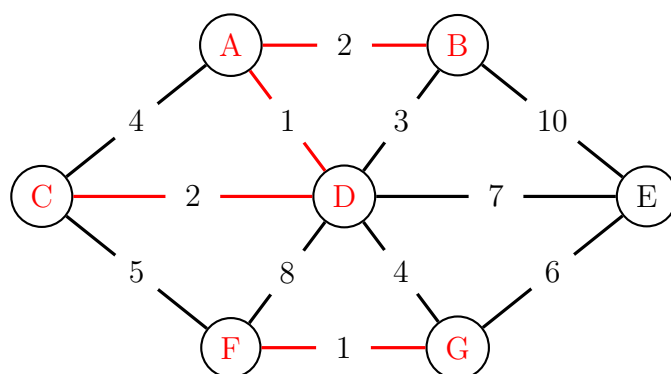
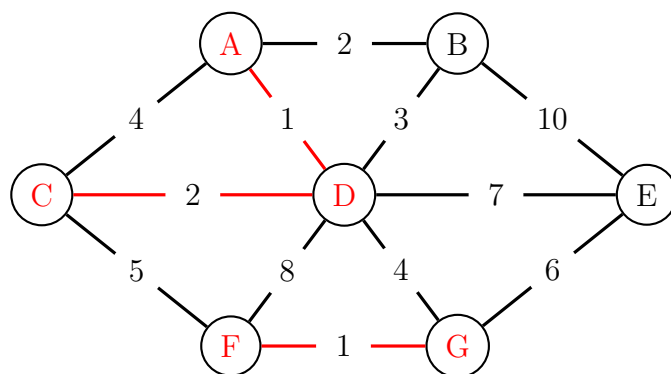
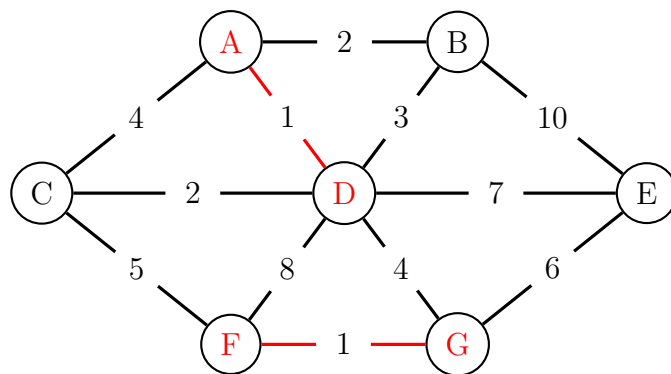
10.6.3 Kruskal

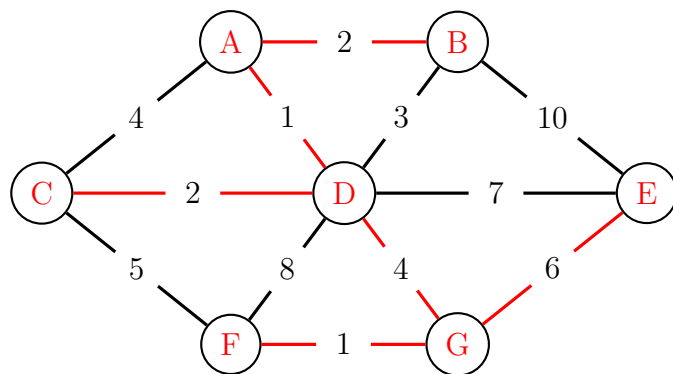
与 Prim 算法不同，Prim 算法是以顶点为关键来获得最小生成树的，而 Kruskal 算法是以边为关键获得最小生成树的。

Kruskal 算法可以理解为将森林合并成树，每次在图中找权值最小的边收录。

例如使用 Kruskal 算法获得一个带权图的最小生成树：







10.7 拓扑排序

10.7.1 拓扑排序 (Topological Sort)

一项大的工程常被分为多个小的子工程，子工程之间可能存在一定的先后顺序，即某些子工程必须在其它的一些子工程完成后才能开始。在现代化管理中，有向图可以用来描述和分析一项工程的计划和实施过程，其中图的顶点表示活动，有向边表示活动之间的先后关系，这样的图称为 AOV (Activity on Vertex) 网。

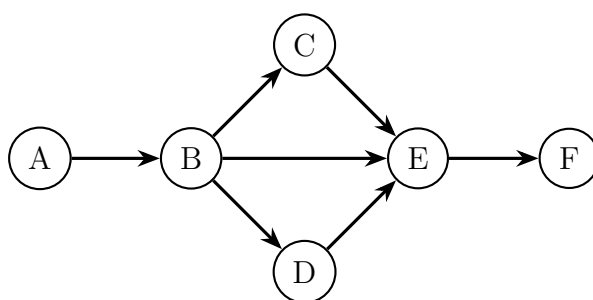


图 10.20: AOV 网

如果图中从顶点 V 到 W 有一条有向路径，则 V 一定排在 W 之前，满足此条件的顶点序列称为一个拓扑序。AOV 网如果有合理的拓扑序，则必定是有向无环图 (DAG, Directed Acyclic Graph)。

可以使用卡恩 Kahn 算法完成拓扑排序。假设列表 L 用于存放拓扑排序的结果，把所有入度为 0 的顶点放入 L 中，然后把这此顶点从图中去掉。重复该操作，直到找不到入度为 0 的顶点。如果此时 L 中的元素个数和图的顶点总数相同，说明拓扑排序完成；如果此时 L 中的元素个数少于图的顶点总数，说明原图中存在环，无法进行拓扑排序。

例如对计算机专业课安排学习顺序：

课程代码	课程名称	预修课程
C1	程序设计基础	无
C2	离散数学	无
C3	数据结构	C1、C2
C4	微积分（一）	无
C5	微积分（二）	C4
C6	线性代数	C5
C7	算法分析与设计	C3
C8	逻辑与计算机设计基础	无
C9	计算机组成	C8
Chapter10	操作系统	C7、C9
C11	编译原理	C7、C9
C12	数据库	C7
C13	计算理论	C2
C14	计算机网络	Chapter10
C15	数值分析	C6

表 10.5: 计算机专业课

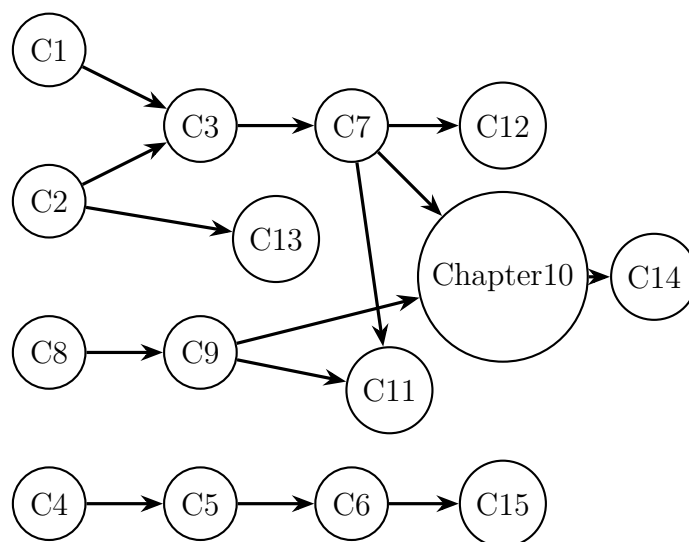
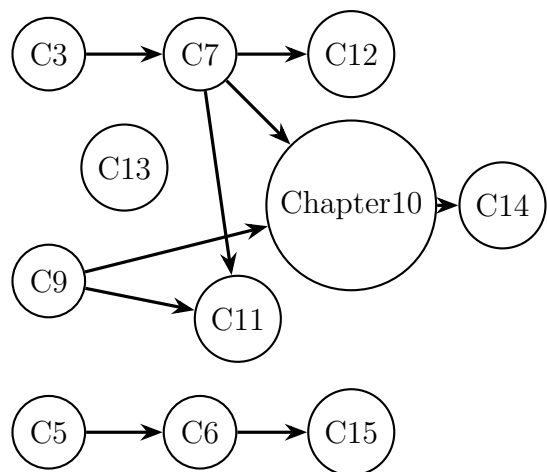
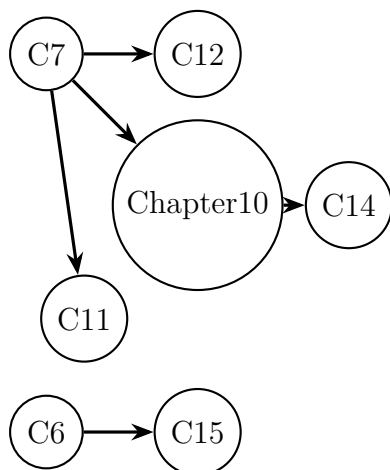


图 10.21: 计算机专业课 AOV 网



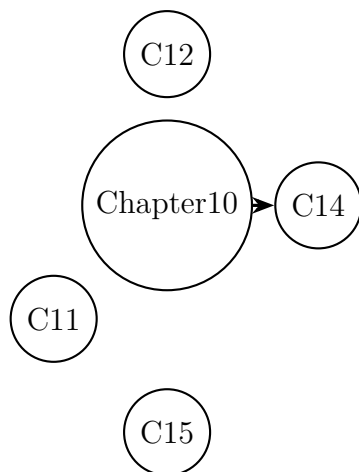
学期	课程
第一学期	C1、C2、C4、C8

表 10.6: 第一学期课程



学期	课程
第一学期	C1、C2、C4、C8
第二学期	C3、C5、C9、C13

表 10.7: 第二学期课程



学期	课程
第一学期	C1、C2、C4、C8
第二学期	C3、C5、C9、C13
第三学期	C6、C7

表 10.8: 第三学期课程



学期	课程
第一学期	C1、C2、C4、C8
第二学期	C3、C5、C9、C13
第三学期	C6、C7
第四学期	Chapter10、C11、C12、C15

表 10.9: 第四学期课程

学期	课程
第一学期	C1、C2、C4、C8
第二学期	C3、C5、C9、C13
第三学期	C6、C7
第四学期	Chapter10、C11、C12、C15
第五学期	C14

表 10.10: 第五学期课程

Chapter 11 分治法

11.1 最近点对

11.1.1 最近点对

在一个平面上有 n 个点，找到所有点对中距离最短的点对。

暴力解法就是计算任意两点之间的距离，找到其中的最小值，因此时间复杂度为 $O(n^2)$ 。

利用分治法，可以根据排序后的横坐标将点集分为左右两个部分，然后递归地对两个子问题进行求解。先求出左半部分的最短距离，再求出右半部分的最短距离。但是最短距离的点对也有可能会跨越边界，因此还需要计算一个点在左半部分、另一个点在右半部分的最短距离。三个距离中最短的就是原问题的最终解。

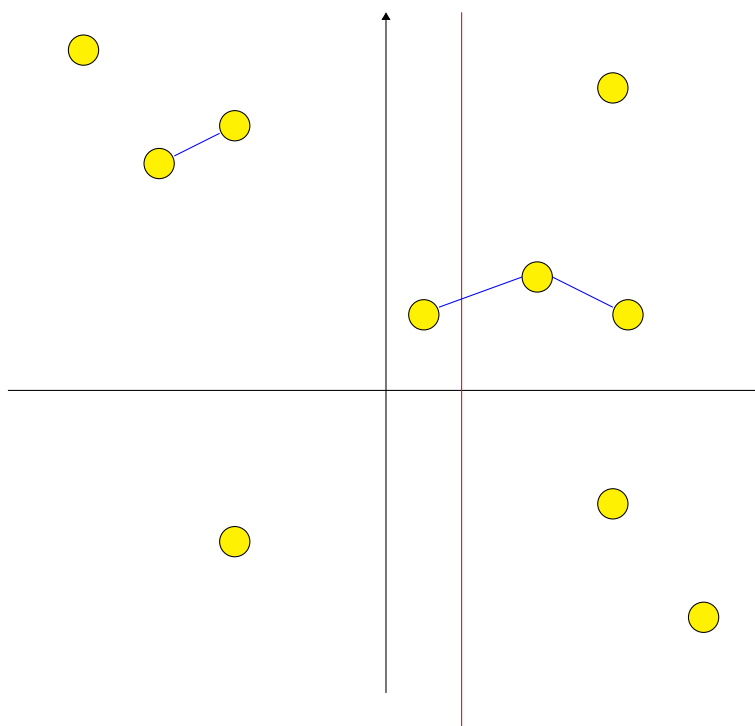


图 11.1: 最近点对

在计算出左右两边的最短距离后，两者较小的值 d 即为跨越边界的范围。

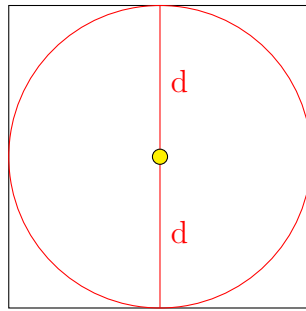


图 11.2: 跨越边界的范围

对于处于左边的每个点而言，在右边的每个小长方形中至多存在 1 个点，即最多只需要比较右边的 6 个点。如果在右半边存在超过 6 个点的话，那么就存在一个比 d 更短的距离，就与之前的最短距离矛盾了。

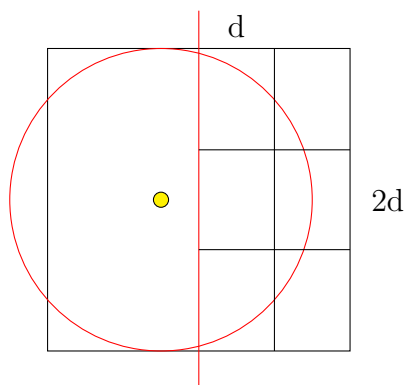


图 11.3: 检查跨越边界的点

检查 1 个点是常数时间，检查 n 个点需要 $O(n)$ 的时间。分治算法中排序需要 $O(n \log n)$ ，递归处理子问题需要 $T(n/2) + O(n)$ 。总体时间复杂度为 $O(n \log n)$ 。

11.2 凸包

11.2.1 凸包 (Convex Hull)

凸包是计算几何中的概念。在大量离散点的集合中，求一个最小的凸多边形，使得所有点都在该多边形的内部或边上。凸包在形状识别、字形识别、碰撞检测中都有所应用。

利用分治算法，连接最大纵坐标和最小纵坐标的两点，将点集划分成左右两部分，并递归对两个子问题进行求解。

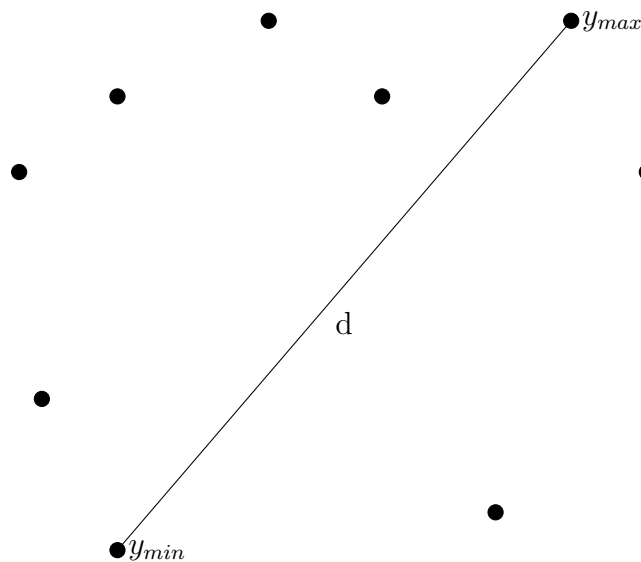


图 11.4: 划分点集

例如在左半边中，先找到距离边界 d 最远的点 P 。落在形成的三角形内部的点全部排除。接着将边 a 外的点与 a 构成作左半边点集的子问题，将边 b 外的点与 b 构成另一个左半边点集的子问题。每次将距离边界最远的点加入凸包即可求解原问题。

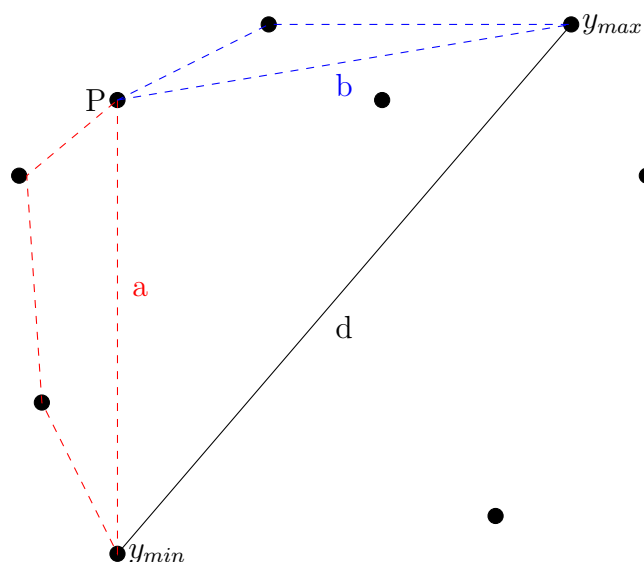


图 11.5: 子问题

每个子问题可以找到一个凸包上的点，使问题规模缩小 1，寻找凸包顶点和划分子问题需要的时间 $O(n)$ 。当子问题的规模小于 3 时，可直接进行求解。因此分治算法的时间复杂度为：

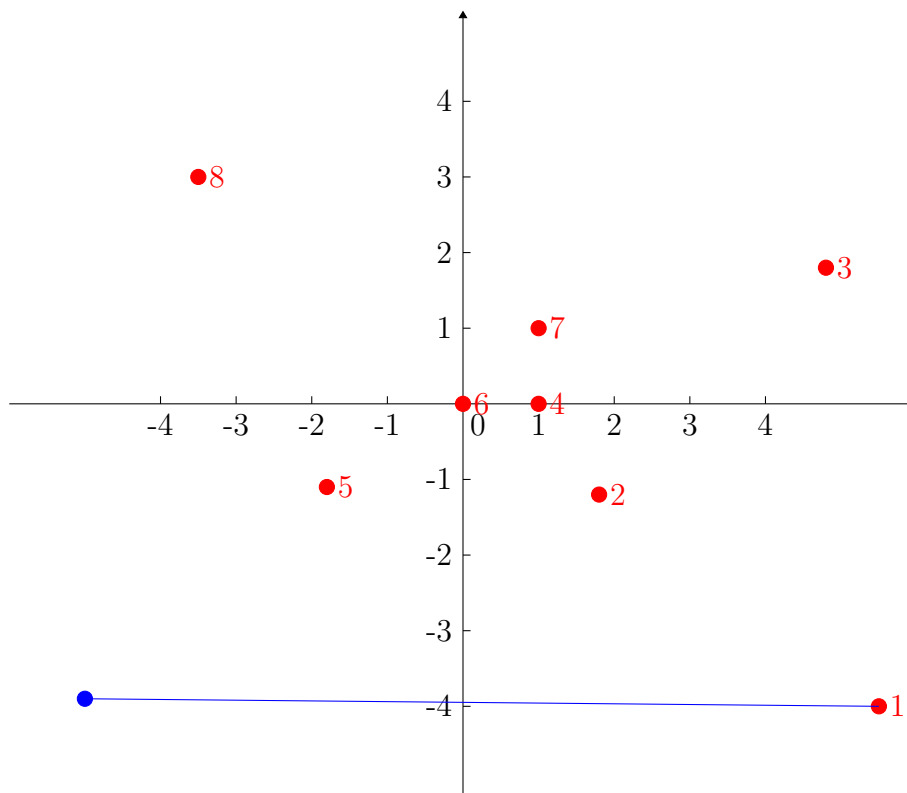
$$W(n) = \begin{cases} O(1) & n = 3 \\ W(n-1) + O(n) & n > 3 \end{cases}$$

求解得到：

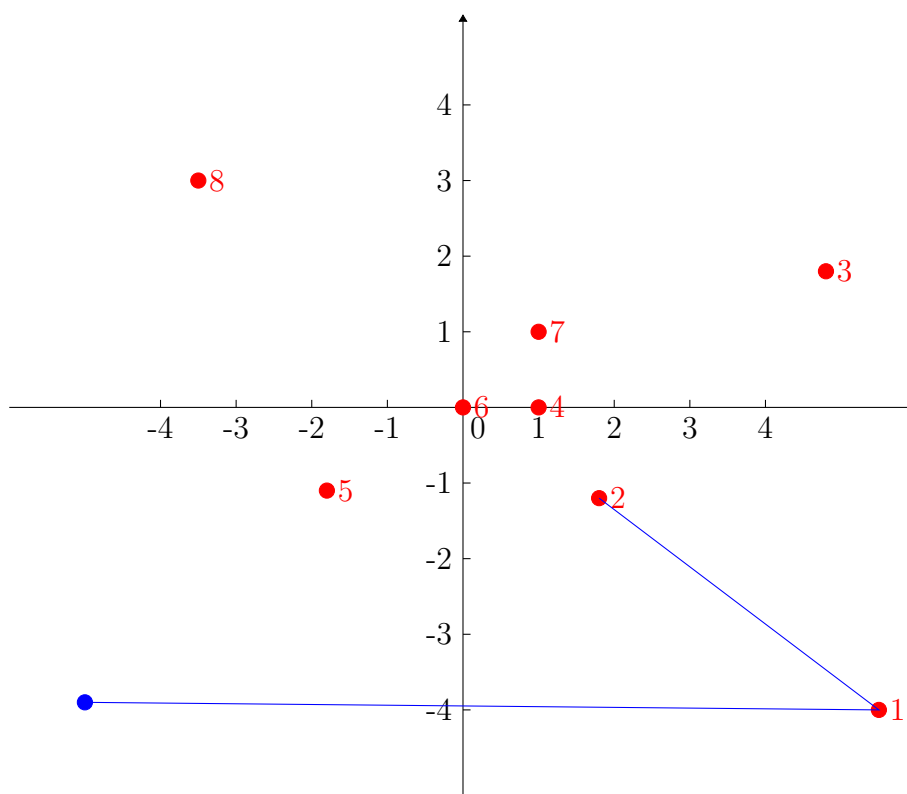
$$W(n) = O(n^2)$$

11.2.2 Graham 扫描法

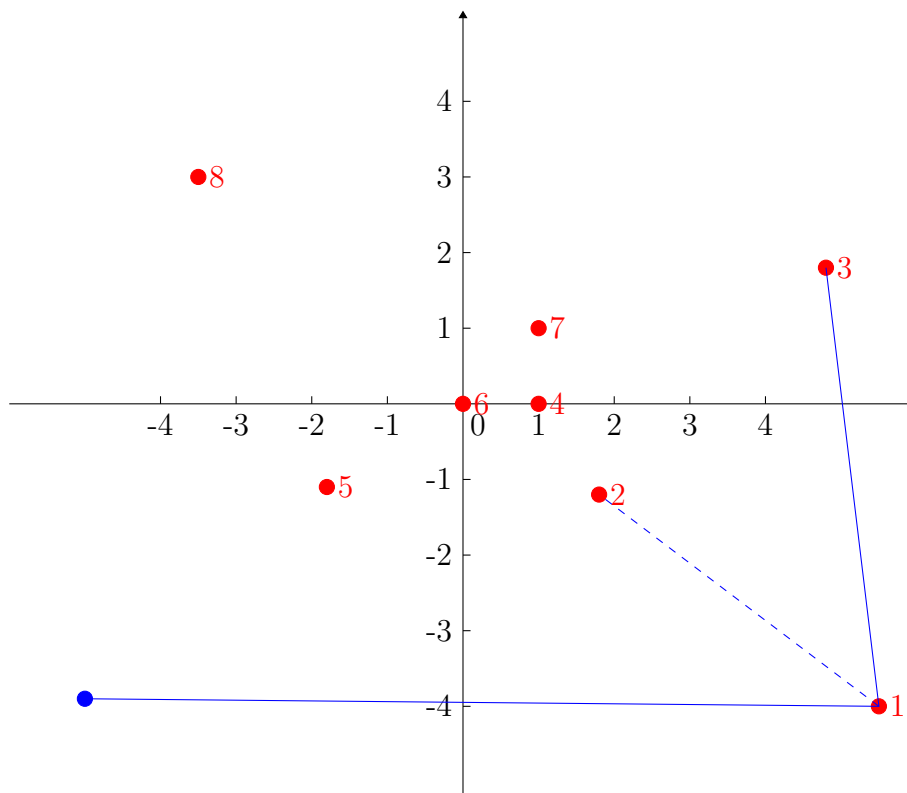
Graham 扫描法首先找到最靠近左下角的点，以这个点为极点，其它点按照极角排序。按照逆时针顺序进行扫描，先把点 1 压入凸包的栈中。



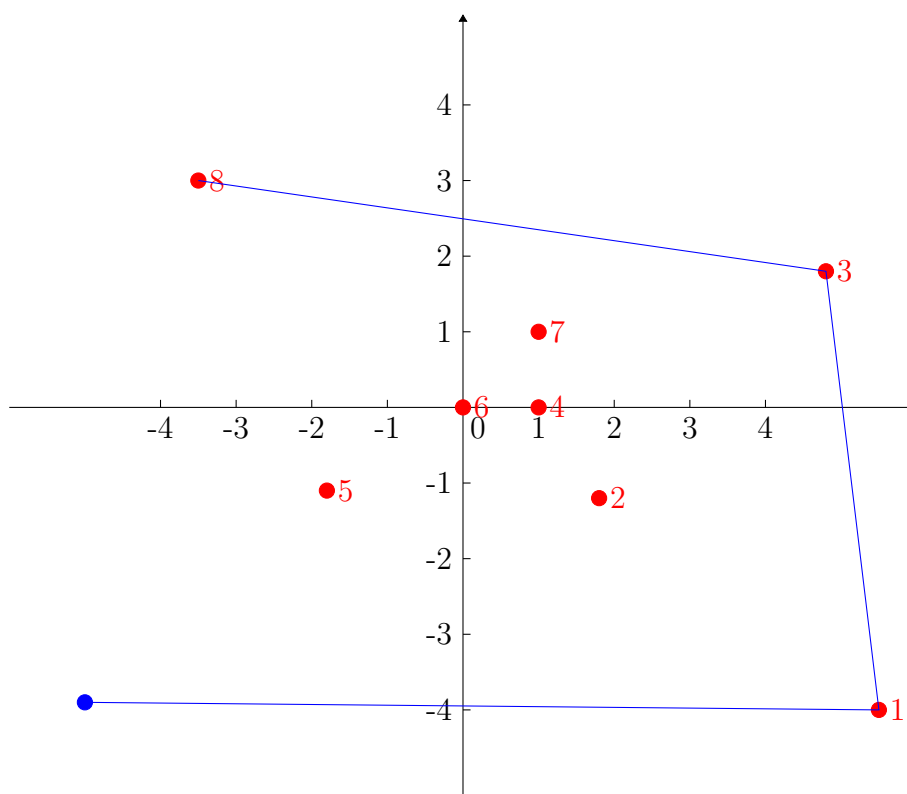
检查点 2 是否在点 1 的一侧，发现点 2 满足，入栈。



点 3 更加靠近外侧，点 2 出栈，点 3 入栈。



依次判断，最终所有凸包上的点都会在栈中。



Graham 扫描法的时间复杂度为 $O(n\log n)$ 。

11.3 芯片测试

11.3.1 芯片测试

有两片芯片 A 和 B，将两片芯片互相进行测试，测试报告是好或者坏。假设好芯片的报告一定是正确的，坏芯片的报告是不确定的（有可能会出错）。

A 报告	B 报告	结论
B 是好的	A 是好的	AB 都好或 AB 都坏
B 是好的	A 是坏的	至少一片是坏的
B 是坏的	A 是好的	至少一片是坏的
B 是坏的	A 是坏的	至少一片是坏的

表 11.1: 测试结果

芯片测试问题给定 n 片芯片，其中好芯片的数量比坏芯片的数量至少多 1 片，要求设计一种方法，用最少的测试次数，通过测试从 n 片芯片中挑出 1 片好芯片。

当需要测试某一片芯片 A 的好坏时，可以用其它 $n - 1$ 片芯片对芯片 A 进行测试。

假设 $n = 7$ （好芯片数 ≥ 4 ）：

- 如果 A 为好芯片，那么 6 个报告中至少 3 个报好。
- 如果 A 为坏芯片，那么 6 个报告状况至少 4 个报坏。

当 n 为奇数时（好芯片数 $\geq (n + 1)/2$ ）：

- 如果 A 为好芯片，至少有 $(n - 1)/2$ 个报好。
- 如果 A 为坏芯片，至少有 $(n + 1)/2$ 个报坏。

假设 $n = 8$ （好芯片数 ≥ 5 ）：

- 如果 A 为好芯片，那么 7 个报告中至少 4 个报好。
- 如果 A 为坏芯片，那么 7 个报告状况至少 5 个报坏。

当 n 为偶数时（好芯片数 $\geq n/2 + 1$ ）：

- 如果 A 为好芯片，至少有 $n/2$ 个报好。
- 如果 A 为坏芯片，至少有 $n/2 + 1$ 个报坏。

使用暴力算法，任意一个芯片进行测试，如果是好芯片，则测试结束。如果是坏芯片，则抛弃，再从剩下的芯片中任取一片进行测试，直到找到好芯片。因此暴力算法的时间复杂度为 $\Theta(n)$ 。

11.3.2 分治算法

假设 n 为偶数，将 n 个芯片两两一组做测试淘汰，剩下的芯片构成子问题，进行下一轮分组淘汰。

淘汰的规则为：

情况	操作
好 + 好	任意保留 1 片
其它情况	全部抛弃

表 11.2: 淘汰规则

递归的截止条件为 $n \leq 3$ ，因为当 $n = 1$ 或 $n = 2$ ，所有芯片一定都为好芯片，无需进行测试。当 $n = 3$ 时，只进行 1 次测试即可得到好芯片。

分治算法要求子问题与原问题性质相同，原问题中在 n 个芯片中好芯片至少比坏芯片多 1 片，因此需要验证分治算法子问题的正确性，即证明：当 n 是偶数时，经过淘汰规则后，剩下的好芯片至少比坏芯片多 1 片。

假设在分组中，A 与 B 都为好芯片的有 i 组，A 与 B 一好一坏的有 j 组，A 与 B 都为坏芯片的有 k 组，初始芯片的总数为 $2i + 2j + 2k = n$ 。

根据淘汰规则：

1. A 与 B 都为好芯片，任留 1 片，因此淘汰后好芯片至少有 i 个。
2. A 与 B 一好一坏，全部抛弃。

3. A 与 B 都为坏芯片，有可能 A 和 B 同时错报为好，因此淘汰后坏芯片至多有 k 个。

$$\text{初始好芯片多于坏芯片: } 2i + j > 2k + j$$

$$\text{子问题好芯片多于坏芯片: } i > k$$

以上算法是基于芯片数量 n 为偶数情况下的，但是如果当 n 为奇数时，可能会出现问題。

例如初始芯片为“好、好、好、好、坏、坏、坏”，两两分组得到 (好, 好)、(好, 好)、(坏、坏)、(坏)，根据规则淘汰后却得到“好、好、坏、坏”。

解决这个问题的处理办法就是额外增加一轮对多余的那个芯片单独测试。如果该芯片为好芯片，则算法结束。如果是坏芯片，则抛弃，将剩余芯片分组淘汰。

每轮分组淘汰后，芯片数量至少减半。对芯片进行测试（包括对额外多余芯片单独测试）的处理时间为 $O(n)$ 。

$$W(n) = \begin{cases} 0 & n = 1 \text{ or } n = 2 \\ 1 & n = 3 \\ W(n/2) + O(n) & n > 3 \end{cases}$$

求解得到：

$$W(n) = O(n)$$

11.4 大整数乘法

11.4.1 逐位相乘

起初对于大整数乘法，认为只要按照大整数相加的思路稍微做一下变形，就可以轻松实现。但是随着深入的学习，才发现事情并没有那么简单。如果沿用大整数加法的思路，通过列竖式求解……

$$\begin{array}{r} 93281 \\ \times 2034 \\ \hline 373124 \\ 279843 \\ 000000 \\ 186562 \\ \hline 189733554 \end{array}$$

乘法竖式的计算过程可以大体分为两步：

1. 整数 B 的每一个数位和整数 A 所有数位依次相乘，得到中间结果。
2. 所有中间结果相加，得到最终结果。

这样的做法确实可以实现大整数乘法，由于两个大整数的所有数位都需要一一彼此相乘。如果整数 A 的长度为 m ，整数 B 的长度为 n ，那么时间复杂度就是 $O(m * n)$ 。如果两个大整数的长度接近，那么时间复杂度也可以写为 $O(n^2)$ 。

那么有没有优化方法，可以让时间复杂度低于 $O(n^2)$ 呢？

11.4.2 分治法

利用分治法可以简化问题的规模，可以把大整数按照数位拆分成两部分。

$$\begin{array}{l} \text{整数 1} = \underbrace{81325}_A \underbrace{79076}_B \\ \text{整数 2} = \underbrace{9213}_C \underbrace{52184}_D \end{array}$$

即：

$$\text{整数 } 1 = A \times 10^5 + B$$

$$\text{整数 } 2 = C \times 10^5 + D$$

如果把两个大整数的长度抽象为 m 和 n ，那么：

$$\text{整数 } 1 = A \times 10^{m/2} + B$$

$$\text{整数 } 2 = C \times 10^{n/2} + D$$

因此：

$$\text{整数 } 1 \times \text{整数 } 2$$

$$= (A \times 10^{m/2} + B) \times (C \times 10^{n/2} + D)$$

$$= AC \times 10^{\frac{m+n}{2}} + AD \times 10^{\frac{n}{2}} + BC \times 10^{\frac{m}{2}} + BD$$

如此一来，原本长度为 n 的大整数的 1 次乘积，被转化成了长度为 $n/2$ 的大整数的 4 次乘积。

通过递归把大整数不断地对半拆分，一直拆分到可以直接计算为止。

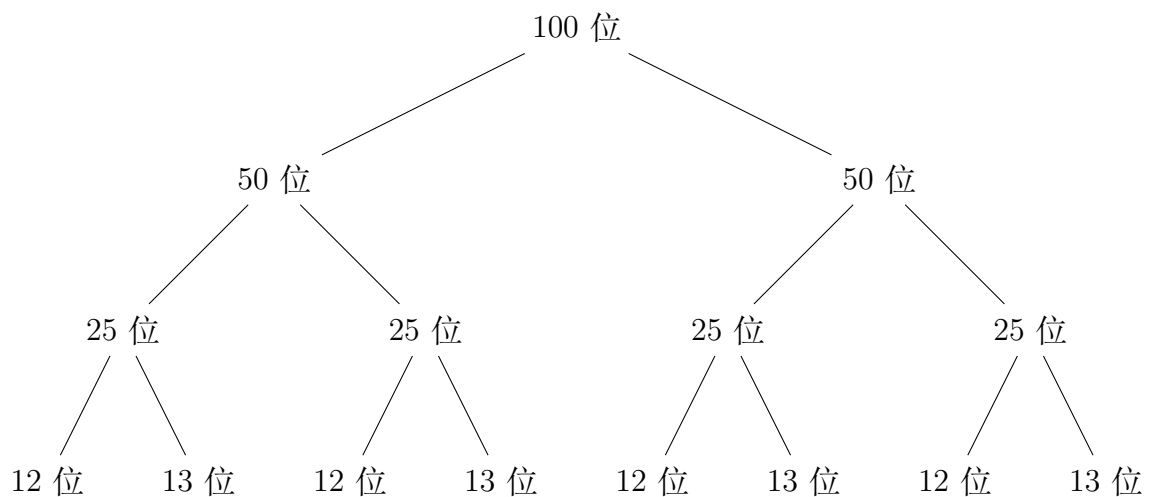


图 11.6: 大整数拆分

但是先别高兴地太早，这个方法真的提高了效率吗？

$$T(n) = \begin{cases} O(1) & n = 1 \\ 4T(n/2) + O(n) & n > 1 \end{cases}$$

求解得到：

$$T(n) = O(n^2)$$

闹了半天，时间复杂度还是 $O(n^2)$ 啊，白高兴一场。但是努力的方向并没有白费，这里还是有优化的方法的。

11.4.3 优化方法

在大整数乘法运算中，如果只简单地利用分治法将大整数的位置减半，并不能降低时间复杂度的阶。分治法把两个大整数相乘到的问题转化为四个较小整数相乘，性能的瓶颈仍旧在乘法上。

分治算法的时间复杂度方程为 $W(n) = aW(n/b) + f(n)$ ，其中 a 为子问题数， n/b 为子问题规模， $f(n)$ 为划分与合并工作量。当 a 较大、 b 较小、 $f(n)$ 不大时， $W(n) = \Theta(n^{\log_b a})$ 。减少 a 是降低 $W(n)$ 的阶的一种途径。利用子问题的依赖问题，可以使某些子问题的解通过组合其它子问题的解而得到。

那么，怎样才能减少乘法运算的次数呢？哪怕由四次乘法变成三次乘法也好呀。通过对之前的乘法等式做一些调整，可以减少乘法的次数。

整数 1 \times 整数 2

$$\begin{aligned} &= (A \times 10^{n/2} + B) \times (C \times 10^{n/2} + D) \\ &= AC \times 10^n + AD \times 10^{n/2} + BC \times 10^{n/2} + BD \\ &= AC \times 10^n + (AD + BC) \times 10^{n/2} + BD \\ &= AC \times 10^n + (AD - AC - BD + BC + AC + BD) \times 10^{n/2} + BD \\ &= AC \times 10^n + ((A - B)(D - C) + AC + BD) \times 10^{n/2} + BD \end{aligned}$$

这样一来，原本的 4 次乘法和 3 次加法，转变成了 3 次乘法和 6 次加法。

骗人！最后式子里明明包含五次乘法啊！

AC 出现了两次，BD 也出现了两次，这两个乘积分别只需计算一次就行了，所以总共只需要三次乘法。

$$T(n) = \begin{cases} O(1) & n = 1 \\ 3T(n/2) + O(n) & n > 1 \end{cases}$$

求解得到：

$$T(n) = O(n^{\log_2 3}) \approx O(n^{1.59})$$

大整数乘法

```
1 def big_int_mul(num1, num2):
2     # 有一个为空，结果为0
3     if not num1 or not num2:
4         return '0'
5     # 终止条件
6     elif len(num1) == 1 and len(num2) == 1:
7         return str(int(num1) * int(num2))
8
9     mid1 = len(num1) // 2
10    mid2 = len(num2) // 2
11
12    # 将num1分成两部分
13    a = num1[:mid1]
14    b = num1[mid1:]
15    # 将num2分成两部分
16    c = num2[:mid2]
17    d = num2[mid2:]
18
19    m = len(b)    # m次幂
20    n = len(d)    # n次幂
```

```

21
22 # 分治计算，分别补上幂次
23 x1 = big_int_mul(a, c) + '0' * (m + n)
24 x2 = big_int_mul(b, c) + '0' * n
25 x3 = big_int_mul(a, d) + '0' * m
26 x4 = big_int_mul(b, d)
27
28 # 将计算结果根据最长的补零，方便之后直接相加
29 max_len = max(len(x1), len(x2), len(x3), len(x4))
30 x1 = '0' * (max_len - len(x1)) + x1
31 x2 = '0' * (max_len - len(x2)) + x2
32 x3 = '0' * (max_len - len(x3)) + x3
33 x4 = '0' * (max_len - len(x4)) + x4
34
35 # 计算x1 + x2 + x3 + x4的值，也就是原问题的解
36 result = ""
37 carry = 0 # 保存进位
38 for i in range(max_len - 1, -1, -1):
39     s = int(x1[i]) + int(x2[i])
40         + int(x3[i]) + int(x4[i])
41         + carry
42     result = str(s % 10) + result
43     carry = s // 10
44 # 判断是否存在进位
45 if carry > 0:
46     result = str(carry) + result
47
48 # 去除结果前面多余的0
49 i = 0
50 while i < len(result) and result[i] == '0':
51     i += 1
52 return result[i:]

```

11.5 快速幂

11.5.1 快速幂 (Fast Exponentiation)

使用传统算法计算 a^n 的时间复杂度为 $\Theta(n)$ ，然而利用快速幂的算法时间复杂度为 $\Theta(\log n)$ 。

$$a^n = \begin{cases} a^{n/2} * a^{n/2} & n \text{ 为偶数} \\ a^{(n-1)/2} * a^{(n-1)/2} * a & n \text{ 为奇数} \end{cases}$$

例如计算 2^{18} 只需要 4 步即可：

$$2^{18} = 2^9 * 2^9$$

$$2^9 = 2^4 * 2^4 * 2$$

$$2^4 = 2^2 * 2^2$$

$$2^2 = 2^1 * 2$$

快速幂

```
1  /**
2   * @brief 快速幂计算a^n
3   */
4  int fastExp(int a, int n) {
5      int result = 1;
6      while(n) {
7          if(n & 1) {
8              result *= a;
9          }
10         a *= a;
11         n >>= 1;
12     }
13     return result;
14 }
```

11.5.2 矩阵快速幂

Fibonacci 数列 $\{0, 1, 1, 2, 3, 5, 8, 13, 21, \dots\}$ 可以通过递归公式 $F_n = F_{n-1} + F_{n-2}$ 计算出第 n 项的值，时间复杂度为 $\Theta(n)$ 。但是利用矩阵快速幂的算法可以在将时间复杂度降低为 $\Theta(\log n)$ 。

令 $M = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$ ，通过计算 M^n 即可计算出 F_n 的值。即：

$$\begin{bmatrix} F_{n+1} & F_n \\ F_n & F_{n-1} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}^n$$

通过数学归纳法可以证明该性质：

当 $n = 1$ 时，

$$\begin{bmatrix} F_2 & F_1 \\ F_1 & F_0 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

当 $n \geq 2$ 时，

$$\begin{aligned} & \begin{bmatrix} F_{n+2} & F_{n+1} \\ F_{n+1} & F_n \end{bmatrix} \\ &= \begin{bmatrix} F_{n+1} & F_n \\ F_n & F_{n-1} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}^n \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}^{n+1} \end{aligned}$$

矩阵快速幂

```
1 N = 2
2
3 def matrix_multiply(a, b):
4     c = [
```

```

5         [0, 0],
6         [0, 0]
7     ]
8     for i in range(N):
9         for j in range(N):
10            for k in range(N):
11                c[i][j] += a[i][k] * b[k][j]
12    return c
13
14 def matrix_fast_exp(n):
15     result = [
16         [1, 1],
17         [1, 0]
18     ]
19     M = [
20         [1, 1],
21         [1, 0]
22     ]
23
24     while n > 0:
25         if n & 1:
26             result = matrix_multiply(result, M)
27             M = matrix_multiply(M, M)
28             n >>= 1
29
30     return result[0][0]

```

11.6 矩阵乘法

11.6.1 矩阵乘法

假设 A 和 B 为 n 阶矩阵 ($n = 2^k$), 计算时, 对于 C 中 n^2 个元素, 每个元素都需要做 n 次乘法, 因此 $W(n) = O(n^3)$ 。

利用简单的分治策略, 可以将矩阵分块计算计算。

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix}$$

其中,

$$C_{11} = A_{11}B_{11} + A_{12}B_{21}$$

$$C_{12} = A_{11}B_{12} + A_{12}B_{22}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21}$$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22}$$

这样就把原问题转换为了 8 个子问题, 递推方程为:

$$W(n) = \begin{cases} 1 & n = 1 \\ 8W(n/2) + cn^2 & n > 1 \end{cases}$$

求解得到:

$$W(n) = O(n^3)$$

简单的分治算法并不能降低时间复杂度的阶, 但是矩阵乘法可以通过减少子问题的个数进行优化。

11.6.2 Strassen 矩阵乘法

让 M_1, M_2, \dots, M_7 分别对应矩阵乘法的 7 个子问题:

$$M_1 = A_{11}(B_{12} - B_{22})$$

$$M_2 = (A_{11} + A_{12})B_{22}$$

$$M_3 = (A_{21} + A_{22})B_{11}$$

$$M_4 = A_{22}(B_{21} - B_{11})$$

$$M_5 = (A_{11} + A_{22})(B_{11} + B_{22})$$

$$M_6 = (A_{12} - A_{22})(B_{21} + B_{22})$$

$$M_7 = (A_{11} - A_{21})(B_{11} + B_{12})$$

利用这些中间矩阵，可以得到结果矩阵：

$$C_{11} = M_5 + M_4 - M_2 + M_6$$

$$C_{12} = M_1 + M_2$$

$$C_{21} = M_3 + M_4$$

$$C_{22} = M_5 + M_1 - M_3 - M_7$$

在这些运算中，一共有 7 个子问题和 18 次矩阵加减法，时间复杂度为：

$$W(n) = \begin{cases} 1 & n = 1 \\ 7W(n/2) + 18(n/2)^2 & n > 1 \end{cases}$$

求解得到：

$$W(n) = O(n^{\log 7}) \approx O(n^{2.8075})$$

Coppersmith-Winograd 算法是目前已知最好的矩阵乘法算法，时间复杂度为 $O(n^{2.376})$ 。矩阵乘法可以应用在科学计算、图形处理、数据挖掘等方面，在回归、聚类、主成分分析、决策树等挖掘算法中常常涉及大规模矩阵运算。

Chapter 12 排序算法

12.1 希尔排序

12.1.1 希尔排序 (Shell Sort)

希尔排序本质上是直接插入排序的升级版。对于插入排序而言，在大多数元素已经有序的情况下，工作量会比较小。这个结论很明显，如果一个数组大部分元素都有序，那么数组中的元素自然不需要频繁地进行比较和交换。

如何能够让待排序的数组中大部分元素有序呢？需要对原始数组进行预处理，使得原始数组的大部分元素变得有序。采用分组的方法，可以将数组进行一定程度地粗略调整。

例如一个有 8 个数字组成的无序序列 {5, 8, 6, 3, 9, 2, 1, 7}，进行升序排序。让元素两两一组，同组两个元素之间的跨度为数组总长度的一半。

接着让每组元素进行独立排序，排序方式使用直接插入排序即可。由于每一组的元素数量很少，所以插入排序的工作量很少。这样一来，仅仅经过几次简单的交换，数组整体的有序程度得到了显著提高，使得后续再进行直接插入排序的工作量大大减少。

但是这样还不算完，还可以进一步缩小分组跨度，重复上述工作。

例如一个有 8 个数字组成的无序序列 {5, 8, 6, 3, 9, 2, 1, 7}，进行升序排序。

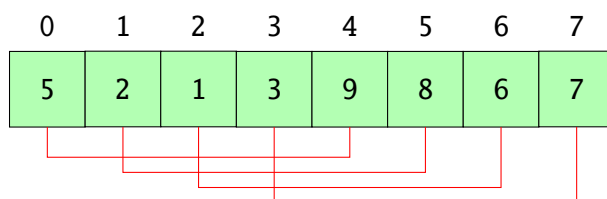


图 12.1: 跨度为 4 分组交换

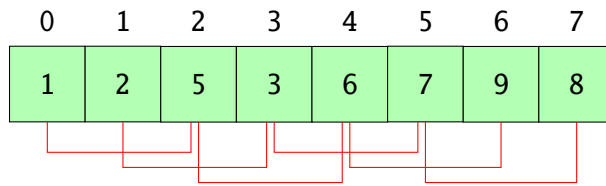


图 12.2: 跨度为 2 分组交换



图 12.3: 跨度为 1 分组交换

希尔排序的发明者是计算机科学家 Donald Shell。希尔排序中说使用的分组跨度被称为希尔排序的增量。增量的选择可以有很多种，最朴素的就是 Donald Shell 在发明希尔排序时所提出的逐步折半的方法。

希尔排序

```

1 void shellSort(int *arr, int n) {
2     int gap = n;
3     while(gap > 1) {
4         gap /= 2;
5         for(int i = 0; i < gap; i++) {
6             for(int j = i+gap; j < n; j += gap) {
7                 int temp = arr[j];
8                 int k = j - gap;
9                 while(k >= 0 && arr[k] > temp) {
10                     arr[k+gap] = arr[k];
11                     k -= gap;
12                 }
13                 arr[k+gap] = temp;
14             }
15         }
16     }
17 }

```

12.1.2 算法分析

希尔排序利用分组粗略调整的方式减少了直接插入排序的工作量，使得算法的平均时间复杂度低于 $O(n^2)$ 。但是在某些极端情况下，希尔排序的最坏时间复杂度仍然是 $O(n^2)$ ，甚至比插入排序更慢。

例如 {2, 1, 5, 3, 7, 6, 9, 8}，无论是以 4 为增量，还是以 2 为增量，每组内部的元素都没有任何交换。直到增量缩减为 1，数组才会按照直接插入排序的方式进行调整。

对于这样的数组，希尔排序不但没有减少直接插入排序的工作量，反而白白增加了分组操作的成本。

这是因为每一轮希尔增量之间都是等比的，这就导致了希尔增量存在盲区。为了避免这样的极端情况，科学家发明了许多更为严谨的增量方式。其中最具有代表性的是 Hibbard 增量和 Sedgewick 增量。

Hibbard 增量序列

Hibbard 增量序列为 1, 3, 7, 15, ..., 通项公式为 $2^i - 1$ 。

利用这种增量方式的希尔排序，最坏时间复杂度是 $O(n^{3/2})$ 。

Hibbard 增量序列

```
1 def get_hibbard_sequence(n):
2     """
3     生成Hibbard序列
4     1, 3, 7, 15, 31, 63, ...
5     """
6     sequence = []
7     i = 1
8     while i <= n:
9         sequence.append(i)
```

```

10         i = (i << 1) + 1
11     sequence.reverse()
12     return sequence
13
14 def shell_sort_hibbard(lst):
15     """
16     希尔排序 (Hibbard增量序列)
17     """
18     n = len(lst)
19     hibbard = get_hibbard_sequence(n)
20     for gap in hibbard:
21         for i in range(gap, n):
22             j = i
23             temp = lst[j]
24             while j >= gap:
25                 if temp < lst[j-gap]:
26                     lst[j] = lst[j-gap]
27                     j -= gap
28             else:
29                 break
30             lst[j] = temp

```

Sedgewick 增量序列

Sedgewick 增量序列为 1, 5, 19, 41, 109, ..., 通项公式为 $9 \times 4^i - 9 \times 2^i + 1$ 和 $4^{i+2} - 3 \times 2^{i+2} + 1$ 。

利用这种增量方式的希尔排序，最坏时间复杂度是 $O(n^{4/3})$ 。

Sedgewick 增量序列

```

1 def get_sedgewick_sequence(n):
2     """
3     生成Sedgewick序列
4     1, 5, 19, 41, 109, ...
5     """

```

```

6     sequence = []
7     i = 0
8     while True:
9         #  $9 * 4^i - 9 * 2^i + 1$ 
10        # ==>  $9 * (2^{(2*i)} - 2^i) + 1$ 
11        item = 9 * ((1 << (2 * i)) - (1 << i)) + 1
12        if item <= n:
13            sequence.append(item)
14        else:
15            break
16
17        #  $4^{(i+2)} - 3 * 2^{(i+2)} + 1$ 
18        # ==>  $2^{(2i+4)} - 3 * 2^{(i+2)} + 1$ 
19        item = (1 << (2 * i + 4)) - 3 * (1 << (i + 2)) + 1
20        if item <= n:
21            sequence.append(item)
22        else:
23            break
24
25        i += 1
26    return sequence
27
28 def shell_sort_sedgewick(lst):
29     """
30     希尔排序 (Sedgewick增量序列)
31     """
32     n = len(lst)
33     sedgewick = get_sedgewick_sequence(n)
34     for gap in sedgewick:
35         for i in range(gap, n):
36             j = i
37             temp = lst[j]
38             while j >= gap:
39                 if temp < lst[j-gap]:
40                     lst[j] = lst[j-gap]
41                     j -= gap
42             else:

```

```
43         break
44     lst[j] = temp
```

这两种增量方式的时间复杂度需要很复杂的数学证明，有些是人们的大致猜想。

时间复杂度	空间复杂度	稳定性
$O(n^{1.3\sim 2})$	$O(1)$	不稳定

表 12.1: 希尔排序算法分析

12.2 归并排序

12.2.1 归并排序优化

简单的归并排序利用分治法，递归地将对小规模子数组进行处理。但是递归会使小规模问题中方法调用太过频繁，因此对于规模较小的子数组可以采用插入排序。一般来说插入排序在小数组中比归并更快，这种优化可以使归并排序的运行时间缩短 10% ~ 15%。

另一个可以优化的地方是对于单次合并的过程，例如将子数组 $arr[start..mid]$ 和 $arr[mid + 1..end]$ 进行合并，如果 $arr[mid] \leq arr[mid + 1]$ 的话，说明 $arr[start..end]$ 已经为有序状态，无序再进行不必要的合并。

归并排序优化

```
1 void mergeSortWorker(int *arr, int start, int end, int *temp) {
2     // 列表长度小于10时，采用二分插入排序
3     if(end - start <= 10) {
4         binaryInsertionSort(arr, start, end);
5         return;
6     }
7     if(start < end) {
8         int mid = start + (end - start) / 2;
9         mergeSortWorker(arr, start, mid, temp);
10        mergeSortWorker(arr, mid+1, end, temp);
11        // 避免不必要的合并
12        if(arr[mid] <= arr[mid+1]) {
13            return;
14        }
15        merge(arr, start, mid, end, temp);
16    }
17 }
```

12.2.2 归并排序迭代实现

递归实现的归并排序是自顶向下的过程，基于循环的归并排序是自底向上进行的。非递归的归并排序避免了递归时深度为 $\log n$ 的栈空间，空间上只用到了长度为 n 的临时空间。

归并排序（迭代）

```
1 public static void mergeSort(int[] arr) {
2     int n = arr.length;
3     int[] temp = new int[n];
4     int pos = 0;           // 临时数组的下表
5     int left1, left2;      // 左子数组边界
6     int right1, right2;    // 右子数组边界
7
8     for (int i = 1; i < n; i *= 2) {
9         for (left1 = 0; left1 < n - i; left1 = right2) {
10             // 设置子数组边界
11             right1 = left2 = left1 + i;
12             right2 = left2 + i;
13
14             // 防止右边界越界
15             right2 = right2 > n ? n : right2;
16
17             pos = 0;
18             while (left1 < left2 && right1 < right2) {
19                 if (arr[left1] < arr[right1]) {
20                     temp[pos++] = arr[left1++];
21                 } else {
22                     temp[pos++] = arr[right1++];
23                 }
24             }
25
26             while (left1 < left2) {
27                 arr[--right1] = arr[--left2];
28             }
```

```

29
30         // 将排好序的部分保存回数组
31         while (pos > 0) {
32             arr[--right1] = temp[--pos];
33         }
34     }
35 }
36 }

```

12.2.3 外部排序

在内存中进行的排序称为内部排序，而在许多实际应用中，经常需要对大文件进行排序。因为文件中的信息量庞大，无法将整个文件拷贝进内存进行排序。因此需要将待排序的记录存储在外存上，排序时再把数据一部分一部分调入内存进行排序，再将排好序的记录写回文件中。

因为磁盘读写的时间远超过内存计算的时间，因此外部排序过程中的时间代价主要是磁盘 I/O 次数。

假如需要对一个包含 40 亿个 int 类型整数的文件进行排序，而计算机的内存只有 2GB。一个 int 占 4 个字节，40 亿个需要 160 亿字节，大概占用 8GB 的内存。因此可以把 8GB 分割成 4 份 2GB 的数据进行排序，然后再把它们凑回去。

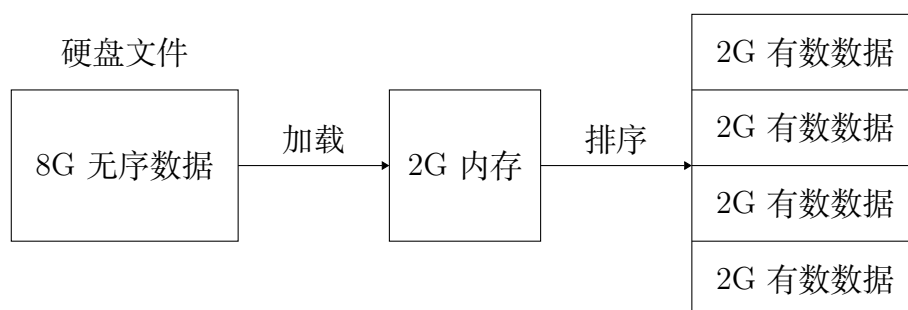


图 12.4: 分割数据

排序的时候可以采用归并排序，每次将两个有序子串合并成一个大的有序子串。

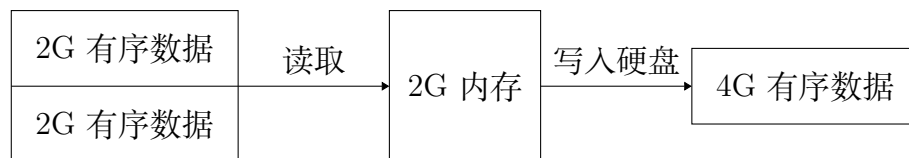


图 12.5: 二路归并

不过硬盘的读写速度比内存要慢得多，通过优化可以降低数据从硬盘读写的次数。

在进行有序数据合并的时候，不采取两两合并的方法，而是可以三组或四组数据一起合并。 n 个有序数据的合并被称为 n 路归并。

12.3 快速排序

12.3.1 随机选择基准值

快速排序利用分治法，通过一趟排序将数组分为两部分，其中一部分小于等于基准值，另一部分大于等于基准值，然后再递归对两个子问题排序。

基本的快速排序采用序列的第一个元素作为基准值，但是这不是一种好方法。当数组已经有序时，这样的分割效率非常糟糕。为了缓解这种极端情况，可以在待排序数组中随机选择一个元素作为基准值。

随机选取基准值

```
1 int selectRandomPivot(int *arr, int start, int end) {  
2     srand(time(NULL));  
3     int pos = rand() % (end - start) + start;  
4     swap(&arr[pos], &arr[start]);  
5     return arr[start];  
6 }
```

12.3.2 三数取中

虽然随机选取基准值可以减少出现分割不好的几率，但是最坏情况下还是 $O(n)$ 。另一种选取基准值的方法就是三数取中，也就是取序列中 start、mid、end 三个元素的中间值作为基准值。

三数取中

```
1 int selectMedianPivot(int *arr, int start, int end) {  
2     int mid = start + (end - start) / 2;  
3     if(arr[mid] > arr[end]) {  
4         swap(&arr[mid], &arr[end]);  
5     }  
}
```

```

6     if(arr[start] > arr[end]) {
7         swap(&arr[start], &arr[end]);
8     }
9     if(arr[mid] > arr[start]) {
10        swap(&arr[mid], &arr[start]);
11    }
12    // 此时arr[mid] <= arr[start] <= arr[end]
13    return arr[start];
14 }

```

12.3.3 三数取中 + 插入排序

对于很小和部分有序的数组，快速排序的效率不如插入排序。因此当待排序数组被分割到一定大小后，可直接采用插入排序。

三数取中 + 插入排序

```

1 void quickSort(int *arr, int start, int end) {
2     if(end - start <= 10) {
3         binaryInsertionSort(arr, start, end);
4         return;
5     }
6
7     if(start < end) {
8         int i = start;
9         int j = end;
10        int pivot = selectMedianPivot(arr, start, end);
11
12        while(i < j) {
13            while(i < j && arr[j] > pivot) {
14                j--;
15            }
16            if(i < j) {
17                arr[i] = arr[j];
18                i++;

```

```

19         }
20         while(i < j && arr[i] < pivot) {
21             i++;
22         }
23         if(i < j) {
24             arr[j] = arr[i];
25             j--;
26         }
27     }
28     arr[i] = pivot;
29     quickSort(arr, start, i-1);
30     quickSort(arr, i+1, end);
31 }
32 }

```

12.3.4 聚集相等基准值

在一次分割结束后，可以把所有与基准值相等的元素聚集在一起，这样在下次分割时，就不用对这些值再分割了。

例如待排序序列为 {1, 4, 6, 7, 6, 6, 7, 6, 8, 6}，选择 6（下标为 4）作为基准值。在进行一次分割后，得到两个子序列 {1, 4, 6} 和 {7, 6, 7, 6, 8, 6}。将所有与基准值相等的元素聚集后，可得到 {1, 4, 6, 6, 6, 6, 6, 7, 8, 7}。这样下一次分割的子序列可以减少为 {1, 4} 和 {7, 8, 7}。

聚集相等基准值

```

1  /**
2   * @brief 聚集相等基准值
3   * @param arr: 待排序数组
4   * @param start: 数组开始位置
5   * @param end: 数组结束位置
6   * @param pivotPos: 基准值下标
7   * @param left: 相等基准值左边界

```

```

8  * @param right: 相等基准值右边界
9  */
10 void gather(int *arr, int start, int end,
11             int pivotPos, int *left, int *right) {
12     if(start >= end) {
13         return;
14     }
15
16     int cnt = pivotPos - 1;
17     for(int i = pivotPos - 1; i >= start; i--) {
18         if(arr[i] == arr[pivotPos]) {
19             swap(&arr[i], &arr[cnt]);
20             cnt--;
21         }
22     }
23     *left = cnt;
24
25     cnt = pivotPos + 1;
26     for(int i = pivotPos + 1; i <= end; i++) {
27         if(arr[i] == arr[pivotPos]) {
28             swap(&arr[i], &arr[cnt]);
29             cnt++;
30         }
31     }
32     *right = cnt;
33 }

```

12.3.5 尾递归优化

快速排序在函数尾部有 2 次递归操作，可以对其中的尾递归进行优化。因为在第一次递归后，start 就没用了，第二次递归可以用循环代替。

尾递归优化

```

1 void quickSort(int *arr, int start, int end) {

```

```

2     if(end - start <= 10) {
3         binaryInsertionSort(arr, start, end);
4         return;
5     }
6
7     while(start < end) {
8         int i = start;
9         int j = end;
10        int pivot = selectMedianPivot(arr, start, end);
11
12        while(i < j) {
13            while(i < j && arr[j] > pivot) {
14                j--;
15            }
16            if(i < j) {
17                arr[i] = arr[j];
18                i++;
19            }
20            while(i < j && arr[i] < pivot) {
21                i++;
22            }
23            if(i < j) {
24                arr[j] = arr[i];
25                j--;
26            }
27        }
28        arr[i] = pivot;
29
30        // 聚集与基准值相等元素
31        int left, right;
32        gather(arr, start, end, i, &left, &right);
33
34        quickSort(arr, start, left);
35        // quickSort(arr, right, end); // 消除尾递归
36        start = right;
37    }
38 }

```

其实这种优化编译器会自己进行优化，因此相比不使用优化的方法，运行时间几乎无异。

12.3.6 快速排序迭代实现

递归实现主要是在划分子区间，因此可以通过利用栈的特性来保存区间即可，因为递归本身就是压栈的过程。

快速排序（迭代）

```
1 int partition(int *arr, int start, int end) {
2     int i = start - 1;
3     int pivot = arr[end];
4
5     for(int j = start; j < end; j++) {
6         if(arr[j] <= pivot) {
7             i++;
8             swap(&arr[i], &arr[j]);
9         }
10    }
11
12    swap(&arr[i+1], &arr[end]);
13    return i + 1;
14 }
15
16 void quickSort(int *arr, int start, int end) {
17     Stack *s = initStack(end - start + 1);
18     push(s, start);
19     push(s, end);
20
21     while(!isEmptyStack(s)) {
22         int right = pop(s);
23         int left = pop(s);
24
25         int index = partition(arr, left, right);
```

```
26     if(index - 1 > left) {
27         push(s, left);
28         push(s, index - 1);
29     }
30     if(index + 1 < right) {
31         push(s, index + 1);
32         push(s, right);
33     }
34 }
35 }
```


12.4 堆排序

12.4.1 堆 (Heap)

二叉堆本质上是一种完全二叉树，分为最大堆和最小堆两个类型。在最大堆中，任何一个父结点的值都大于等于它左右孩子结点的值。在最小堆中，任何一个父结点的值都小于等于它左右孩子结点的值。

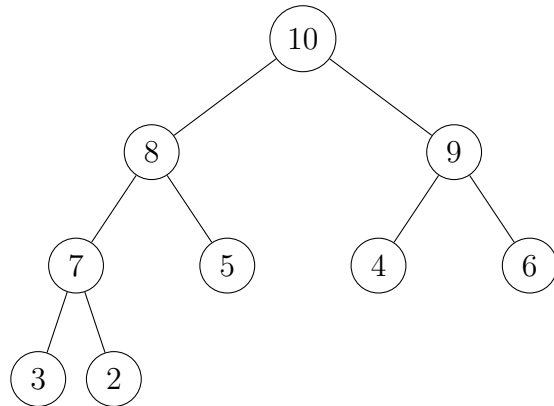


图 12.6: 大顶堆

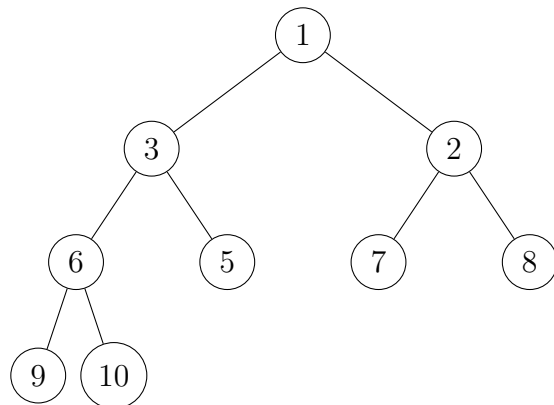


图 12.7: 小顶堆

二叉堆的根结点称为堆顶，在最大堆中堆顶是整个堆中的最大元素，在最小堆中堆顶是整个堆中的最小元素。

在二叉堆中插入结点、删除结点、构造二叉堆的操作都基于堆的自我调整。

二叉堆虽然是一棵完全二叉树，但它的存储方式并不是链式存储，而是顺序存储。

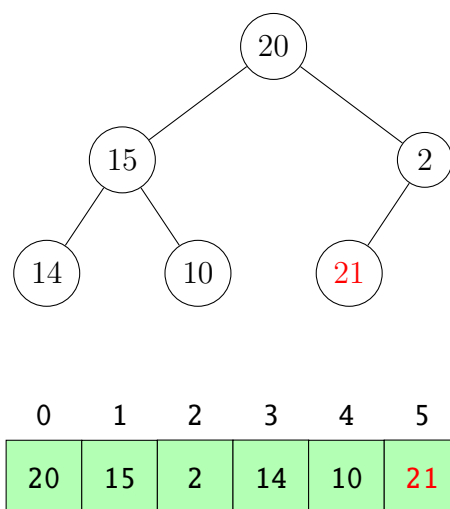
数组中，通过下标可以定位到结点的左右孩子，假设父结点的下标是 parent ，那么它的左孩子下标为 $2 * \text{parent} + 1$ 、右孩子下标为 $2 * \text{parent} + 2$ 。

12.4.2 插入结点

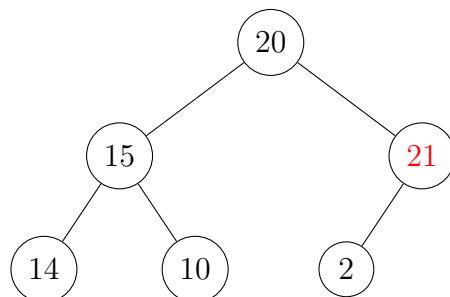
二叉堆的插入操作可以看成是结点上浮，当在堆中插入一个结点时，必须满足完全二叉树的标准，那么被插入结点的位置是完全二叉树的最后一个位置。在最大堆中，如果新结点的值大于它的父结点的值，则让新结点上浮，即和父结点交换位置。

堆的插入时间复杂度取决于树高为 $O(\log n)$ 。

例如在大顶堆 $\{20, 15, 2, 14, 10\}$ 中插入 21：

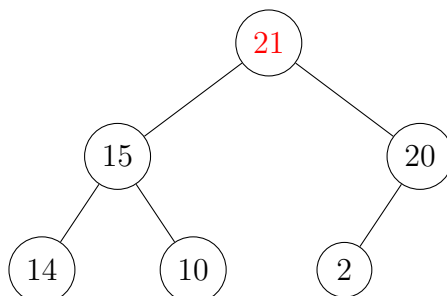


将新元素 21 与其父结点 2 比较，因为 $21 > 2$ ，将 21 和 2 的位置交换：



0	1	2	3	4	5
20	15	21	14	10	2

因为 $21 > 20$ ，将 21 与 20 的位置交换：



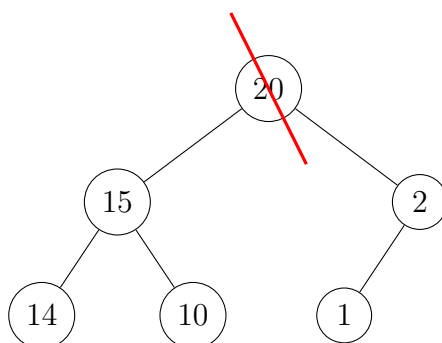
0	1	2	3	4	5
21	15	20	14	10	2

12.4.3 删除结点

二叉堆的删除操作总是从堆的根结点删除元素。根结点被删除之后为了保证该树还是一棵完全二叉树，需要将完全二叉树的最后一个结点补到根结点的位置，让其继续符合完全二叉树的定义。二叉堆的删除结点操作可以看作是结点下沉。在最大堆中，如果新堆顶元素小于它的左右孩子中较大的那个结点，则与它的较大的子结点交换位置。

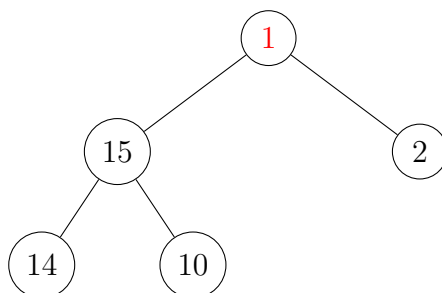
堆的删除时间复杂度取决于树高为 $O(\log n)$ 。

例如删除大顶堆的堆顶元素 20：



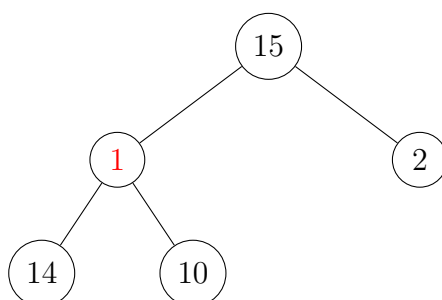
0	1	2	3	4	5
20	15	2	14	10	1

移动最后一个结点到堆顶，使其满足二叉树的性质：



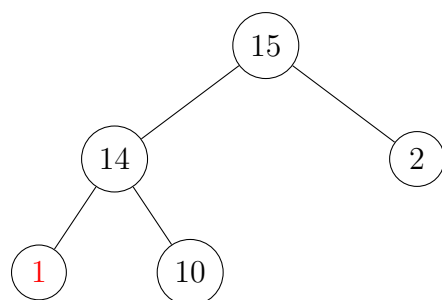
0	1	2	3	4
1	15	2	14	10

将堆顶元素 1 与其子结点比较，因为 $15 > 2$ ，交换较大子结点 15 与 1 的位置：



0	1	2	3	4
15	1	2	14	10

继续将元素 1 与其子结点比较，因为 $14 > 10$ ，交换较大子结点 14 与 1 的位置：

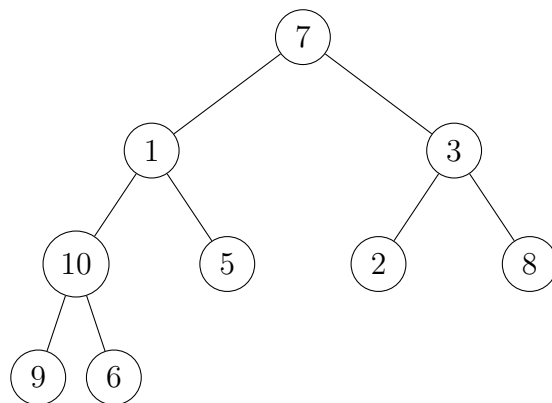


0	1	2	3	4
15	1	2	1	10

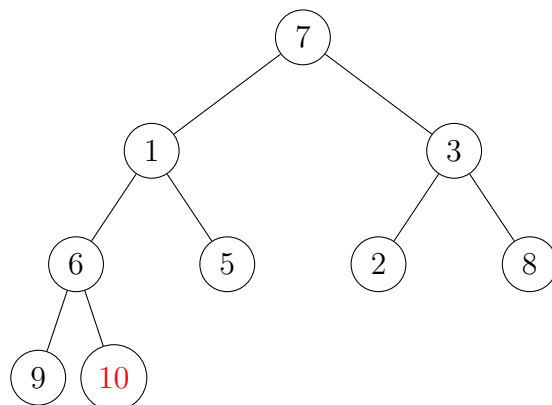
12.4.4 构建二叉堆

构建二叉堆，就是把一个无序的完全二叉树调整为二叉堆，本质上就是让所有非叶子结点依次下沉。

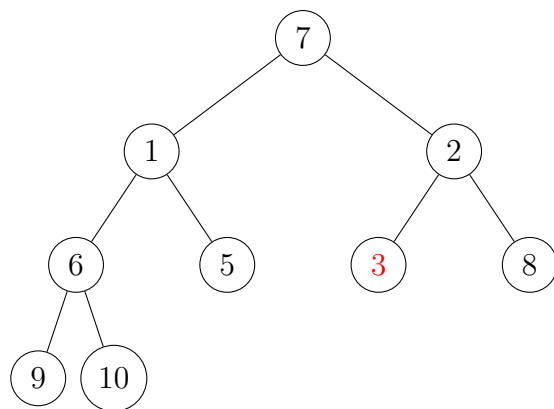
例如将一个无序的完全二叉树构建成最小堆：



首先从最后一个非叶子结点开始，结点 10 下沉：

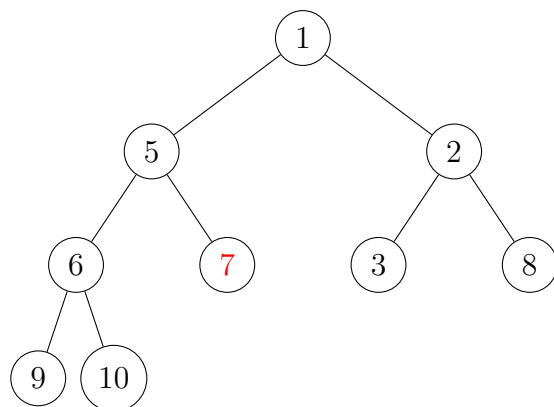


接着处理倒数第二个非叶子结点，结点 3 下沉：



倒数第三个非叶子结点 1 无需移动。

最后处理倒数第四个非叶子结点 7 下沉：

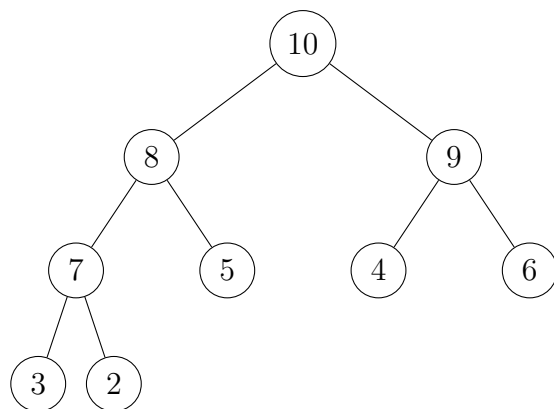


最终一棵无序完全二叉树就调整成了一个最小堆。

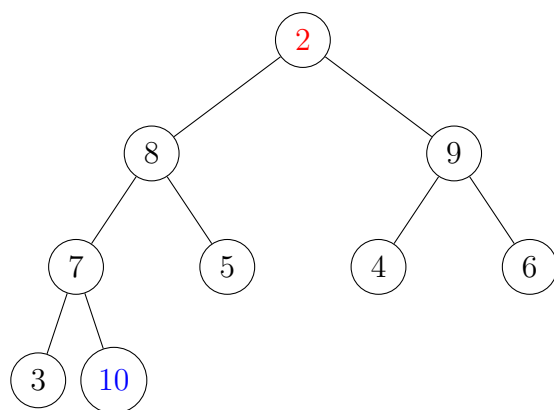
12.4.5 堆排序 (Heap Sort)

有了二叉堆的构建、删除和自我调节，实现堆排序就是水到渠成了。当删除一个最大堆的堆顶后（并不是完全删除，而是替换到堆的最后面），经过自我调节，第二大的元素就会被交换上来，成为最大堆的新堆顶。

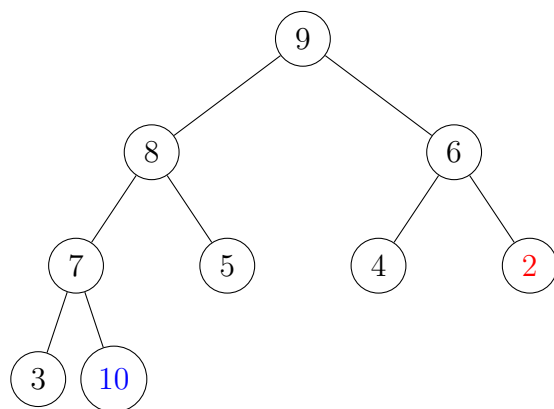
首先将待排序数组构建成大顶堆：



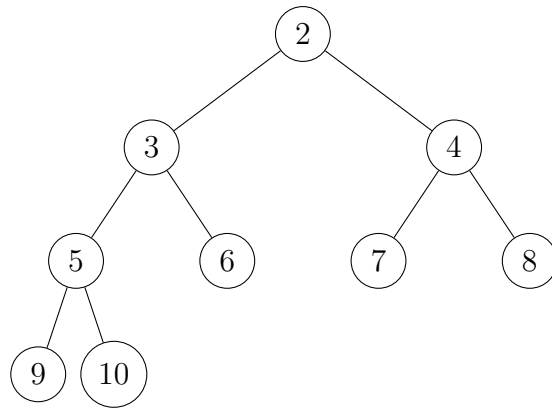
移除堆顶元素（与最后元素交换）：



将新的堆顶元素进行下沉，重新调整为大顶堆：



只要反复删除堆顶，反复调节二叉堆，所得到的集合就成为了一个有序集合。



堆排序

```
1 public static void downAdjust(int[] arr, int parentIndex, int len) {
2     // 保存父结点的值，用于最后的赋值
3     int temp = arr[parentIndex];
4     int childIndex = 2 * parentIndex + 1;
5
6     while(childIndex < len) {
7         // 如果有右孩子，且右孩子大于左孩子的值，则定位到右孩子
8         if(childIndex + 1 < len
9             && arr[childIndex + 1] > arr[childIndex]) {
10             childIndex++;
11         }
12         // 如果父结点小于任何一个孩子的值，直接跳出
13         if(temp >= arr[childIndex]) {
14             break;
15         }
16         // 无需真正交换，单向赋值即可
17         arr[parentIndex] = arr[childIndex];
18         parentIndex = childIndex;
19         childIndex = 2 * childIndex + 1;
20     }
21     arr[parentIndex] = temp;
22 }
23
24 public static void heapSort(int[] arr) {
25     // 把无序数组构建成二叉堆
```



```

26     for(int i = (arr.length-2) / 2; i >= 0; i--) {
27         downAdjust(arr, i, arr.length);
28     }
29
30     // 循环删除堆顶元素，移到数组尾部，调节堆产生新的堆顶
31     for(int i = arr.length - 1; i > 0; i--) {
32         // 最后一个元素和第一个元素交换
33         int temp = arr[i];
34         arr[i] = arr[0];
35         arr[0] = temp;
36         // 下沉调整最大堆
37         downAdjust(arr, 0, i);
38     }
39 }

```

堆排序的空间复杂度为 $O(1)$ ，因为算法并没有开辟额外的集合空间。

至于空间复杂度，假设二叉堆总共有 n 个元素，那么下沉调整的最坏时间复杂度就等同于二叉堆的高度 $O(\log n)$ 。

堆排序的算法步骤分为两部分：

1. 把无序数组构建成二叉堆：进行 $n/2$ 次循环，每次循环进行一次下沉调节，因为此步骤的计算规模为 $n/2 * \log n$ ，时间复杂度为 $O(n \log n)$ 。
2. 循环删除堆顶元素，移到数组尾部，调节堆产生新堆顶：进行 $n - 1$ 次循环，每次循环进行一次下沉调节，因此此步骤的计算规模为 $(n - 1) * \log n$ ，时间复杂度为 $O(n \log n)$ 。

综合堆排序的两个步骤，整体时间复杂度为 $O(n \log n)$ 。

Chapter 13 贪心算法

13.1 贪心算法

13.1.1 贪心算法 (Greedy Algorithm)

贪心算法，又称贪婪算法，是指在对问题求解时，总是做出在当前看来是最好的选择，也就是说不从整体最优上加以考虑。算法得到的是在某种意义上的局部最优解。贪心算法不是对所有问题都能得到整体最优解，关键在于贪心策略的选择。

贪心算法没有固定的算法框架，算法设计的关键是贪心策略的选择。必须注意的是，贪心算法不是对所有问题都能得到整体最优解，选择的贪心策略必须具备无后效性，即某个状态以后的过程不会影响以前的状态，只与当前状态有关。

在利用贪心算法求解问题之前，必须需要清楚什么样的问题适合用贪心算法。一般而言，能够利用贪心算法求解的问题都会具备以下两点性质：

1. 贪心选择：当某一个问题的整体最优解可通过一系列局部最优解的选择达到，并且每次做出的选择可以依赖以前做出的选择，但不需要依赖后面需要做出的选择。
2. 最优子结构：如果一个问题的最优解包含其子问题的最优解，则此问题具备最优子结构的性质。

贪心算法的基本思路分为：

1. 建立数学模型描述问题。
2. 把求解的问题分成若干个子问题。
3. 对每个子问题求解，得到子问题的局部最优解。
4. 把子问题的局部最优解合成为原问题的解。

买卖股票的最佳时期

给定一个数组，它的第 i 个元素是一支给定股票第 i 天的价格。设计一个算法来计算所能获取的最大利润。（提示：可以尽可能地完成更多的交易，但不能同时参与多笔交易，在再次购买前出售掉之前的股票）。

示例

输入: [7, 1, 5, 3, 6, 4]

输出: 7

解释: 在第 2 天（股票价格 = 1）的时候买入，在第 3 天（股票价格 = 5）的时候卖出，这笔交易所能获得利润 = $5 - 1 = 4$ 。随后，在第 4 天（股票价格 = 3）的时候买入，在第 5 天（股票价格 = 6）的时候卖出，这笔交易所能获得利润 = $6 - 3 = 3$ 。

```
1 public static int maxProfit(int[] prices) {  
2     int profit = 0;  
3     for(int i = 0; i < prices.length - 1; i++) {  
4         if(prices[i] < prices[i+1]) {  
5             profit += prices[i+1] - prices[i];  
6         }  
7     }  
8     return profit;  
9 }
```

分发饼干

有一群孩子和一堆饼干，每个孩子有一个饥饿度，每个饼干都有大小。每个孩子最多只能吃一个饼干，且只有饼干的大小大于等于孩子的饥饿度时，这个孩子才能吃饱。求解最多有多少孩子可以吃饱。

示例 1

输入: children = [1, 2, 3], cookies = [1, 1]

输出: 1

解释：有三个孩子和两块饼干，3 个孩子的饥饿度分别是 1, 2, 3。虽然有两块饼干，由于它们的大小都是 1，只能让饥饿度是 1 的孩子满足。所以应该输出 1。

示例 2

输入：children = [1, 2], cookies = [1, 2, 3]

输出：2

解释：有两个孩子和三块饼干，2 个孩子的饥饿度分别是 1, 2。拥有的饼干数量和大小都足以让所有孩子满足。所以应该输出 2。

因为饥饿度最小的孩子最容易吃饱，所以先考虑这个孩子。为了尽量使得剩下的饼干可以满足饥饿度更大的孩子，所以应该把大于等于这个孩子饥饿度的、且大小最小的饼干给这个孩子。满足了这个孩子之后，采取同样的策略，考虑剩下孩子里饥饿度最小的孩子，直到没有满足条件的饼干存在。

这里的贪心策略是，给剩余孩子里最小饥饿度的孩子分配最小的能饱腹的饼干。因为需要获得大小关系，一个便捷的方法就是把孩子和饼干分别排序，这样就可以从饥饿度最小的孩子和大小最小的饼干出发，计算有多少个孩子可以吃饱。

```
1 public static int distribute(int[] children, int[] cookies) {
2     Arrays.sort(children);
3     Arrays.sort(cookies);
4     int child = 0;
5     int cookie = 0;
6     while (child < children.length && cookie < cookies.length) {
7         if (children[child] <= cookies[cookie++]) {
8             child++;
9         }
10    }
11    return child;
12 }
```

硬币找零

假设硬币的面值分别为 1 元、5 元、10 元、25 元和 100 元，有一个小店拥有各

面值硬币数量无限个，顾客买东西之后需要给顾客找零，如何能够使找零的硬币个数最少？

示例

输入：36

输出：3

解释：找零 [25, 10, 1]

为了尽量减少硬币的数量，首先得尽可能地多使用面值大的硬币，也就是优先使用大面值的硬币。

```
1 def get_min_coins(coins, price):
2     """
3     最少硬币找零
4     Args:
5         coins (list): 硬币面值
6         price (int): 找零数量
7     """
8     solution = []
9     coins = sorted(coins, reverse=True)
10    for coin in coins:
11        num = price // coin
12        solution += [coin] * num
13        price -= coin * num
14        if price <= 0:
15            break
16    return solution
```

13.2 活动安排

13.2.1 活动安排

假设有 n 个活动的集合 $S = \{a_1, a_2, a_3, \dots, a_n\}$ ，其中每个活动都要使用同一种资源（如场馆等），而在同一时间内只能有一个活动使用这个资源。每个活动 a_i 都有一个开始时间 s_i 和结束时间 f_i ($s_i < f_i$)。如果某个活动 a_i 被选择，那么活动 a_i 的发生时间在半开区间 $[s_i, f_i)$ 内。如果两个活动 a_i 和 a_j 满足 $[s_i, f_i)$ 和 $[s_j, f_j)$ 不重叠，则称它们是兼容的。活动安排问题就是要求选出最大兼容活动集。

假设活动已经按照结束时间的单调递增顺序排序好：

i	1	2	3	4	5	6	7	8	9	10	11
s_i	1	3	0	5	3	5	6	8	8	2	12
f_i	4	5	6	7	9	9	10	11	12	14	16

表 13.1: 活动安排

对于活动安排问题的贪心选择，应该选择一个这样的活动，选择它之后剩下的资源可以被尽量多的其它活动所占用。活动安排问题中的贪心选择就是每次选择最先结束的活动。

一个贪心选择的日常应用场景就比如有多个不同的兼职岗位，每个岗位都有一个开始时间和结束时间，小灰在同一时间内只能做一份兼职，小灰每天最多可以做多少份兼职。面对这样的问题，在日常生活中，我们的第一选择肯定是先把结束时间早的兼职做了，这样就可以留出更多的时间做其它兼职。

活动安排

```
1 public static List<Integer> selectActivity(int[] a, int[] s, int[] f){
2     List<Integer> activity = new ArrayList<Integer>();
3     int n = s.length;
4     activity.add(a[0]);
5     int k = 0;
```

```
6   for(int i = 1; i < n; i++) {  
7       if(s[i] >= f[k]) {  
8           activity.add(a[i]);  
9           k = i;  
10      }  
11  }  
12  return activity;  
13 }
```

13.3 部分背包

13.3.1 部分背包 (Knapsack)

假设有一个小偷，背着一个容量为 c 的背包（最多可以放的物品重量不能超过 c ）。商店中一共有 n 种物品，每种物品 i 的重量和价值分别为 w_i 和 v_i 。如何选择物品放入背包才能使得总价值最大？

假设背包的容量为 30，各类物品对应的重量和价格如下：

物品	1	2	3	4	5
重量 w	10	5	15	10	20
价值 v	20	30	15	25	10

表 13.2: 物品信息

对于背包问题，很显然它是满足最优子结构性质的，因为一个容量为 c 的背包问题必然包含容量小于 c 的背包问题的最优解。如果要使最终的价值最大，那么必定需要使得选择的单位中物品的价值最大。所以背包问题的贪心策略是优先选择单位重量价值最大的物品，当这个物品选择完之后，继续选择其它价值最大的物品。

部分背包

```
1 class Item:
2     def __init__(self, type, weight, value):
3         """
4             初始化物品
5             Args:
6                 type (int): 物品类型
7                 weight (float): 重量
8                 value (float): 价值
9         """
10        self.type = type
11        self.weight = weight
```



```

12         self.value = value
13         self.unit_value = value / weight    # 单位价值
14
15 def main():
16     capacity = 30        # 背包容量
17     item_type = [1, 2, 3, 4, 5]
18     item_weight = [10, 5, 15, 10, 30]
19     item_value = [20, 30, 15, 25, 10]
20
21     items = []
22     for i in range(len(item_type)):
23         items.append(
24             Item(item_type[i], item_weight[i], item_value[i])
25         )
26
27     # 物品按照单价降序排列
28     items.sort(key=lambda x: x.unit_value, reverse=True)
29
30     # 背包选择
31     selected_items = []
32     cur_weight = 0        # 当前背包重量
33     for item in items:
34         if cur_weight + item.weight <= capacity:
35             cur_weight += item.weight
36             selected_items.append(item)
37         else:
38             item.weight = capacity - cur_weight
39             selected_items.append(item)
40             break
41
42     # 输出选择结果
43     for item in selected_items:
44         print("类型: %d, 重量: %.2f" % (item.type, item.weight))
45
46 if __name__ == "__main__":
47     main()

```

13.4 贪心算法局限性

13.4.1 局限性

贪心算法可以用来解决优化问题，与其他算法相比，贪心算法的最大优点是在大多数情况下易于实现且非常高效，但也存在一些问题。

贪心算法不能保证解是最佳的，因为它总是从局部出发，并没从整体考虑。尽管贪心算法给出了接近最优的解决方案，但它未能产生最优的解决方案。背包问题和旅行商问题是贪心算法无法产生最佳解决方案的问题示例。

当需要实时解决方案且近似答案足够好时，贪心算法最适用。显然，贪心算法可在确保产生最佳解决方案的同时最大程度地减少时间，因此更适用于需要较少时间的情况。

13.4.2 硬币找零

假设有 1 元、5 元、11 元这三种面值的硬币，给定一个找零金额，比如 28 元，最少使用的硬币组合是什么？

使用贪心算法可以得到结果为 [11, 11, 5, 1]，总共是 4 个硬币。对于这个例子而言，4 个硬币的确是最优解。但是假如找零是 15 元呢？使用贪心算法结果为 [11, 1, 1, 1, 1]，共用了 5 个硬币，然而最优解是 [5, 5, 5]，共 3 个硬币。硬币找零问题可以使用动态规划计算出最优解。

13.4.3 0-1 背包

部分背包问题可以用贪心算法实现，是因为选择放入的物品可以进行拆分，即并不需要放入整个物品。与之对应的另一种背包问题为 0-1 背包问题，这个时候整个物品不可以拆分，只可以选择放入或者不放入。0-1 背包问题用贪心算法并不能求得准确的解，需要用动态规划算法求解。

假设 0-1 背包的容量为 8，物品的重量和价值分别为：

物品	1	2	3	4
重量 w	2	3	4	5
价值 v	3	4	5	6
单位价值	1.5	1.33	1.25	1.2

表 13.3: 物品信息

使用贪心算法会优先选择单位重量价值最高的物品，当拿完物品 1 和物品 2 后，背包剩余容量为 $8 - 2 - 3 = 3$ 。由于物品不可拆分，背包无法装入物品 3 和物品 4。最终总价值为 $3 + 4 = 7$ 。

但是使用动态规划可以计算出最佳方案是选择物品 2 和物品 4，最大价值为 $4 + 6 = 10$ 。

Chapter 14 动态规划

14.1 动态规划

14.1.1 动态规划 (Dynamic Programming)

动态规划在数学上属于运筹学的分支，是求解决策过程最优化的数学方法，同时也是计算机科学与技术领域中一种常见的算法思想。

动态规划算法的基本思想与分治法类似，也是将带求解的问题分解为若干个子问题，按顺序求解子问题。前一子问题的解，为后一子问题的求解提供了有用的信息。

在求解任一子问题时，列出各种可能的局部解，通过决策保留那些有可能达到最优的局部解，丢弃其它局部解。依次解决各子问题，最后一个子问题就是初始问题的解。

动态规划的本质是对问题状态的定义和状态转移方程的定义。动态规划通过拆分问题，定义问题状态和状态之间的关系，使得问题能够以递推的方式去解决。因此在一个典型的动态规划问题上，需要定义问题状态以及写出状态转移方程，这样对于问题的解答就会一目了然。

14.1.2 爬楼梯

有一座高度是 10 级台阶的楼梯，从下往上走，每跨一步只能向上 1 级或者 2 级台阶，要求求出一共有多少种走法。

比如，每次走 1 级台阶，一共走 10 步，这是其中一种走法，可以简写成 [1, 1, 1, 1, 1, 1, 1, 1, 1, 1]。再比如，每次走 2 级台阶，一共走 5 步，这是另一种走法，可以简写成 [2, 2, 2, 2, 2]。当然，除此之外，还有很多很多种走法。

暴力枚举的算法利用排列组合的思想，通过多重循环遍历出所有的可能性。但是暴力枚举的时间复杂度是指数级的，有没有更高效的解法呢？

要不找个楼梯走一下试试吧！正好能减肥！

动态规划是一种分阶段求解决策问题的数学思想，它不止用于编程领域，也应用于管理学、经济学、生物学等。总的来说就是大事化小，小事化了。

在爬楼梯问题中，假设你只差最后一步就走到第 10 级台阶，这时候会出现几种情况？

当然是两种喽，因为每一步只许走 1 级或 2 级，所以最后一步要么是从第 9 级走到第 10 级，要么是从第 8 级走到第 10 级。

接下来就引申出了一个新的问题，如果已知从第 0 级走到第 9 级的走法有 X 种，从第 0 级走到第 8 级的走法有 Y 种，那么从第 0 级走到第 10 级的走法就有 $X + Y$ 种。

为了方便表达，我们把 10 级台阶的走法数量简写为 $F(10)$ ，此时 $F(10) = F(9) + F(8)$ 。那么如何计算 $F(9)$ 和 $F(8)$ 呢？

利用刚才的思路可以很容易地推断出 $F(9) = F(8) + F(7)$ ， $F(8) = F(7) + F(6)$ 。这样我们就把一个复杂的问题分阶段进行简化，逐步简化成简单的问题，这就是动态规划的思想。

当只有 1 级台阶和 2 级台阶的时候，显然分别只有 1 种和 2 种走法。由此可以归纳出递推公式：

$$F(n) = \begin{cases} 1 & n = 1 \\ 1 & n = 2 \\ F(n-1) + F(n-2) & n \geq 3 \end{cases}$$

动态规划当中包含了三个重要的概念：

- 最优子结构
- 边界
- 状态转移方程

刚才分析得出的 $F(10) = F(9) + F(8)$ 中, $F(9)$ 和 $F(8)$ 就是 $F(10)$ 的最优子结构。

当只有 1 级台阶或 2 级台阶时，我们可以直接得出结果，无需继续简化。因此 $F(1)$ 和 $F(2)$ 就是问题的边界。如果一个问题没有边界，将永远无法得到有限的结果。

$F(n) = F(n-1) + F(n-2)$ 是阶段与阶段之间的状态转移方程，这是动态规划的核心，决定了问题的每一个阶段与下一阶段的关系。

既然已经归纳出了 $F(n) = F(n-1) + F(n-2)$ ，又知道了边界，那就可以直接用递归的思路实现。

爬楼梯（递归）

```

1 int climbStairs(int n) {
2     if(n <= 0) {
3         return 0;
4     } else if(n == 1) {
5         return 1;
6     } else if(n == 2) {
7         return 2;

```

```
8 |     }
9 |     return climbStairs(n-1) + climbStairs(n-2);
10| }
```

递归的确可以计算出最终答案，可是其时间复杂度却是指数级的。

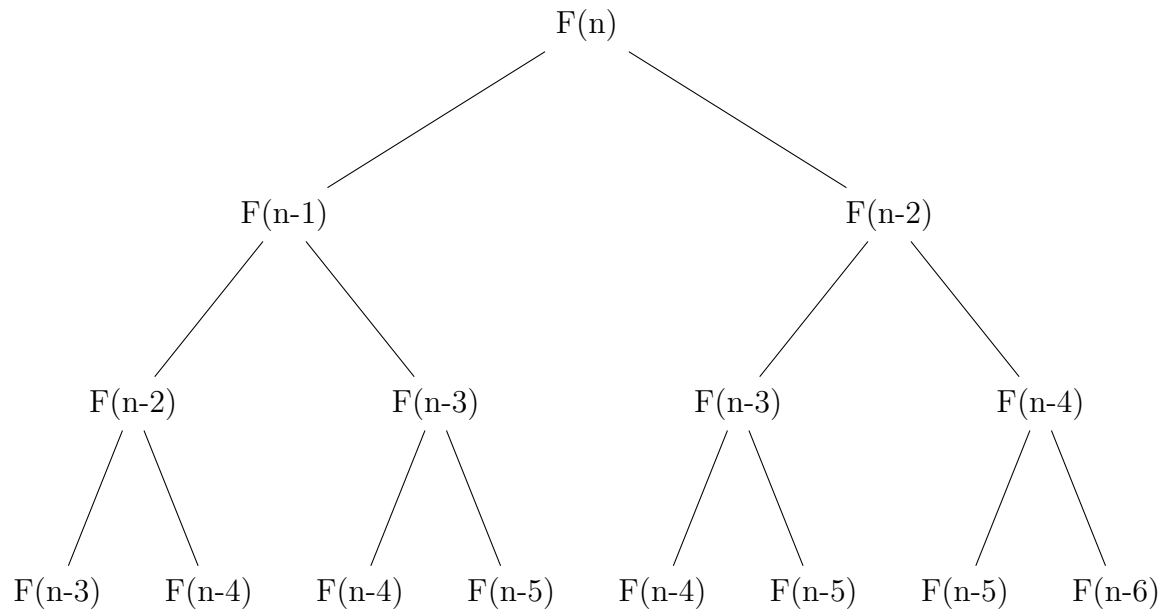


图 14.1: 递归树

二叉树的结点个数就是递归方法所需要计算的次数，因此递归方法的时间复杂度为 $O(2^n)$ 。

可以发现在递归树中，有些相同的参数被重复计算了，越往下走，重复的越多。

可是一一定要对 $F(n)$ 自顶向下做递归运算吗？采用自底向上，用迭代的方法也可以推导出结果。

之前已经得出结论， $F(1) = 1$ ， $F(2) = 2$ ：

台阶数	1	2	3	4	5	6
走法数	1	2				

因为 $F(3) = F(1) + F(2)$ ，经过一次迭代可以计算出 $F(3)$ ：

台阶数	1	2	3	4	5	6
走法数	1	2	3			

第二次迭代，由于 $F(4)$ 只依赖于 $F(2)$ 和 $F(3)$ ，而 $F(3)$ 已经在上一次迭代计算得出，无需重复计算：

台阶数	1	2	3	4	5	6
走法数	1	2	3	5		

爬楼梯（动态规划）

```

1 int climbStairs(int n) {
2     if(n <= 0) {
3         return 0;
4     } else if(n == 1) {
5         return 1;
6     } else if(n == 2) {
7         return 2;
8     }
9     int num1 = 1;
10    int num2 = 2;
11    int sum;
12    for(int i = 3; i <= n; i++) {
13        sum = num1 + num2;
14        num1 = num2;
15        num2 = sum;
16    }
17    return sum;
18 }
```


14.2 硬币找零

14.2.1 硬币找零

有三种硬币，面值分别是 2 元、5 元、7 元，每种硬币都有足够多。买一个物品需要 27 元，如何用最少的硬币组合正好付清？

要让硬币最少，应该尽量用面值大的硬币，也就是贪心算法的结果为 $7 + 7 + 7 + 5 = 26$ ，呃……

那么需要改变一下策略，尽量用面值大的硬币，最后如果可以用一种硬币付清就行。 $7 + 7 + 7 + 2 + 2 + 2 = 27$ ，一共 6 枚硬币，应该对了吧……

但是正确答案是 $7 + 5 + 5 + 5 + 5 = 27$ ，一共 5 枚硬币。

状态在动态规划中的作用属于定海神针，确定状态需要两个意识：最后一步和子问题。

虽然目前不知道最优策略是什么，但是最优策略肯定是 k 枚硬币 a_1, a_2, \dots, a_k 面值加起来是 27，所以一定存在最后一枚硬币 a_k 。除了这枚硬币，前面硬币的面值加起来是 $27 - a_k$ 。因为是最优策略，所以拼出 $27 - a_k$ 的硬币数一定要最少。

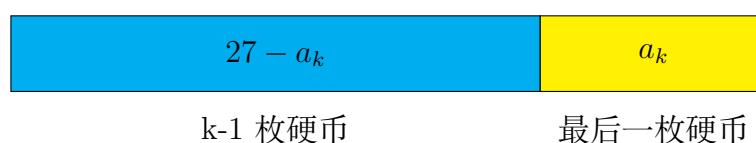


图 14.2: 最后一步

因此问题就变成了最少用多少枚硬币可以拼出 $27 - a_k$ 。这样就将原问题转化成了一个子问题，而且规模更小。

假设状态 $F(x)$ 表示拼出 x 元的所需的最少硬币数，最后那枚硬币 a_k 只可能是 2 元、5 元或 7 元：

- 如果 a_k 是 2 元, $F(27) = F(27-2) + 1$ (加上最后一枚 2 元硬币)。
- 如果 a_k 是 5 元, $F(27) = F(27-5) + 1$ (加上最后一枚 5 元硬币)。
- 如果 a_k 是 7 元, $F(27) = F(27-7) + 1$ (加上最后一枚 7 元硬币)。

由此可得到递归公式:

$$F(27) = \min\{F(27-2) + 1, F(27-5) + 1, F(27-7) + 1\}$$

硬币找零 (递归)

```

1 public static int getMinCoins(int price) {
2     // 0元钱只需要0枚硬币
3     if(price == 0) {
4         return 0;
5     }
6     // 初始化为无穷大
7     int coinNum = Integer.MAX_VALUE - 1;
8     // 最后一枚硬币是2元
9     if(price >= 2) {
10         coinNum = Math.min(getMinCoins(price-2) + 1, coinNum);
11     }
12     // 最后一枚硬币是5元
13     if(price >= 5) {
14         coinNum = Math.min(getMinCoins(price-5) + 1, coinNum);
15     }
16     // 最后一枚硬币是7元
17     if(price >= 7) {
18         coinNum = Math.min(getMinCoins(price-7) + 1, coinNum);
19     }
20     return coinNum;
21 }

```

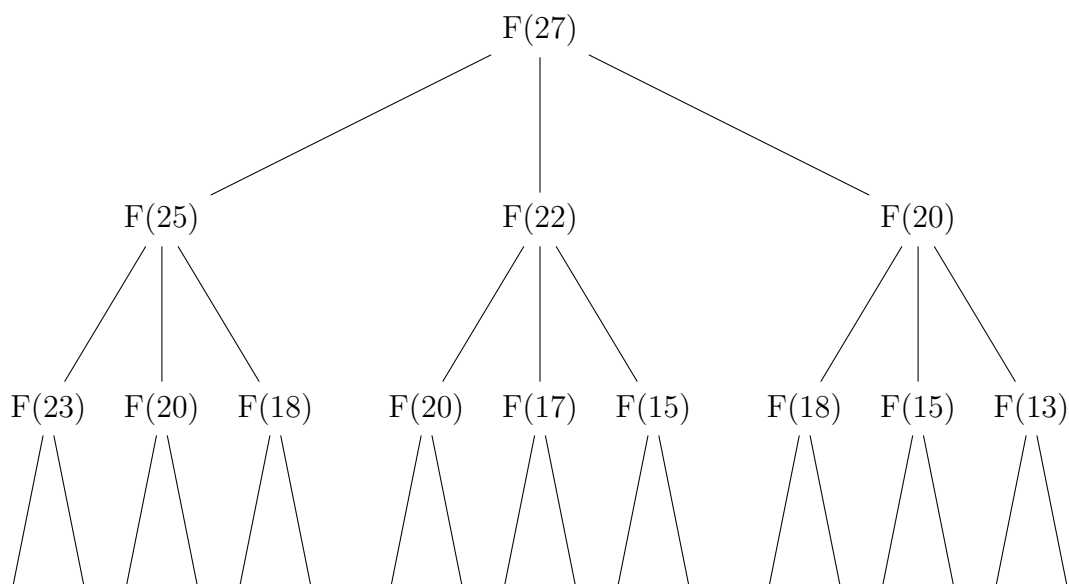


图 14.3: 递归树

动态规划的另一个组成部分就是状态转移方程：

$$F[X] = \min\{F[X - 2] + 1, F[X - 5] + 1, F[X - 7] + 1\}$$

其次需要考虑初始条件和边界情况。如果不能拼出 i 元，则定义 $F[i] = \infty$ 。例如当 $x - 2$ 、 $x - 5$ 、 $x - 7$ 小于 0 时：

$$F[-1] = F[-2] = \dots = \infty$$

状态转移方程的初始状态为 $F[0] = 0$ ，因为无需任何硬币就能拼成 0 元。

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0								

表 14.1: 初始状态

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0	∞							

表 14.2: $F[1] = \min\{F[-1] + 1, F[-4] + 1, F[-6] + 1\} = \infty$

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0	∞	1						

表 14.3: $F[2] = \min\{F[0] + 1, F[-3] + 1, F[-5] + 1\} = 1$

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0	∞	1	∞					

表 14.4: $F[3] = \min\{F[1] + 1, F[-2] + 1, F[-4] + 1\} = \infty$

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0	∞	1	∞	2				

表 14.5: $F[4] = \min\{F[2] + 1, F[-1] + 1, F[-3] + 1\} = 2$

...	F[-1]	F[0]	F[1]	F[2]	F[3]	F[4]	F[5]	F[6]	...	F[27]
∞	∞	0	∞	1	∞	2	1	3	...	5

表 14.6: $F[27] = \min\{F[25] + 1, F[-22] + 1, F[-20] + 1\} = 5$

硬币找零（动态规划）

```

1 def get_min_coins(coins, price):
2     f = [INF] * (price + 1)
3
4     f[0] = 0
5     for i in range(1, price+1):
6         for j in range(len(coins)):
7             if i >= coins[j] and f[i - coins[j]] != INF:
8                 f[i] = min(f[i - coins[j]] + 1, f[i])
9
10    if f[price] == INF:
11        f[price] = -1
12    return f[price]

```

14.3 路径问题

14.3.1 路径问题

有一个机器人位于一个 m 行 n 列的网格的左上角 $(0,0)$ ，机器人每次只能向下或向右移动一步，问有多少种方法可以走到右下角。

	0	1	2	3	4	5	6	7
0	R							
1								
2								
3								

表 14.7: 起点

	0	1	2	3	4	5	6	7
0								
1								
2								↓
3							→	R

表 14.8: 终点

无论机器人用何种方式到达右下角，总有最后挪动的一步。右下角的坐标为 $(m-1, n-1)$ ，那么前一步机器人一定在 $(m-2, n-1)$ 或 $(m-1, n-2)$ 的位置。

如果机器人有 X 种方式从左上角走到 $(m-2, n-1)$ ，有 Y 种方式从左上角走到 $(m-1, n-2)$ ，那么机器人一共有 $X + Y$ 种方式从左上角走到 $(m-1, n-1)$ 。原问题就转换为了机器人有多少种方式从左上角走到 $(m-2, n-1)$ 和 $(m-1, n-2)$ 。

那么可以得出转移方程 $f[i][j] = f[i-1][j] + f[i][j-1]$ ，其中 $f[i][j]$ 表示机器人有多少种方式走到 (i, j) 。

初始条件为 $f[0][0] = 1$ ，因为机器人只有 1 种方式到达左上角。

边界情况为当 $i = 0$ 或 $j = 0$ ，则前一步只能有一个方向到达，因此 $f[i][j] = 1$ 。

	0	1	2	3	4	5	6	7
0		→	→	→	→	→	→	R
1	↓							
2	↓							
3	R							

表 14.9: 边界情况

路径问题

```
1 int uniquePath(int m, int n) {
2     int f[m][n];
3     for(int i = 0; i < m; i++) {
4         for(int j = 0; j < n; j++) {
5             if(i == 0 || j == 0) {
6                 f[i][j] = 1;
7             } else {
8                 f[i][j] = f[i-1][j] + f[i][j-1];
9             }
10        }
11    }
12    return f[m-1][n-1];
13 }
```

14.4 跳跃游戏

14.4.1 跳跃游戏

有 n 块石头分别在 $0, 1, \dots, n-1$ 的位置，一只青蛙在石头 0，想跳到石头 $n-1$ 。如果青蛙在第 i 块石头上，每块石头的元素表示可以跳跃的最长距离。问青蛙能否跳到石头 $n-1$ 。

示例 1

输入: $s = [2, 3, 1, 1, 4]$

输出: True

解释: 可以先跳 1 步，从石头 0 达到石头 1，然后再从石头 1 跳 3 步到达目标。

示例 2

输入: $s = [3, 2, 1, 0, 4]$

输出: False

解释: 无论怎样，总会达到石头 3，但该石头的最大跳跃长度是 0，所以永远不可能到达目标。

如果青蛙能跳到最后一块石头 $n-1$ ，那么它一定是从石头 i 跳过来的 ($i < n-1$)。

这需要两个条件同时满足：

1. 青蛙可以跳到石头 i 。
2. 最后一跳不能超过可以跳跃的最大距离。

那么问题转化为了青蛙能否跳到石头 i 。假设 $f[j]$ 表示青蛙能否跳到石头 j ，可以得出转移方程：

$$f[j] = OR_{0 \leq i < j} (f[i] \text{ AND } i + a[i] \geq j)$$

- $f[j]$: 青蛙能否跳到石头 j 。
- $OR_{0 \leq i < j}$: 枚举上一个跳到的石头 i 。

- $f[i]$: 青蛙能否跳到石头 i 。
- $a[i]$: 最后一步的距离不能超过 a_i 。

初始条件为 $f[0] = \text{True}$ ，因为青蛙一开始就在石头 0。

跳跃游戏

```
1 public static boolean canJump(int[] stone) {
2     int n = stone.length;
3     boolean[] f = new boolean[n];
4     f[0] = true;
5
6     for(int j = 1; j < n; j++) {
7         f[j] = false;
8         for(int i = 0; i < j; i++) {
9             if(f[i] && i + stone[i] >= j) {
10                 f[j] = true;
11                 break;
12             }
13         }
14     }
15     return f[n-1];
16 }
```


14.5 0-1 背包

14.5.1 0-1 背包 (0-1 Knapsack)

有一个小偷带了一个能够装 $C = 20$ 公斤物品的背包到商店里面偷东西，请问他要怎么偷才能使价值最高？

物品	重量 W	价格 V
0	2	3
1	3	4
2	4	5
3	5	8
4	9	10

表 14.10: 物品信息

假设用 $B(k, C)$ 表示当背包容量还剩下 C 的时候，在前 k 件物品中能偷到的最大价值。

$$B(k, C) = \begin{cases} B(k-1, C) & \text{当第 } k \text{ 件太重} \\ \max \begin{cases} B(k-1, C - w_k) + v_k & \text{偷} \\ B(k-1, C) & \text{不偷} \end{cases} & \end{cases}$$

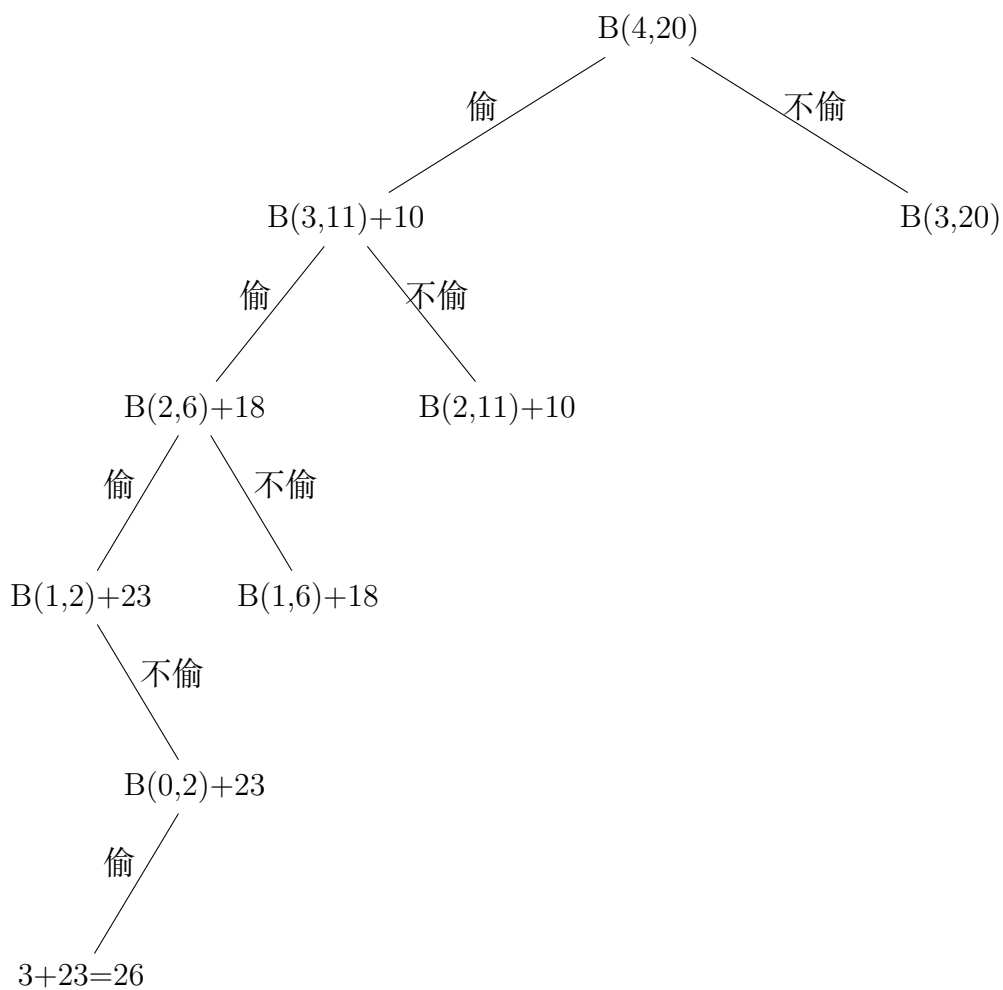


图 14.4: 决策树

边界情况为 $B(i, 0) = B(0, j) = 0$ ，即当背包容量为 0 或者没有物品可偷的情况下，最大价值为 0。

Capacity	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
No Item	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Item 0	0	0	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
Item 1	0	0	3	4	4	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
Item 2	0	0	3	4	5	7	8	9	9	12	12	12	12	12	12	12	12	12	12	12	12
Item 3	0	0	3	4	5	8	8	11	12	13	15	16	17	17	20	20	20	20	20	20	20
Item 4	0	0	3	4	5	8	8	11	12	13	15	16	17	17	20	20	21	22	23	25	26

表 14.11: 动态规划表格法

```

1  #define ITEM_NUM 5
2  #define CAPACITY 20
3
4  int getMaxValue(int *weight, int *value) {
5      int b[ITEM_NUM+1][CAPACITY+1] = {{0}};
6
7      for(int k = 1; k <= ITEM_NUM; k++) {
8          for(int c = 1; c <= CAPACITY; c++) {
9              if(weight[k] > c) {
10                 b[k][c] = b[k-1][c];
11             } else {
12                 b[k][c] = max(
13                     b[k-1][c-weight[k]] + value[k],
14                     b[k-1][c]
15                 );
16             }
17         }
18     }
19     return b[ITEM_NUM][CAPACITY];
20 }

```

Chapter 15 模式匹配

15.1 BF

15.1.1 BF (Brute Force)

在一个字符串中查找另一个字符串的操作称为模式匹配。其中，被查找的字符串被称为文本串，需要查找的子串称为模式串。

查找子串的最简单的算法就是采用暴力匹配的方式。暴力匹配的基本思想就是逐个比较相应位置的字符。

假设文本串为 S，模式串为 P：

1. 如果当前字符匹配成功（即 $S[i] == P[j]$ ），则 $i++$ ， $j++$ ，继续匹配下一个字符。
2. 如果失配（即 $S[i] != P[j]$ ），令 $i = i - j + 1$ ， $j = 0$ ，相当于每次匹配失败时， i 回溯， j 被置为 0。

例如当文本串 $S = \text{"BBC ABCDAB ABCDABCDABDE"}$ ， $P = \text{"ABCDABD"}$ 。

$S[0]$ 为 B， $P[0]$ 为 A，不匹配（情形 2）。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
A	B	C	D	A	B	D																

然后再判断 $S[1]$ 和 $P[0]$ 是否匹配，相当于模式串向右移动一位。 $S[1]$ 与 $P[0]$ 还是不匹配（情形 2）。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
	A	B	C	D	A	B	D															

然后再判断 S[2] 和 P[0] 是否匹配，模式串再向右移动一位。直到 S[4] 与 P[0] 匹配成功（情形 1）。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
				A	B	C	D	A	B	D												

此时 $i = 5$, $j = 1$ ，接下来判断 S[5] 与 P[1] 是否匹配。S[5] 与 P[1] 匹配成功（情形 1）。可得 $i = 6$, $j = 2$ 。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
				A	B	C	D	A	B	D												

直到 S[10] 为空格，P[6] 为 D，不匹配（情形 2）。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
				A	B	C	D	A	B	D												

此时 $i = 5$, $j = 0$ ，相当于判断 S[5] 跟 P[0] 是否匹配。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	B	C		A	B	C	D	A	B		A	B	C	D	A	B	C	D	A	B	D	E
					A	B	C	D	A	B	D											

至此可以看出，按照暴力匹配算法的思路。同理，直到找到匹配的字符串或文本串遍历结束退出。

BF 暴力匹配在最坏情况下时间复杂度是 $O(mn)$ 。

BF

```
1 def brute_force(s, p):
```

```
2     s_len = len(s)
3     p_len = len(p)
4     i = 0
5     j = 0
6
7     while i < s_len and j < p_len:
8         if s[i] == p[j]:
9             i += 1
10            j += 1
11        else:
12            i = i - j + 1
13            j = 0
14
15    if j == p_len:
16        return i - j
17    else:
18        return -1
```

15.2 Sunday

15.2.1 Sunday

Sunday 是从前往后匹配的算法，在匹配失败时重点关注的是文本串中参加匹配的最末位字符的下一位字符。如果该字符没有在模式串中出现则直接跳过，移动位数为模式串长度 + 1，否则移动位数为模式串长度 - 该字符最右出现的下标。

Sunday 算法巧妙的地方在于它发现匹配失败之后可以直接考察文本串中参加匹配的最末尾字符的下一个字符。

例如当文本串 $S = \text{"bcaitsnaxzfnihao"}$ ， $P = \text{"nihao"}$ 。

$S[0]$ 为 b， $P[0]$ 为 n，匹配失败。关注文本串中参加匹配的最末位字符的下一位字符 s，该字符并没在模式串中出现，因此将模式串向右移动模式串长度 + 1，即 $5 + 1 = 6$ 。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
b	c	a	i	t	s	n	a	x	z	f	i	n	i	h	a	o
n	i	h	a	o												

$S[7]$ 为 a， $P[1]$ 为 i，匹配失败。关注文本串中参加匹配的最末位字符的下一位字符 i，该字符在模式串中最右出现出现下标为 1，因此将模式串向右移动模式串长度 - 最右下标，即 $5 - 1 = 4$ 。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
b	c	a	i	t	s	n	a	x	z	f	i	n	i	h	a	o
						n	i	h	a	o						

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
b	c	a	i	t	s	n	a	x	z	f	i	n	i	h	a	o
										n	i	h	a	o		

S[10] 为 f, P[0] 为 n, 匹配失败。关注文本串中参加匹配的最末位字符的下一位字符 a, 该字符在模式串中最右出现下标为 3, 因此将模式串向右移动模式串长度 - 最右下标, 即 $5 - 3 = 2$ 。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
b	c	a	i	t	s	n	a	x	z	f	i	n	i	h	a	o
										n	i	h	a	o		

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
b	c	a	i	t	s	n	a	x	z	f	i	n	i	h	a	o
												n	i	h	a	o

此时, 模式串匹配成功。

Sunday

```

1 def sunday(s, p):
2     s_len = len(s)
3     p_len = len(p)
4     i = 0
5     j = 0
6     result = 0
7
8     while i < s_len and j < p_len:
9         if s[i] == p[j]:
10             i += 1
11             j += 1
12             continue
13
14         idx = result + p_len
15         if idx >= s_len:
16             return -1
17
18         k = p_len - 1
19         while k >= 0 and s[idx] != p[k]:

```



```
20         k -= 1
21
22     i = result
23     i += p_len - k
24     result = i
25     j = 0
26
27     if result + p_len > s_len:
28         return -1
29
30     return result
```

15.3 RK

15.3.1 RK (Rabin-Karp)

RK 算法的命名由 Rabin 和 Karp 两位发明者的名字而来，它的实现方式有点与众不同。BF 算法只是简单粗暴地对两个字符串的所有字符依次比较，而 RK 算法比较的是两个字符串的哈希值。

RK 算法的基本思想就是将模式串 P 的哈希值跟文本串 S 中每一个长度为 $|S|$ 的子串的哈希值进行比较。如果两个字符串哈希值不相同，则它们肯定不匹配。如果它们的哈希值相同，它们有可能匹配（因为可能存在哈希冲突）。

由于哈希函数有可能会产生哈希冲突，哈希值相等的两个字符串不一定相同。因此如果两个字符串的哈希值相等，就把这两个字符串本身进行一次比较即可。这种方法的前提是要控制冲突概率，达到可以接受的状态。

不过每次 hash 的时间复杂度为 $O(n)$ ，如果把全部子串都进行 hash，总的时间复杂度不是和 BF 算法一样，都是 $O(mn)$ 了吗？

其实对子串的 hash 计算并不是独立的，从第二个子串开始，每一次子串的 hash 都可以由上一次子串进行简单的加减得到（减去第一个字符的 hash，加上新字符的 hash）。因此通过设计特殊的哈希算法，只需扫描一遍文本串就能计算出所有子串的哈希值。期间最多需要比较 $m - n + 1$ 个子串的哈希值。

相比于 BF 算法，RK 算法采用哈希值比较的方式，免去了许多无谓的字符比较，所以时间复杂度大大提高了。RK 算法的缺点在于哈希冲突，每一次哈希冲突的时候，RK 算法都要对子串和模式串进行逐个字符的比较。如果冲突太多了，RK 算法就退化成了 BF 算法。

RK

```
1 public static int rk(String s, String p) {
```

```

2     int sLen = s.length();
3     int pLen = p.length();
4
5     int sHash = hash(s.substring(0, pLen)); // 文本串子串哈希值
6     int pHash = hash(p);                  // 模式串哈希值
7
8     for(int i = 0; i < sLen - pLen + 1; i++) {
9         if(sHash == pHash) {
10             if(match(s.substring(i, i + pLen), p)) {
11                 return i;
12             }
13         }
14         if(i < sLen - pLen) {
15             sHash = nextHash(s, sHash, i, pLen);
16         }
17     }
18
19     return -1;
20 }
21
22 public static int hash(String s) {
23     int hashCode = 0;
24     for(int i = 0; i < s.length(); i++) {
25         hashCode += s.charAt(i) - 'a';
26     }
27     return hashCode;
28 }
29
30 public static int nextHash(String s, int hash, int start, int n) {
31     hash -= s.charAt(start) - 'a';
32     hash += s.charAt(start + n) - 'a';
33     return hash;
34 }
35
36 public static boolean match(String s, String p) {
37     return s.equals(p);
38 }

```

15.4 BM

15.4.1 BM (Boyer-Moore)

BM 算法的名字取自于它的两位发明者，计算机科学家 Bob Boyer 和 J Strother Moore。为了能减少比较，BM 算法制定了两条规则：

- 坏字符规则 (bad character)
- 好后缀规则 (good suffix)

15.4.2 坏字符规则

坏字符是指文本串与模式串不匹配的字符。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
G	T	A	G	C	G	G	C	G												

咦？为什么坏字符不是主串中下标为 2 的字符 T 呢？那个位置不是先被检测到的吗？

BM 算法的检测顺序是从字符串的最右侧向最左侧进行的。当检测到第一个坏字符后，并没有必要让模式串一位一位向后挪动并比较。因为只有模式串与坏字符对齐的位置相同的情况下，两者才有匹配的可能。由于模式串的第 1 位字符也是 T，这样就可以直接把模式串中的 T 与文本串的坏字符对齐，进行下一轮比较。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
							G	T	A	G	C	G	G	C	G					

坏字符的位置越靠右，下一轮模式串的挪动跨度就可能越长，节省的比较次数也就越多。这就是 BM 算法从右向左检测的好处。

接着，从右向左成功匹配 GCG，并遇到坏字符 A。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
							G	T	A	G	C	G	G	C	G					

按照类似的方式，找到模式串的第 2 位字符 A，将它与文本串的坏字符对齐。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
										G	T	A	G	C	G	G	C	G		

如果出现坏字符在模式串中不存在的情况，就直接把模式串挪到主串坏字符的下一位。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
G	C	A	I	C	G	G	C	G												

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	T	A	T	A	G	C	T	G	G	T	A	G	C	G	G	C	G	A	A
									G	C	A	I	C	G	G	C	G			

15.4.3 好后缀规则

好后缀是指文本串与模式串当中相匹配的后缀。

例如对于这个例子，如果继续使用坏字符规则，模式串只能向后挪动一位。

0	1	2	3	4	5	6	7	8	9	10	11	12	13
C	T	G	G	G	C	G	A	G	C	G	G	A	A
G	C	G	A	G	C	G							

为了能真正减少比较次数，就需要使用好后缀规则。在第一轮比较中，文本串和模式串都有共同的后缀 GCG，这就是所谓的好后缀。如果模式串的其它位置也包含与 GCG 相同的子串，那么就可以挪动模式串，让这个子串与好后缀对齐。

0	1	2	3	4	5	6	7	8	9	10	11	12	13
C	T	G	G	G	C	G	A	G	C	G	G	A	A
G	C	G	A	G	C	G							

0	1	2	3	4	5	6	7	8	9	10	11	12	13
C	T	G	G	G	C	G	A	G	C	G	G	A	A
				G	C	G	A	G	C	G			

如果模式串中不存在其它与好后缀相同的片段，是不是可以直接把模式串挪到好后缀之后？

使不得！这里还要判断一种特殊情况，模式串的前缀是否和好后缀的后缀相匹配，免得挪过头了。

0	1	2	3	4	5	6	7	8	9	10	11	12
T	G	G	G	C	G	A	G	C	G	G	A	A
C	G	A	G	C	G							

0	1	2	3	4	5	6	7	8	9	10	11	12
T	G	G	G	C	G	A	G	C	G	G	A	A
				C	G	A	G	C	G			

那应该什么时候用坏字符规则，什么时候用好后缀规则呢？

可以在每一轮字符比较之后，按照坏字符和好后缀规则分别计算相应的挪动距离，哪一种距离更长，就把模式串挪动对应的长度。比如坏字符可以让模式串在下一轮挪动 3 位，好后缀可以让模式串移动 5 位，那么就应该采用好后缀规则。

15.5 KMP

15.5.1 KMP (Knuth-Morris-Pratt)

KMP 算法是一个里程碑似的算法，它的出现宣告了人类找到了线性时间复杂度的字符串匹配算法。在此之后才出现了其它线性时间的字符串匹配算法，比如 BM 算法和 Sunday 算法。

KMP 算法由三位计算机科学家 D. E. Knuth、J. H. Morris 和 V. R. Pratt 提出，KMP 这个算法名字正是取自这三个人的姓氏首字母。

与 BM 算法类似，KMP 算法也在试图减少无谓的字符比较，但 KMP 算法把专注点放在了已匹配的前缀。

在每一次匹配过程中，其实可以判断出后续几次匹配是否会成功。算法的核心就是每次匹配过程中推断出后续完全不可能匹配成功的部分，从而减少比较的次数。

例如主串中已匹配末尾的 GTG 是最长可匹配后缀，模式串中开头的 GTG 是最长可匹配前缀。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	G	T	G	A	G	C	T	G	G	T	G	T	G	T	G	C	F	A	A
G	T	G	T	G	C	F														

将最长可匹配后缀与最长可匹配前缀对齐。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	G	T	G	A	G	C	T	G	G	T	G	T	G	T	G	C	F	A	A
		G	T	G	T	G	C	F												

主串中已匹配末尾的 G 是最长可匹配后缀，模式串中开头的 G 是最长可匹配前缀。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	G	T	G	A	G	C	T	G	G	T	G	T	G	T	G	C	F	A	A
		G	T	G	T	G	C	F												

将最长可匹配后缀与最长可匹配前缀对齐。

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
G	T	G	T	G	A	G	C	T	G	G	T	G	T	G	T	G	C	F	A	A
				G	T	G	T	G	C	F										

KMP 算法的整体思路就是在已匹配的前缀当中寻找最长可匹配后缀子串和最长可匹配前缀子串，在下一轮直接把两者对齐，从而实现模式串的快速移动。

那么如何找到一个字符串前缀的最长可匹配后缀子串和最长可匹配前缀子串呢？难道在每一轮都要重新遍历吗？

要找到这两个子串没有必要每次都去遍历，可以事先缓存到一个集合当中，用的时候再去集合里面取。这个集合被称为 next 数组，如何生成 next 数组是 KMP 算法的最大难点。

15.5.2 next 数组

next 数组实质上就是找出模式串前后字符重复出现的个数，next[i] 表示模式串 T[0] 到 T[i] 这个子串使得前 k 个字符等于后 k 个字符的最大值，其中 k 不能取 i + 1，因为子串一共才 i + 1 个字符，自己跟自己相等毫无意义。

例如当文本串 S = "ABABAABAABAC", P = "ABAABAC"。

模式串	A	B	A	A	B	A	C
next 数组	0	0	1	1	2	3	0

表 15.1: 模式串 ABAABAC 的 next 数组

0	1	2	3	4	5	6	7	8	9	10	11
A	B	A	B	A	A	B	A	A	B	A	C
A	B	A	A	B	A	C					

已经成功匹配 3 个，查看 $\text{next}[3-1] = 1$ ，表示最长匹配前后缀长度为 1。模式串需要跳过的长度为匹配上的长度减去最长匹配前后缀长度，即 $3 - 1 = 2$ 。

0	1	2	3	4	5	6	7	8	9	10	11
A	B	A	B	A	A	B	A	A	B	A	C
		A	B	A	A	B	A	C			

已经成功匹配 6 个，查看 $\text{next}[6-1] = 3$ ，表示最长匹配前后缀长度为 3。模式串需要跳过的长度为匹配上的长度减去最长匹配前后缀长度，即 $6 - 3 = 3$ 。

0	1	2	3	4	5	6	7	8	9	10	11
A	B	A	B	A	A	B	A	A	B	A	C
					A	B	A	A	B	A	C

在跳过不可能匹配的趟数后，并非再从头开始匹配，而是从之前不匹配的位置开始。

15.5.3 算法分析

假设模式串的长度为 m ，主串的长度为 n 。KMP 算法唯一的额外空间是 next 数组，那么空间复杂度就是 $O(m)$ 。

至于时间复杂度，KMP 算法包括两步。第一步生成 next 数组，时间复杂度为 $O(m)$ ；第二步是对主串的遍历，时间复杂度为 $O(n)$ 。因此 KMP 算法的整体时间复杂度就是 $O(m + n)$ 。

KMP

```
1 public static int[] getNexts(String p) {
```

```

2     int n = p.length();
3     int[] next = new int[n];
4     int j = 0;
5     for(int i = 2; i < n; i++) {
6         while(j != 0 && p.charAt(j) != p.charAt(i-1)) {
7             j = next[j];
8         }
9         if(p.charAt(j) == p.charAt(i-1)) {
10             j++;
11         }
12         next[i] = j;
13     }
14     return next;
15 }
16
17 public static int kmp(String s, String p) {
18     int[] next = getNexts(p);
19     int j = 0;
20     for(int i = 0; i < s.length(); i++) {
21         while(j > 0 && s.charAt(i) != p.charAt(j)) {
22             j = next[j];
23         }
24         if(s.charAt(i) == p.charAt(j)) {
25             j++;
26         }
27         if(j == p.length()) {
28             return i - p.length() + 1;
29         }
30     }
31     return -1;
32 }

```

Chapter 16 树

16.1 AVL 树

16.1.1 AVL 树

AVL 树的命名是取自两位发明者的首字母 G. M. Adelson-Velsky 和 E. M. Landis, AVL 树也称为平衡二叉树。AVL 树能够调整自身的平衡性, AVL 树遵循高度平衡, 任何结点的两个子树高度差不会超过 1。

对于 AVL 树的每一个结点, 平衡因子 (balance factor) 是它左子树高度和右子树高度的差值。只有当二叉树所有结点的平衡因子都是-1、0、1 这三个值时, 这棵树才是一个合格的 AVL 树。

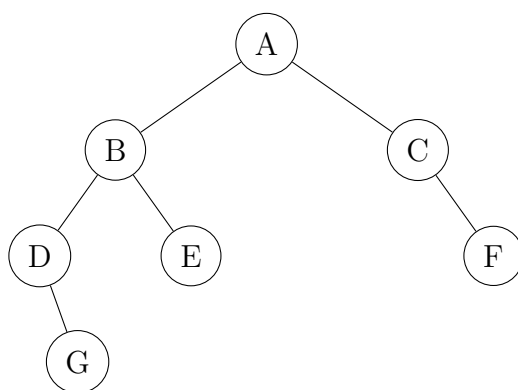


图 16.1: AVL 树

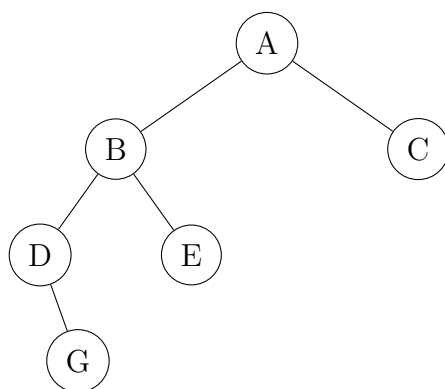


图 16.2: 非 AVL 树

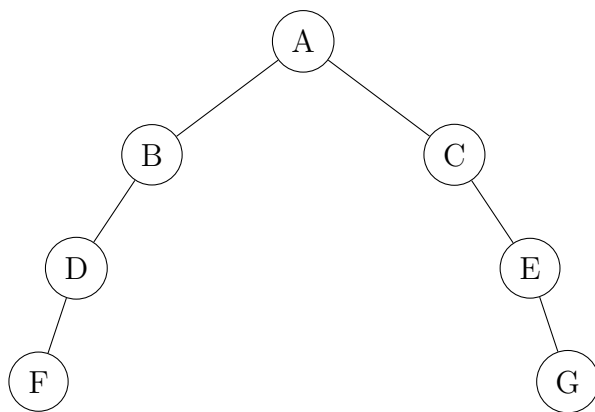


图 16.3: 非 AVL 树

16.1.2 失衡调整

当 AVL 树插入或删除结点时，平衡有可能被打破。

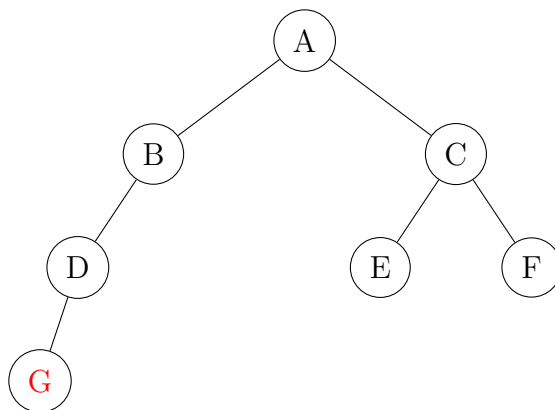


图 16.4: 失衡

通过对 AVL 树进行左旋转、右旋转的操作，就能使其重新恢复平衡。

左旋转

逆时针旋转 AVL 树的两个结点 X 和 Y，使得父结点被自己的右孩子取代，而自己成为左孩子。

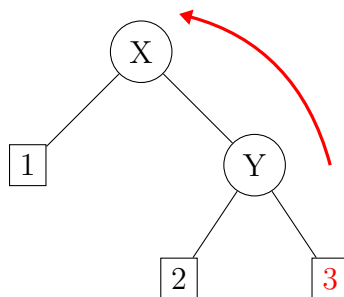


图 16.5: 左旋转 (前)

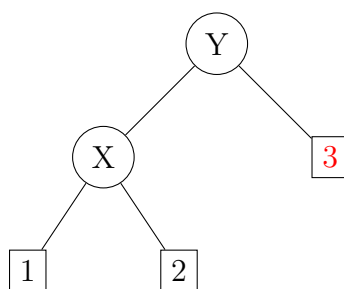


图 16.6: 左旋转 (后)

右旋转

顺时针旋转 AVL 树的两个结点 X 和 Y，使得父结点被自己的左孩子取代，而自己成为右孩子。

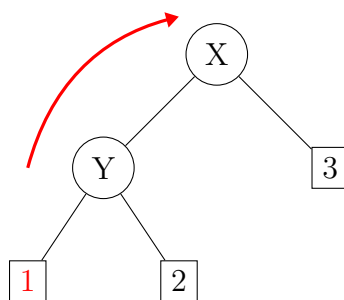


图 16.7: 右旋转 (前)

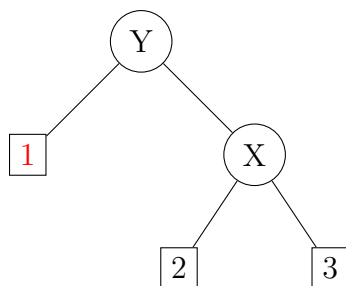


图 16.8: 右旋转 (后)

AVL 树的失衡调整可以分为四种情况：

左左局面 (LL)：右旋转

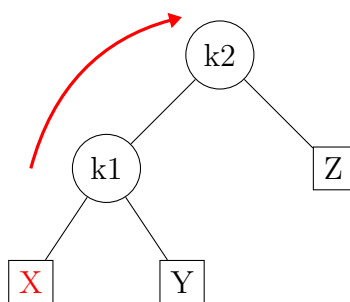


图 16.9: 左左局面

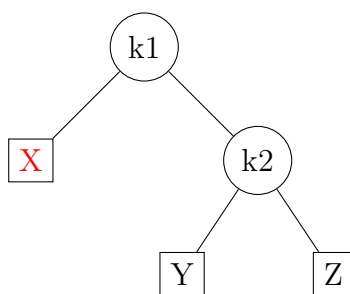


图 16.10: 右旋转

LL 旋转

```

1 static AVLNode* LLRotation(AVLTree *k2) {
2     AVLTree *k1 = k2->left;
3     k2->left = k1->right;
4     k1->right = k2;
  
```

```

5
6     k2->height = max(height(k2->left), height(k2->right)) + 1;
7     k1->height = max(height(k1->left), k2->height) + 1;
8     return k1;
9 }

```

右右局面 (RR): 左旋转

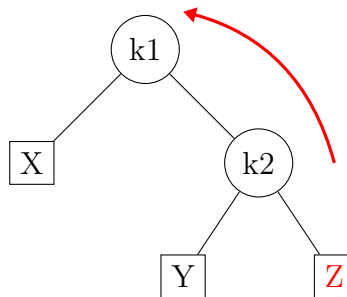


图 16.11: 右右局面

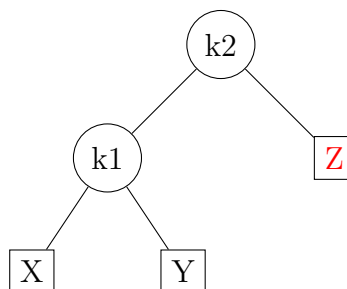


图 16.12: 左旋转

RR 旋转

```

1 static AVLNode* RRRotation(AVLTree *k1) {
2     AVLTree *k2 = k1->right;
3     k1->right = k2->left;
4     k2->left = k1;
5
6     k1->height = max(height(k1->left), height(k1->right)) + 1;
7     k2->height = max(k1->height, height(k2->right)) + 1;
8     return k2;

```

左右局面 (LR)：先左旋转、再右旋转

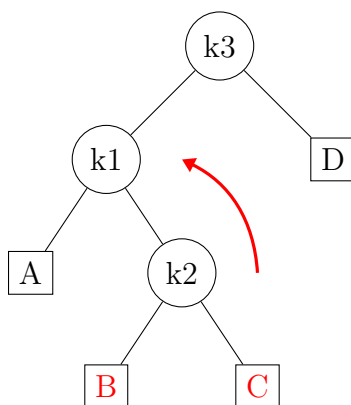


图 16.13: 左右局面

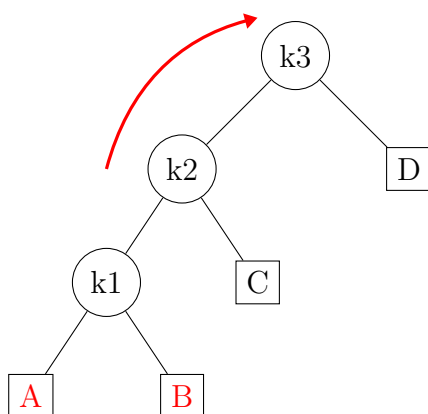


图 16.14: 左旋转

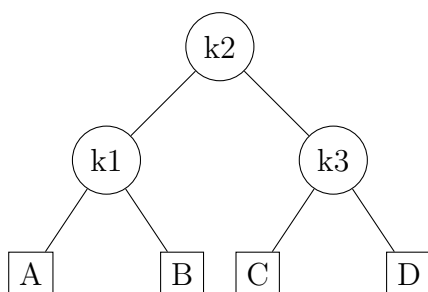


图 16.15: 右旋转

LR 旋转


```

1 static AVLNode* LRRotation(AVLTree *k3) {
2     k3->left = RRRotation(k3->left);
3     return LLRotation(k3);
4 }

```

右左局面 (RL): 先右旋转、再左旋转

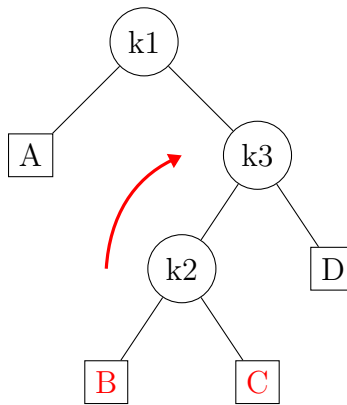


图 16.16: 右左局面

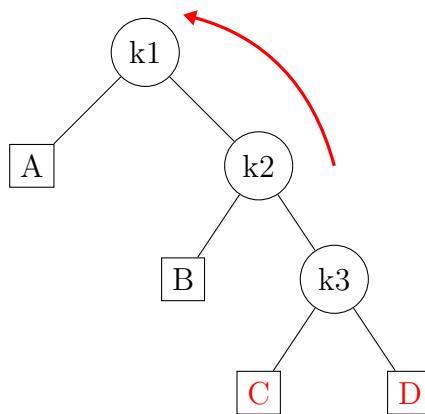


图 16.17: 右旋转

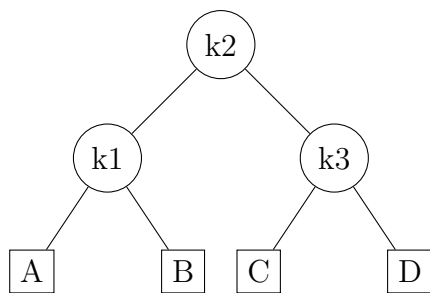


图 16.18: 左旋转

RL 旋转

```

1 static AVLNode* RLRotation(AVLTree *k1) {
2     k1->right = LLRotation(k1->right);
3     return RRRotation(k1);
4 }
  
```

16.1.3 插入/删除结点

例如依次向 AVL 树添加结点 3, 2, 1, 4, 5, 6, 7, 16, 15, 14。

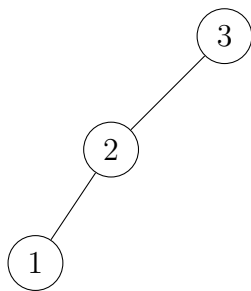


图 16.19: 左左局面：插入 3、2、1

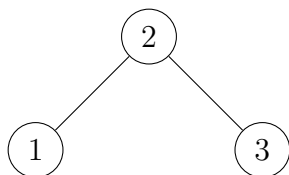


图 16.20: 右旋转

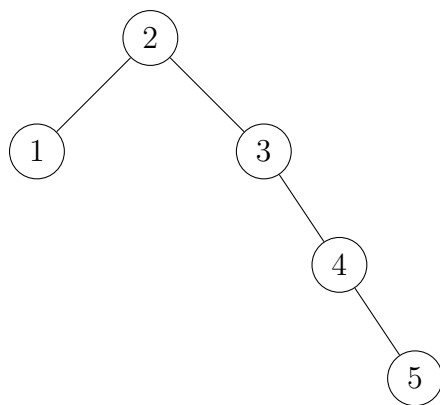


图 16.21: 右右局面: 插入 4、5

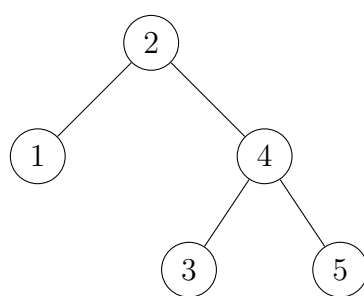


图 16.22: 左旋转

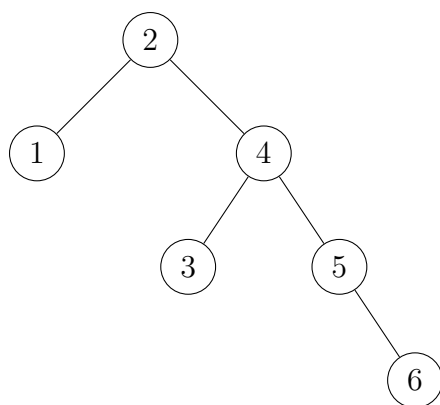


图 16.23: 右右局面: 插入 6

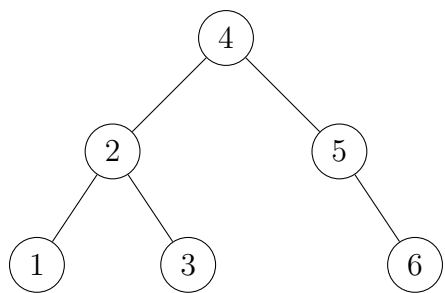


图 16.24: 左旋转

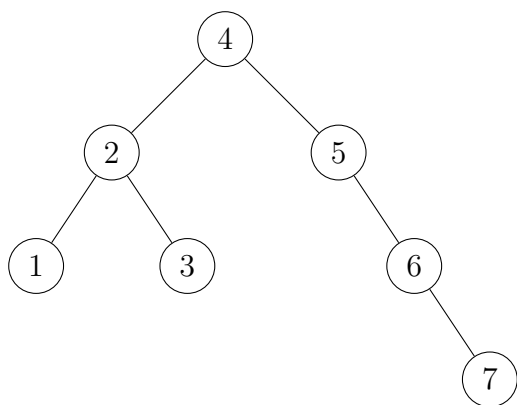


图 16.25: 右右局面: 插入 7

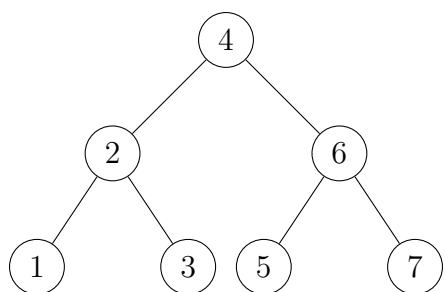


图 16.26: 左旋转

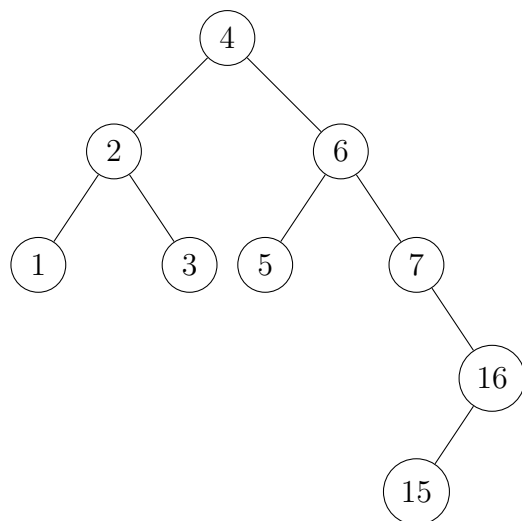


图 16.27: 右左局面: 插入 16、15

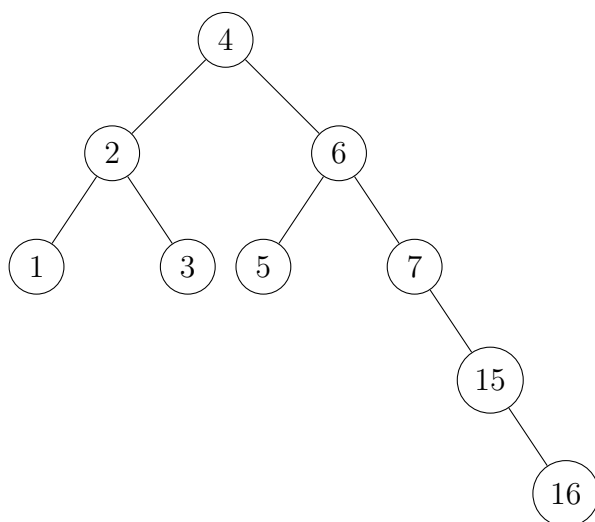


图 16.28: 右旋转

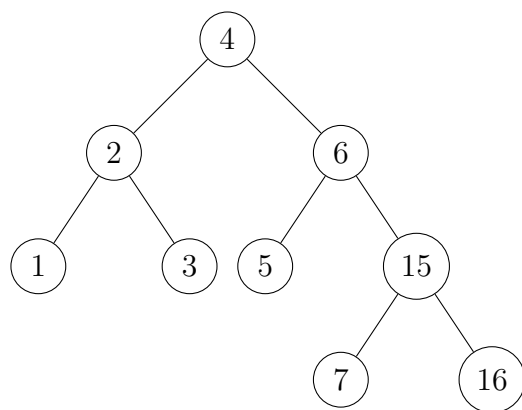


图 16.29: 左旋转

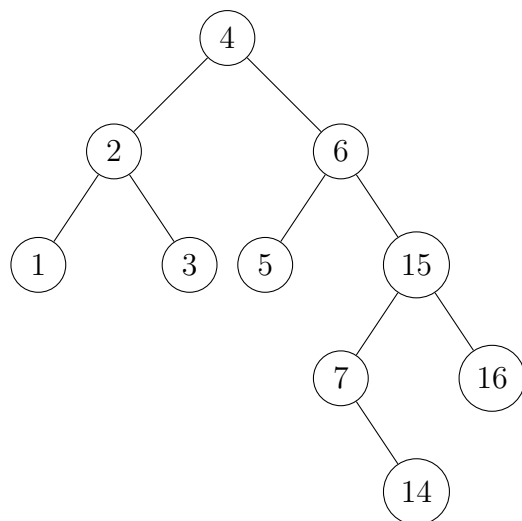


图 16.30: 右左局面: 插入 14

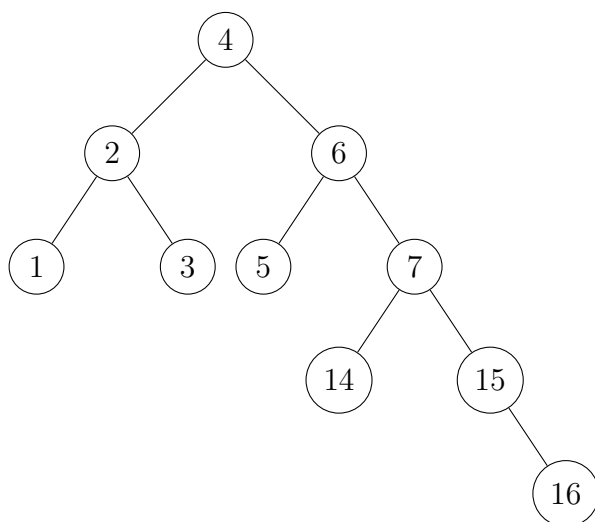


图 16.31: 右旋转

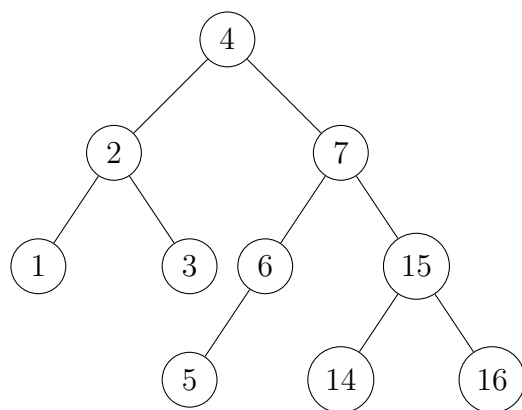


图 16.32: 左旋转

插入结点

```
1 AVLNode* insert(AVLTree *tree, dataType val) {
2     if(!tree) {
3         tree = createNode(val, NULL, NULL);
4         tree->height = 1;
5         return tree;
6     }
7
8     if(val < tree->data) {
9         tree->left = insert(tree->left, val);
10        if(height(tree->left) - height(tree->right) == 2) {
11            if(val < tree->left->data) {
12                tree = LLRotation(tree);
13            } else {
14                tree = LRRotation(tree);
15            }
16        }
17    } else {
18        tree->right = insert(tree->right, val);
19        if(height(tree->right) - height(tree->left) == 2) {
20            if(val > tree->right->data) {
21                tree = RRRotation(tree);
22            } else {
23                tree = RLRotation(tree);
24            }
25        }
26    }
27
28    tree->height = max(height(tree->left), height(tree->right)) + 1;
29    return tree;
30 }
```

删除结点

```

1 static AVLNode* deleteNode(AVLTree *tree, AVLNode *del) {
2     if(!tree || !del) {
3         return NULL;
4     }
5
6     if(del->data < tree->data) {
7         tree->left = deleteNode(tree->left, del);
8         if(height(tree->right) - height(tree->left) == 2) {
9             AVLNode *rightNode = tree->right;
10            if(height(rightNode->left) > height(rightNode->right)) {
11                tree = RLRotation(tree);
12            } else {
13                tree = RRRotation(tree);
14            }
15        }
16    } else if(del->data > tree->data) {
17        tree->right = deleteNode(tree->right, del);
18        if(height(tree->left) - height(tree->right) == 2) {
19            AVLNode *leftNode = tree->left;
20            if(height(leftNode->right) > height(leftNode->left)) {
21                tree = LRRotation(tree);
22            } else {
23                tree = LLRotation(tree);
24            }
25        }
26    } else {
27        if(tree->left && tree->right) {
28            if(height(tree->left) > height(tree->right)) {
29                // 如果左子树比右子树高：
30                // 1. 找出左子树的最大结点
31                // 2. 将最大结点的值赋给tree
32                // 3. 删除最大结点
33                AVLNode *max = getMax(tree->left);
34                tree->data = max->data;
35                tree->left = deleteNode(tree->left, max);
36            } else {
37                // 如果右子树比左子树高（或相等）：

```



```

38         // 1. 找出右子树的最小结点
39         // 2. 将最小结点的值赋给tree
40         // 3. 删除最小结点
41         AVLNode *min = getMin(tree->right);
42         tree->data = min->data;
43         tree->right = deleteNode(tree->right, min);
44     }
45     } else {
46         AVLNode *temp = tree;
47         tree = tree->left ? tree->left : tree->right;
48         free(temp);
49     }
50 }
51 return tree;
52 }

```

16.2 红黑树

16.2.1 红黑树 (Red Black Tree)

红黑树是一种自平衡的二叉查找树，除了符合二叉查找树的基本特性外，它还具有如下附加特性：

1. 结点是红色或黑色的。
2. 根结点是黑色的。
3. 叶子结点都是黑色的空结点 NIL。
4. 红色结点的两个子结点都是黑色的，即从叶子到根的所有路径上不能有连续的两个红色结点。
5. 从任一结点到其每个叶子的所有路径都包含相同数目的黑色结点。

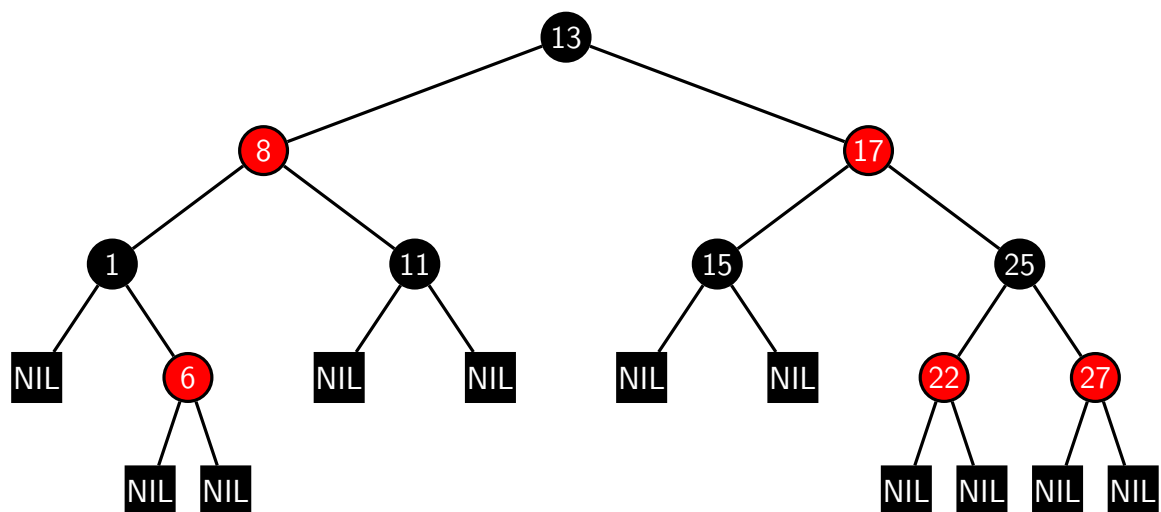


图 16.33: 红黑树

天呐，这条条框框的太多了吧！

正是因为这些规则限制，才保证了红黑树的自平衡，红黑树从根到叶子的最长路径不会超过最短路径的 2 倍。

红黑树的应用有很多，其中 JDK 的集合类 TreeMap 和 TreeSet 底层就是红黑树实现的。在 Java8 中，连 HashMap 也用到了红黑树。

16.2.2 失衡调整

当插入或删除结点时，红黑树的规则可能被破坏，需要调整使其重新符合规则。

例如向红黑树中插入新结点 14，由于父结点 15 是黑色结点，这种情况不会破坏红黑树的规则，无需做任何调整。

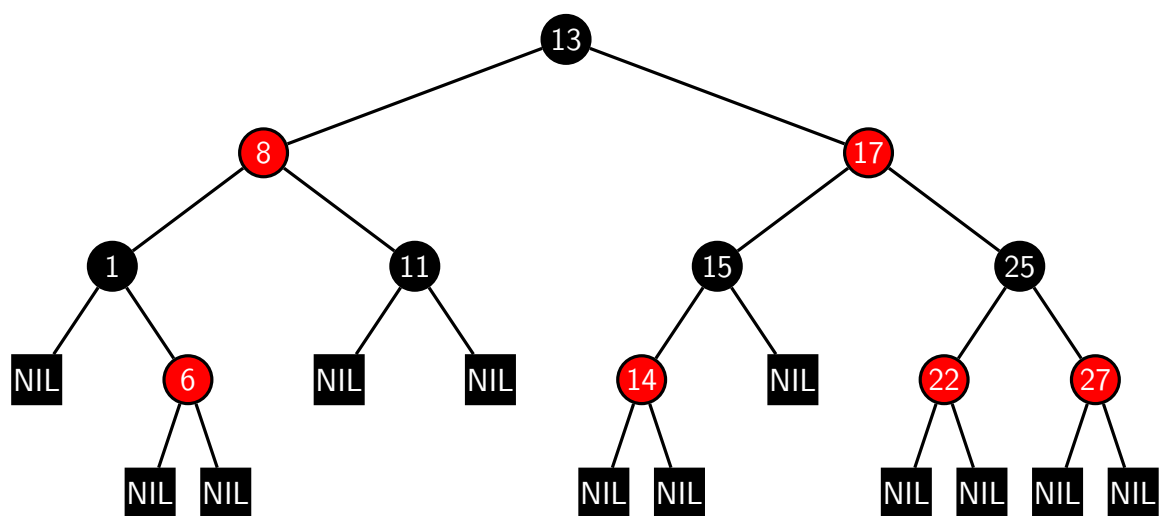


图 16.34: 插入 14

向红黑树中插入新结点 21，由于父结点 22 是红色结点，违反了红黑树的规则 4 (红色结点的两个子结点都是黑色的)。

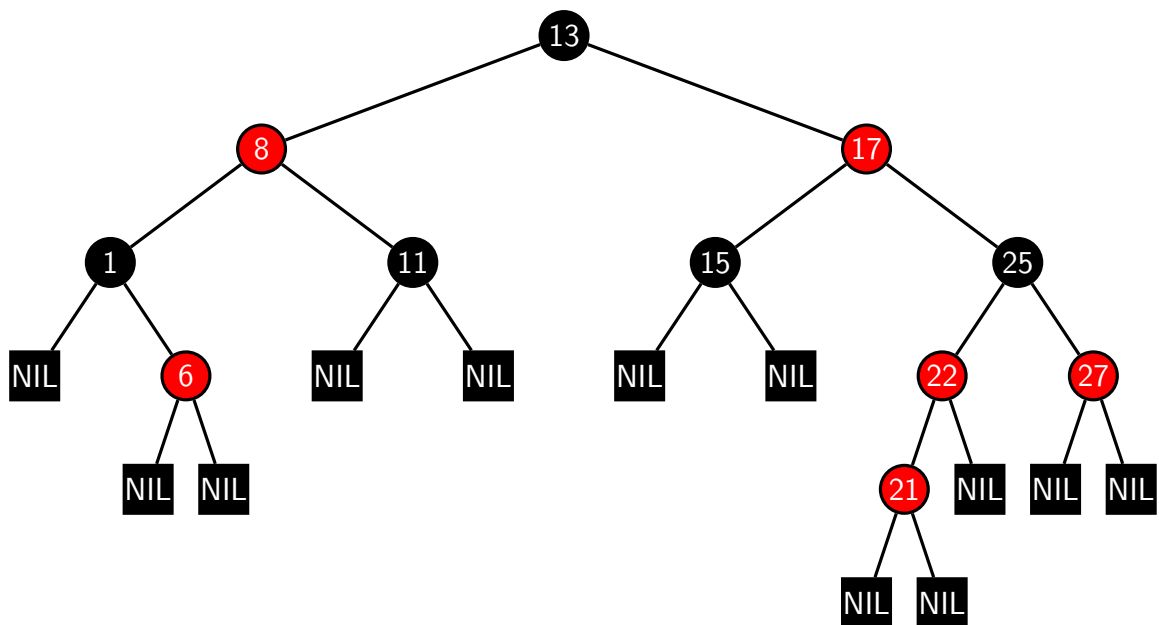


图 16.35: 插入 21

调整的方法有变色和旋转两种，而旋转又包含左旋转和右旋转两种方式。

为了重新符合红黑树的规则，有时需要把红色结点变为黑色，或是把黑色结点变为红色。

例如对于红黑树的一部分（子树），新插入的结点 Y 是红色结点，它的父结点 X 也是红色结点，不符合规则 4（红色结点的两个子结点都是黑色的），因此可以把结点 X 变为黑色。

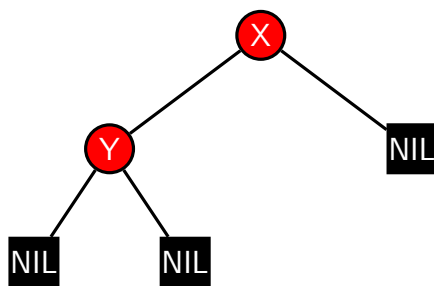


图 16.36: 违反规则 4

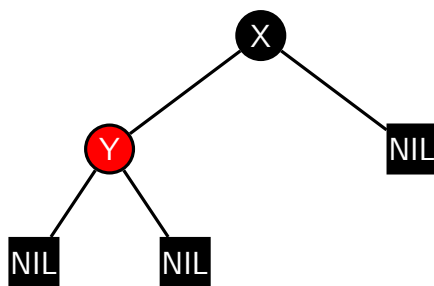


图 16.37: 变色

但是，如果这是简单的把一个结点变色，会导致相关路径凭空多出一个黑色结点，这样就会打破规则 5（从任一结点到其每个叶子的所有路径都包含相同数目的黑色结点），因此还需要其它的调整策略。

16.2.3 红黑树插入结点

红黑树插入新结点时，可以分为五种不同的局面。每一种局面有不同的调整方法。

局面 1

新结点（A）位于树根，没有父结点。

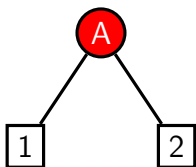
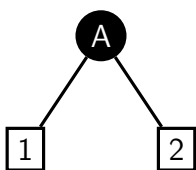


图 16.38: 局面 1

这种局面，直接让新结点变色为黑色，规则 2（根结点是黑色的）满足。同时黑色的根结点使每条路径上的黑色结点数目都增加了 1，因此并没有打破规则 5（从任一结点到其每个叶子的所有路径都包含相同数目的黑色结点）。



局面 2

新结点 (B) 的父结点 A 是黑色的。新插入的红色结点 B 并没有打破规则，无需调整。

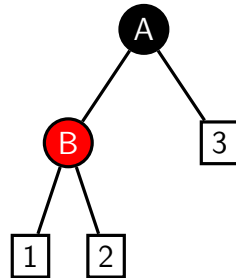


图 16.39: 局面 2

局面 3

新结点 (D) 的父结点和叔叔结点都是红色。

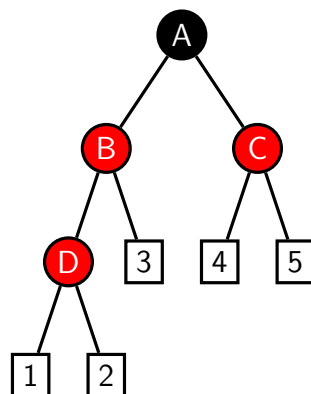
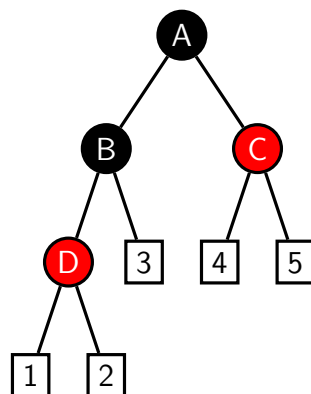
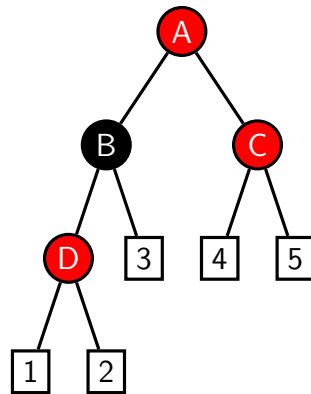


图 16.40: 局面 3

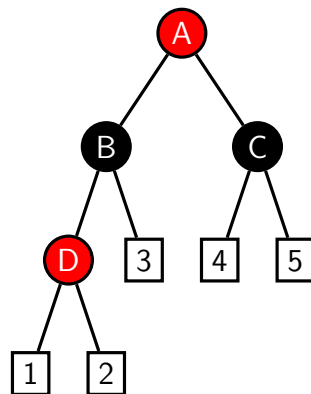
这种局面，两个红色结点 B 和 D 连续，违反了规则 4（红色结点的两个子结点都是黑色的），因此需要先让结点 B 变为黑色。



但是这样一来，结点 B 所在路径凭空多出了一个黑色结点，打破了规则 5（从任一结点到其每个叶子的所有路径都包含相同数目的黑色结点），因此再让结点 A 变为红色。



这时结点 A 和 C 又成为了连续的红色结点，再将结点 C 变为黑色。



局面 4

新结点（D）的父结点是红色，叔叔结点是黑色或者没有叔叔，且新结点是父结点的右孩子，父结点是祖父结点的左孩子。

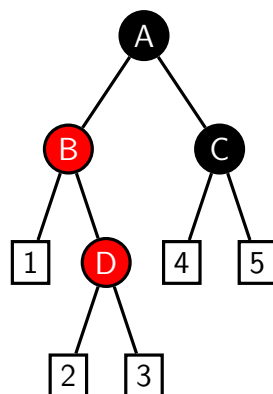
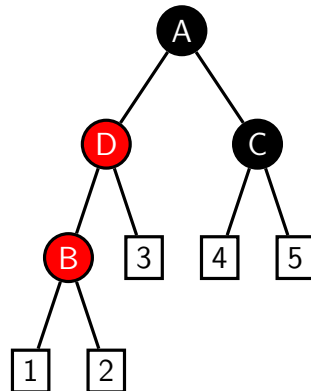


图 16.41: 局面 4

这个局面可以以结点 B 为轴，做一次左旋转，使得新结点 D 成为父结点，结点 B 成为 D 的左孩子。这样一来进入了局面 5。



局面 5

新结点 (D) 的父结点是红色，叔叔结点是黑色或者没有叔叔，且新结点是父结点的左孩子，父结点是祖父结点的左孩子。

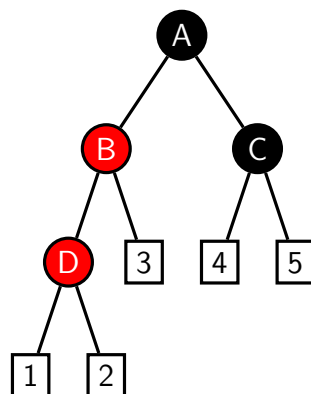
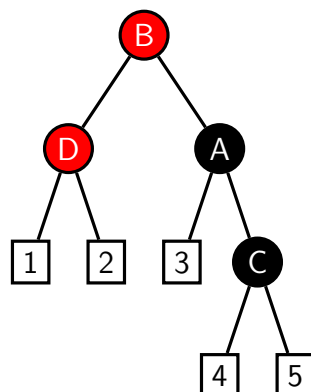
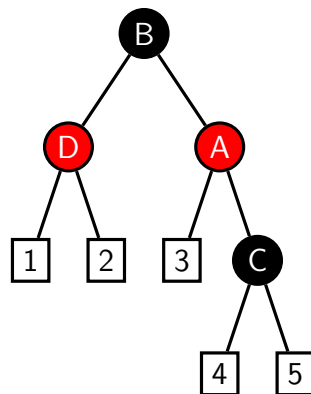


图 16.42: 局面 5

这一局面可以以结点 A 为轴，做一次右旋转，使得结点 B 成为祖父结点，结点 A 成为 B 的右孩子。

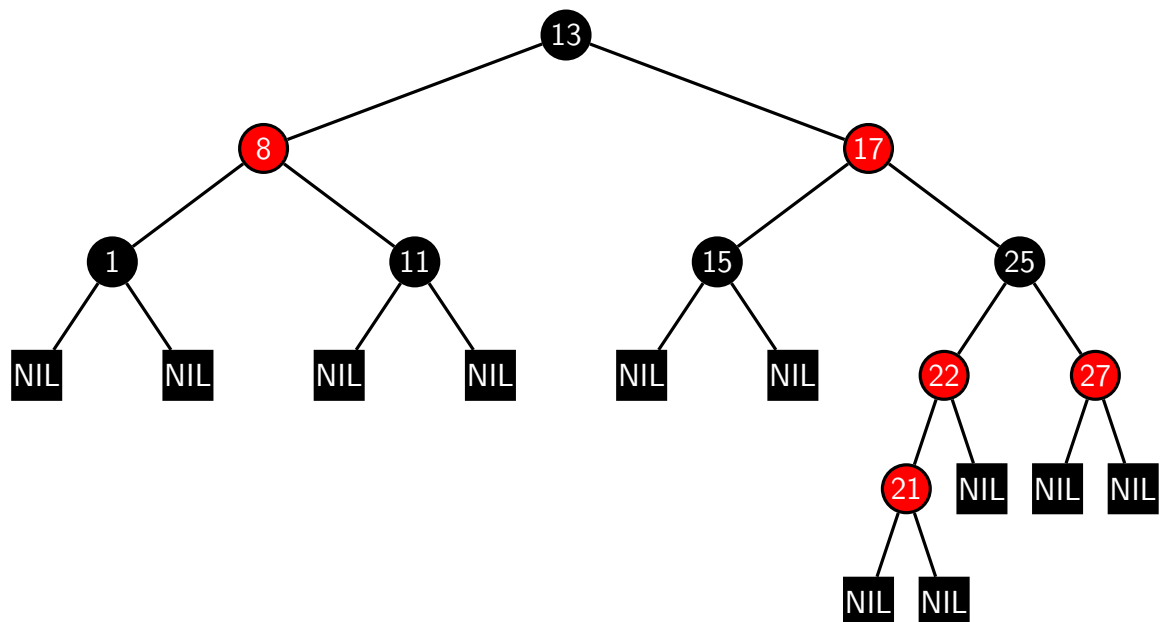


再将结点 B 变为黑色，结点 A 变为红色。

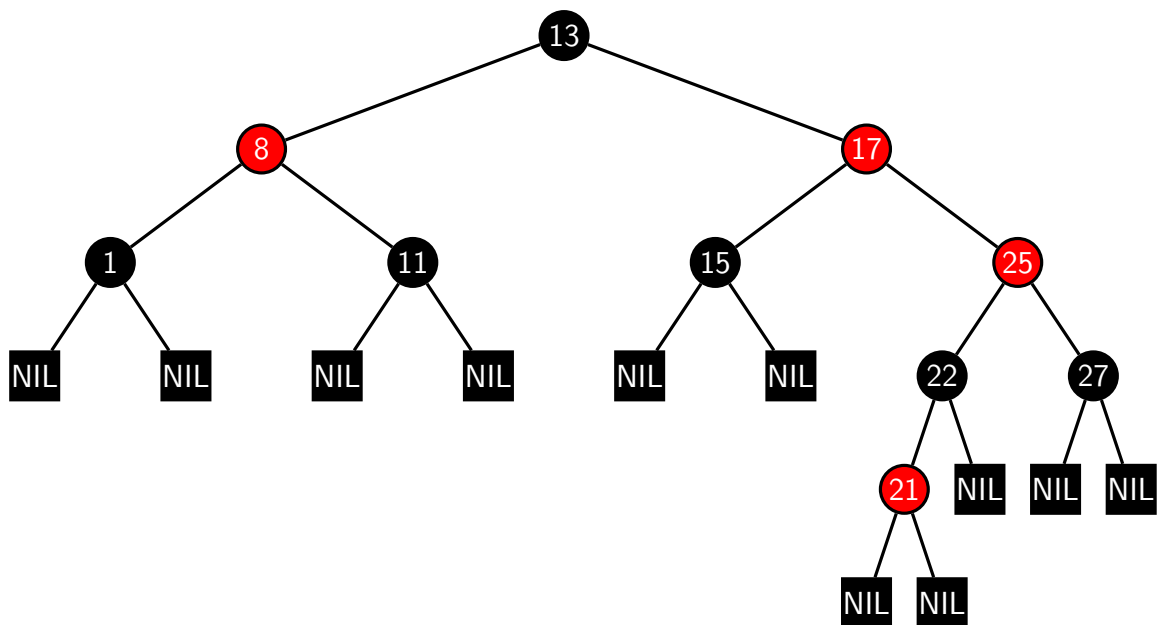


红黑树的插入操作设计到这 5 种局面。如果局面 4 中父结点 B 是右孩子，则成为了局面 5 的镜像，原本的右旋转改为左旋转；如果局面 5 中父结点 B 是右孩子，则成为了局面 4 的镜像，原本的左旋转改为右旋转。

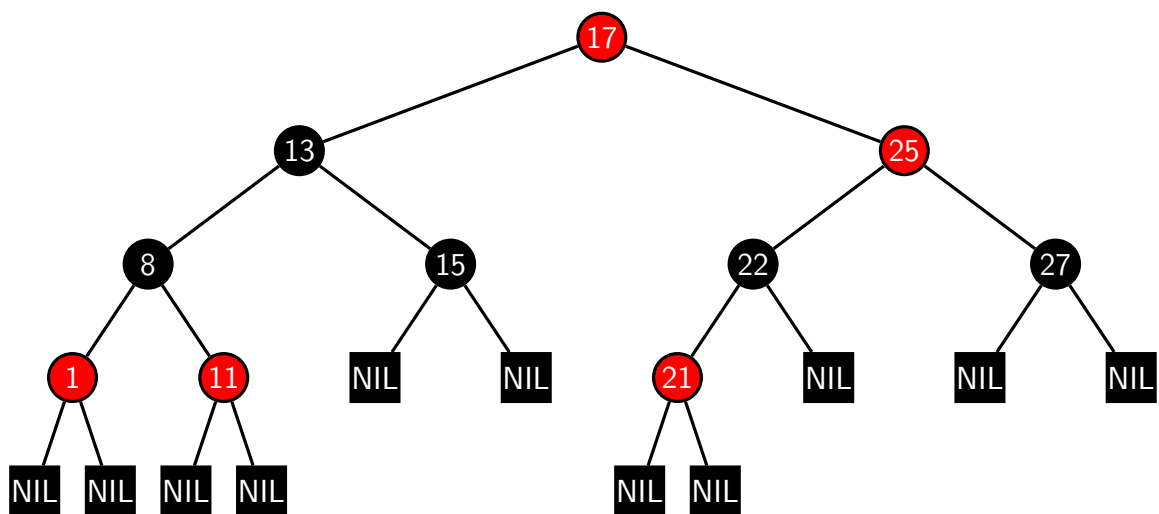
例如在一个红黑树中插入新结点 21，需要根据不同局面进行调整。



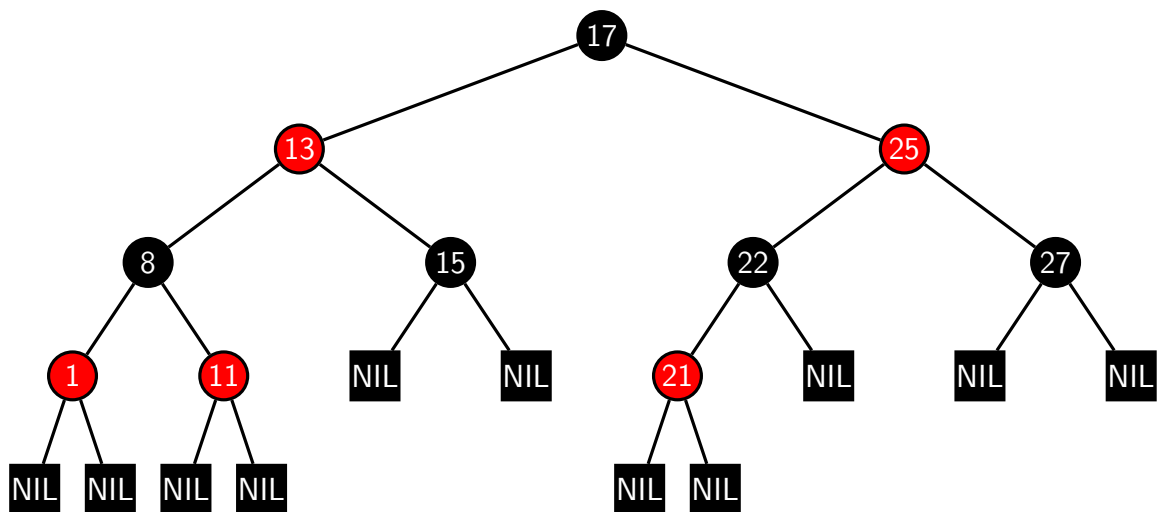
新结点 21 和它的父结点 22 是连续红色结点，违背了规则 4。当前情况符合局面 3（新结点的父结点和叔叔结点都是红色）。于是经过三次变色（22 变为黑色，25 变为红色，27 变为黑色），将以结点 25 为根的子树符合了红黑树的规则。



但结点 25 和结点 17 成为了连续红色结点，违背了规则 4。于是可以将结点 25 看做一个新结点，当前正好符合局面 5 的镜像（新结点的父结点是红色，叔叔是黑色或者没有叔叔，且新结点是父结点的右孩子，父结点是祖父结点的右孩子）。因此可以以根结点 13 为轴进行左旋转，使得结点 17 成为新的根结点。



再让结点 17 变为黑色，13 变为红色，使红黑树重新符合规则。



16.2.4 二叉查找树删除结点

在介绍红黑树的删除操作之前，需要先理解二叉查找树的删除操作。

二叉查找树的删除可分为三种情况：

待删除结点无子结点

如果待删除结点没有子结点，直接删除即可。

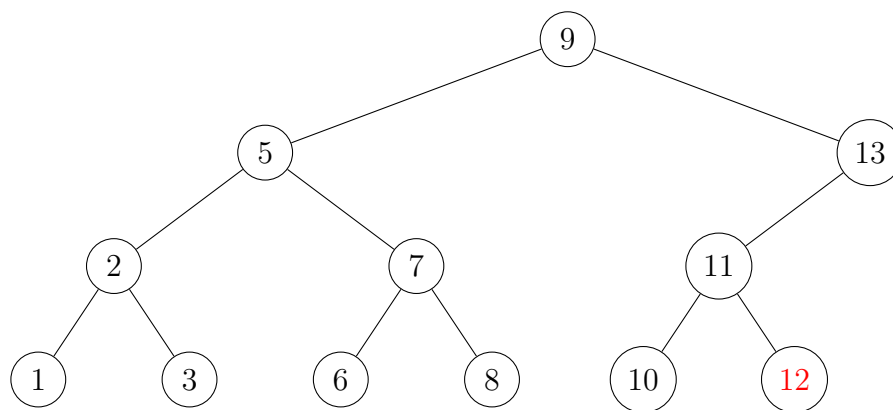


图 16.43: 删除结点 12

待删除结点只有一个孩子

如果待删除结点只有一个孩子，让孩子取代待被删除结点。

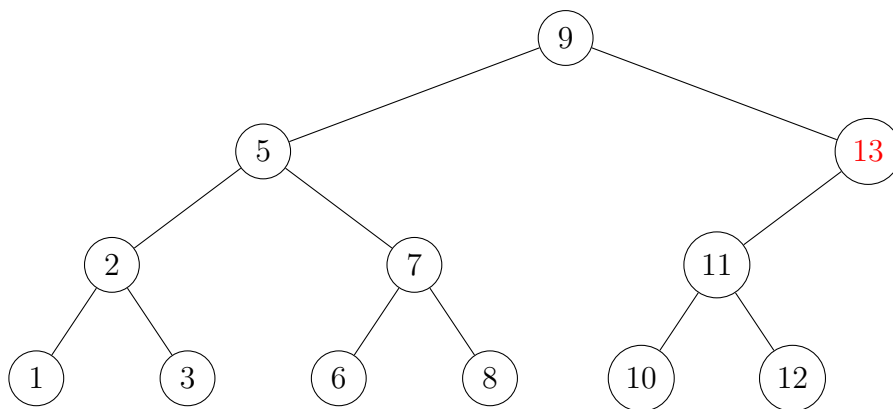
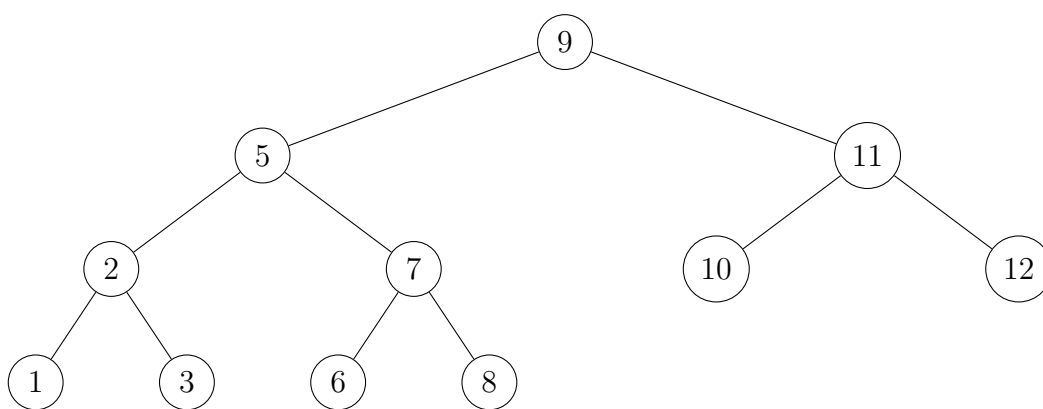


图 16.44: 删除结点 13



待删除结点有两个孩子

选择仅小于或仅大于待删除结点的结点取代，习惯上更多地会选择仅大于待删除结点的结点。

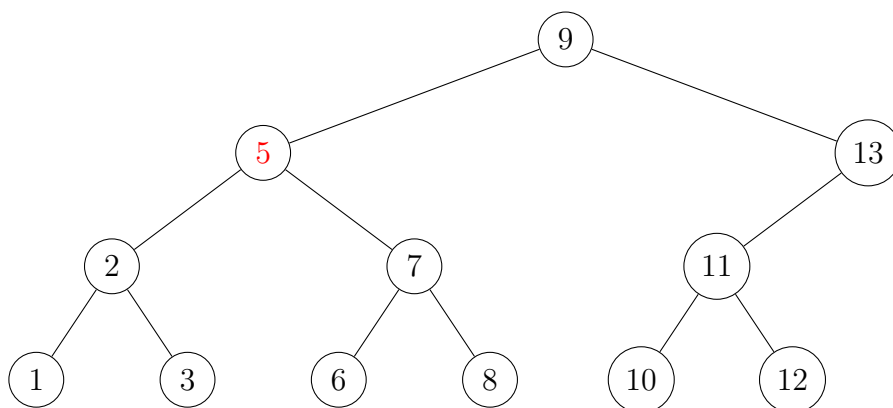
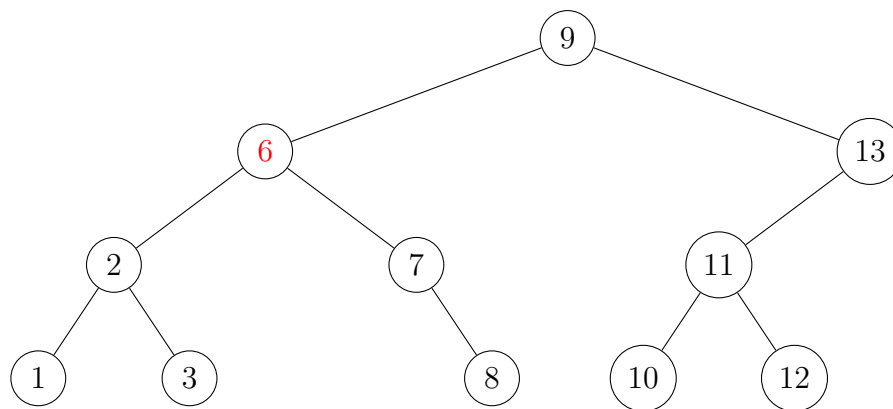


图 16.45: 删除结点 5



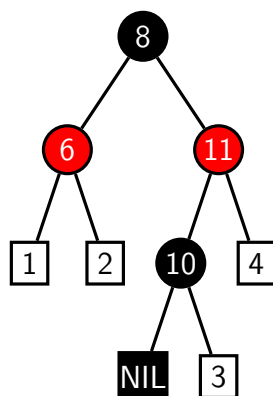
16.2.5 红黑树删除结点

红黑树的删除操作要比插入操作复杂得多。

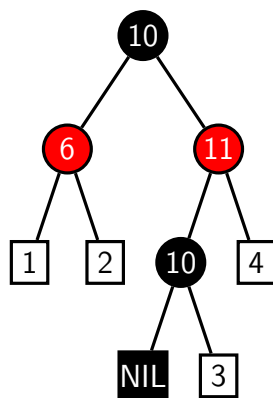
第一步

如果待删除结点有两个非空的孩子结点，转化成待删除结点只有一个孩子（或没有孩子）的情况。

例如删除结点 8：



因为结点 8 有两个孩子，可以选择仅大于 8 的结点 10 复制到 8 的位置，结点颜色变成待删除结点的颜色。

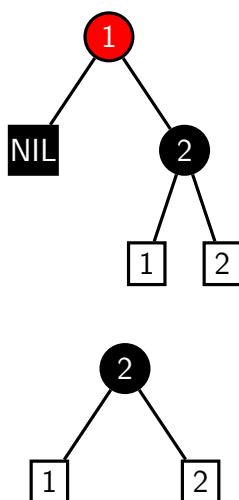


结点 10 能成为仅大于 8 的结点，必定没有左孩子结点，所以问题转换成了待删除结点只有一个右孩子（或者没有孩子）的情况。

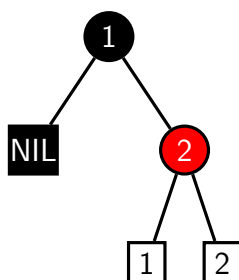
第二步

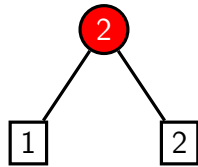
根据待删除结点和其唯一子结点的颜色，分情况处理。

情况 1：自身是红色，子结点是黑色。直接按照二叉查找树的删除操作，删除结点 1 即可。

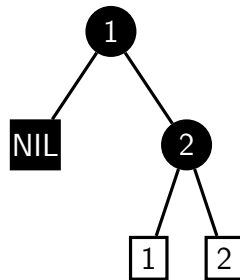


情况 2：自身是黑色，子结点是红色。按照二叉查找树的删除操作，删除结点 1。此时这条路径凭空少了一个黑色结点，因此需要将结点 2 变成黑色即可。





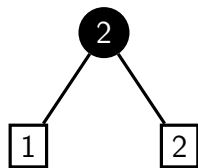
情况 3：自身是黑色，子结点也是黑色，或者子结点是空叶子结点。这种情况最为复杂，涉及到很多变化。首先还是按照二叉查找树的删除操作，删除结点 1。此时这条路径凭空少了一个黑色结点，而且并不能改变结点 2 的颜色来解决问题。这时需要进入第三步，专门解决父子双黑的情况。



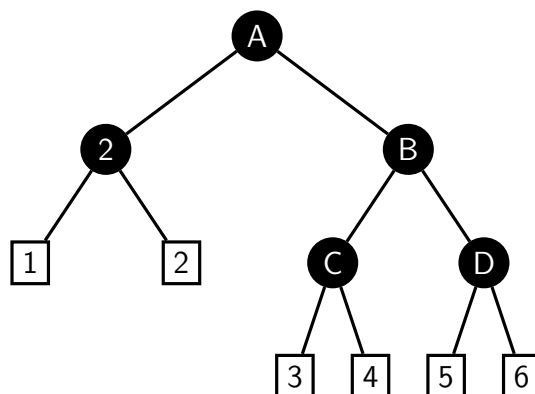
第三步

遇到双黑结点，在子结点顶替父结点后，可分为 6 种情况处理。

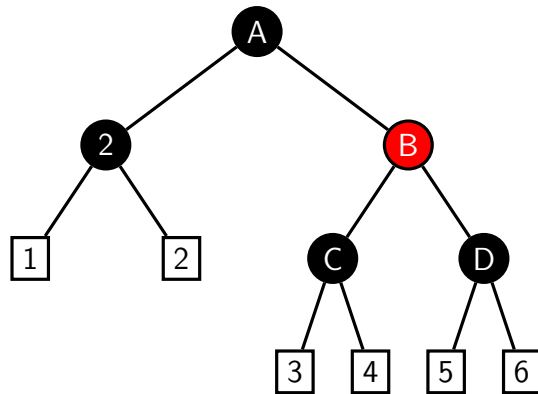
情况 1：结点 2 是红黑树的根。此时所有路径都减少了一个黑色结点，并未打破规则，无需调整。



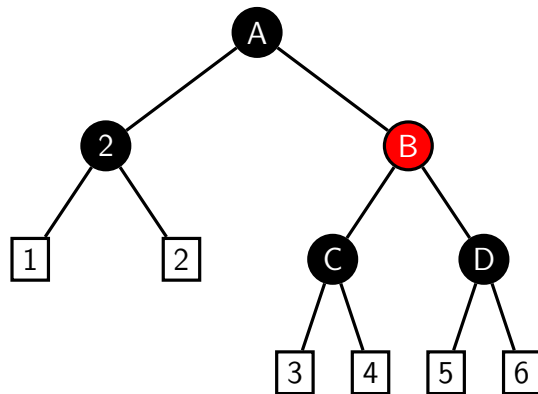
情况 2：结点 2 的父亲、兄弟和侄子结点都是黑色。



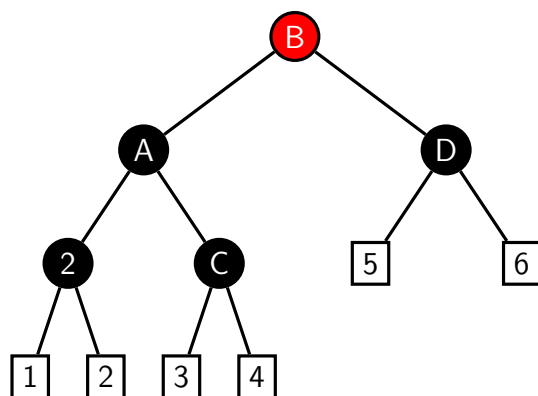
直接把结点 2 的兄弟结点变为红色。这样结点 B 所在路径少了一个黑色结点，两边扯平了。



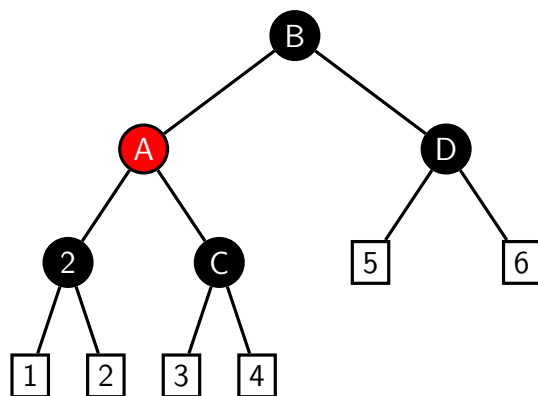
情况 3： 结点 2 的兄弟结点是红色。



以结点 2 的父结点 A 为轴进行左旋转。

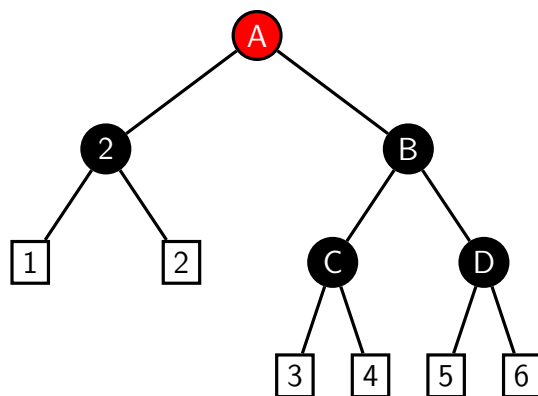


然后结点 A 变为红色，结点 B 变为黑色。

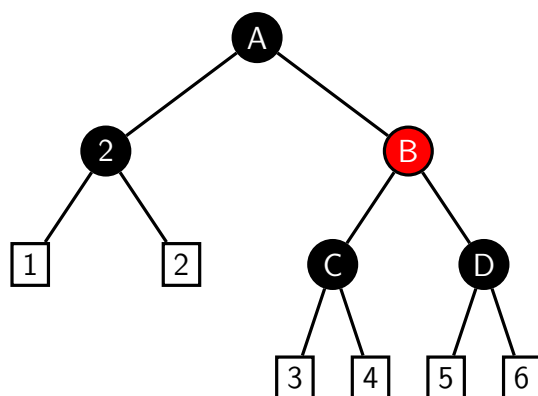


这样的变化就有可能转换成情况 4、5、6 中的任意一种。

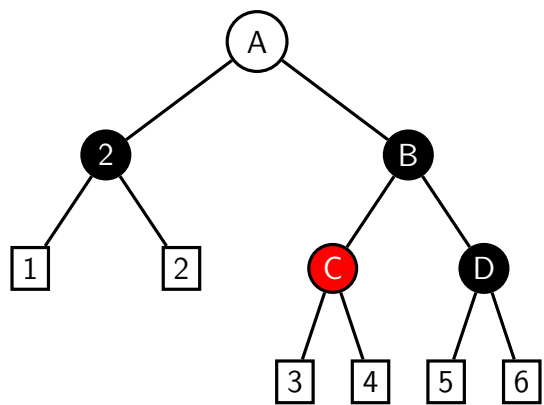
情况 4： 结点 2 的父结点是红色，兄弟和侄子结点是黑色。



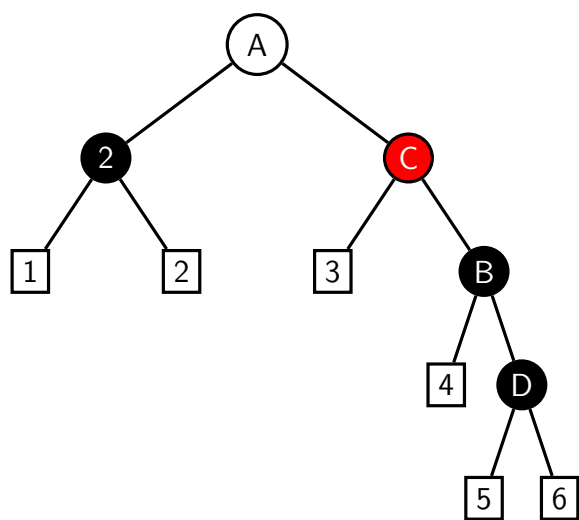
将结点 2 的父结点 A 变为黑色，兄弟结点 B 变为红色。这样结点 2 的路径补充了黑色结点，而结点 B 的路径并没有减少黑色结点。



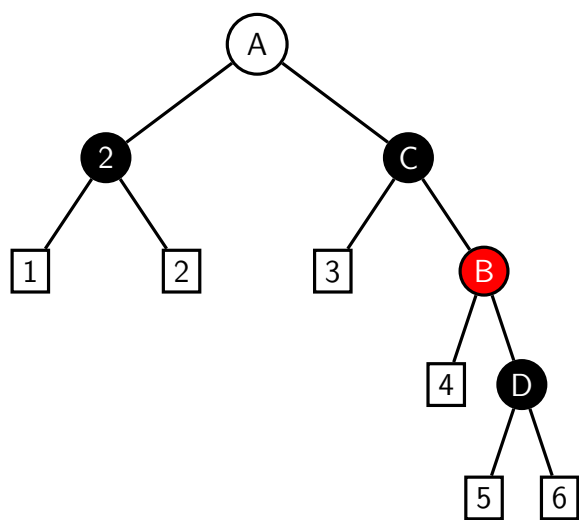
情况 5： 结点 2 的父结点随意，兄弟结点 B 是黑色右孩子，结点 2 的左侄子是红色，右侄子是黑色。



以结点 2 的兄弟结点 B 为轴进行右旋转。

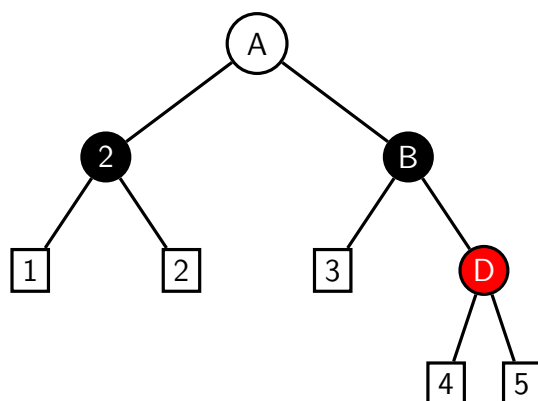


将结点 B 变为红色，结点 C 变为黑色。

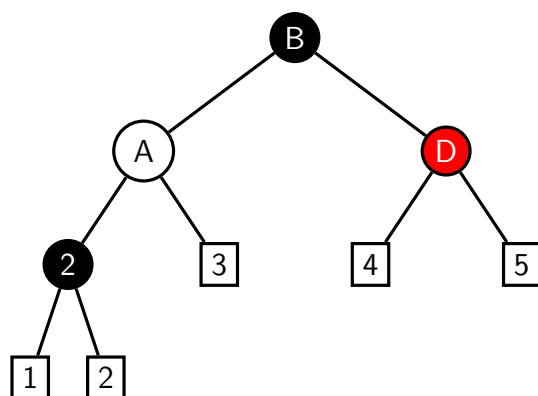


这样的变化就转换成了情况 6。

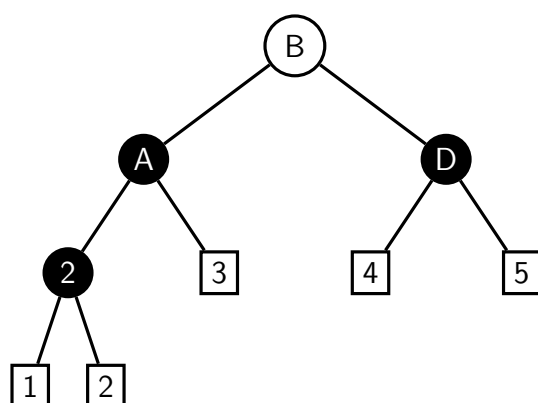
情况 6： 结点 2 的父结点随意，兄弟结点 B 是黑色右孩子，结点 2 的右侄子的红色。



以结点 2 的父结点 A 为轴进行左旋转。

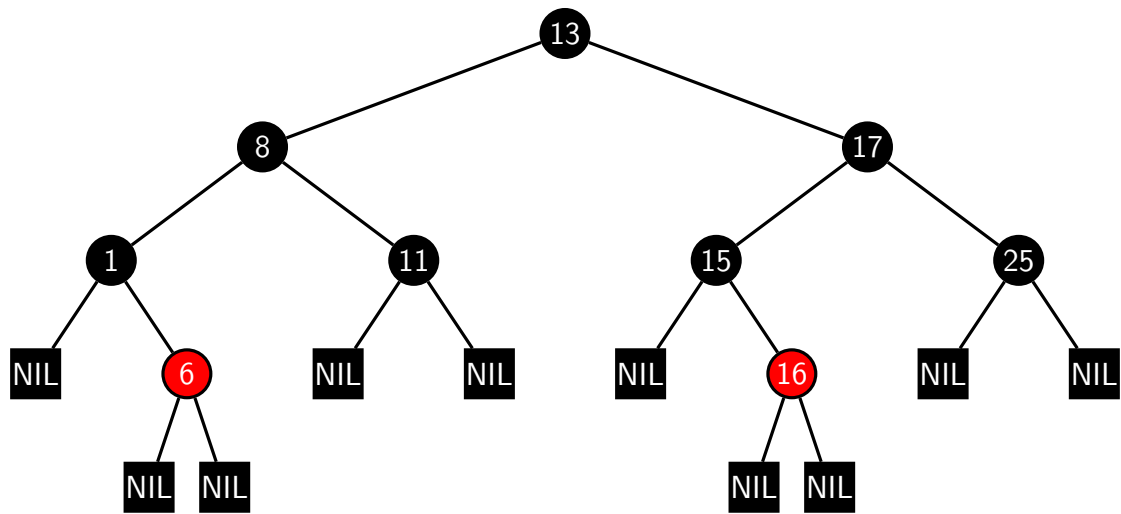


将结点 A 和 B 的颜色交换，让结点 D 变为黑色。

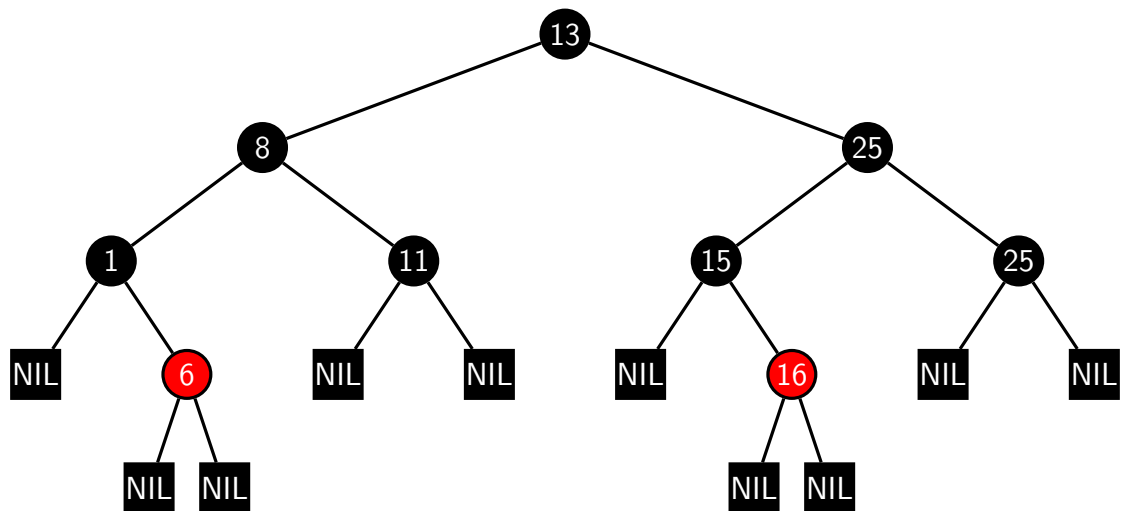


这样经过结点 2 的路径由之前的随机 + 黑变成了随机 + 黑 + 黑，补充了一个黑色结点。经过结点 D 的路径由之前的随机 + 黑 + 红变成了随机 + 黑，黑色结点并没有减少。

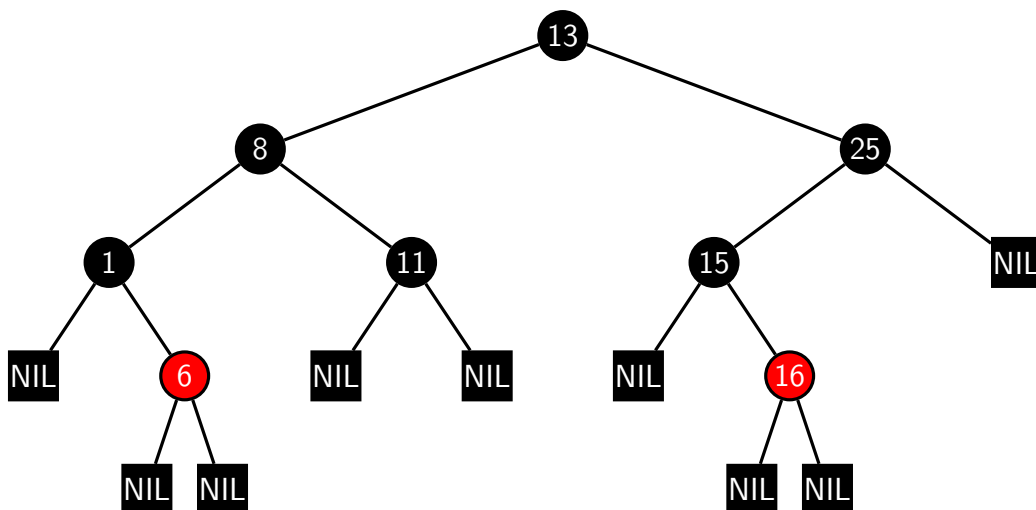
例如在一个红黑树中删除结点 17，需要根据不同情况进行调整。



由于待删除结点 17 有两个孩子，子树当中仅大于 17 的结点是 25，所以把结点 25 复制到 17 的位置，保持黑色。

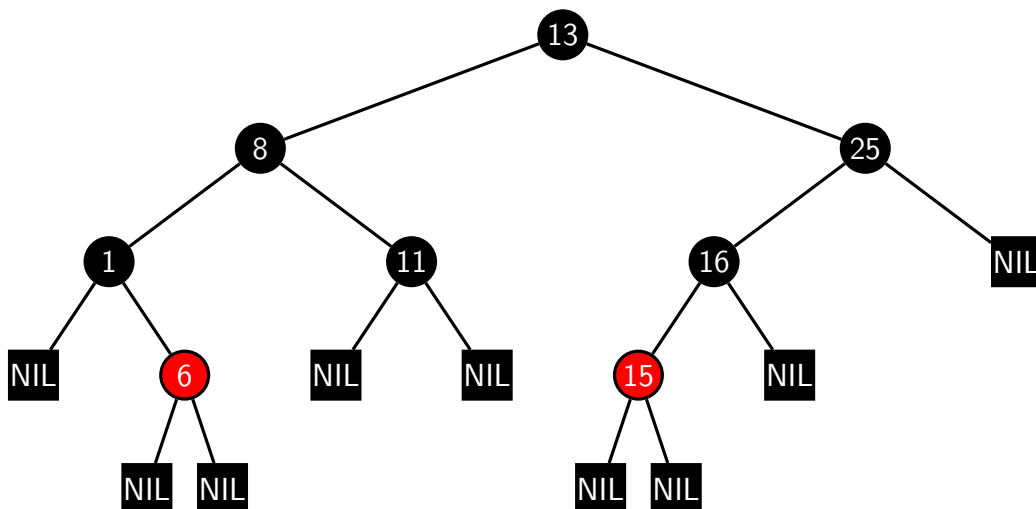


接着需要删除原本的结点 25，这个情况对应的是删除结点操作中的第二步情况 3（待删除结点是黑色，子结点是空叶子结点）。

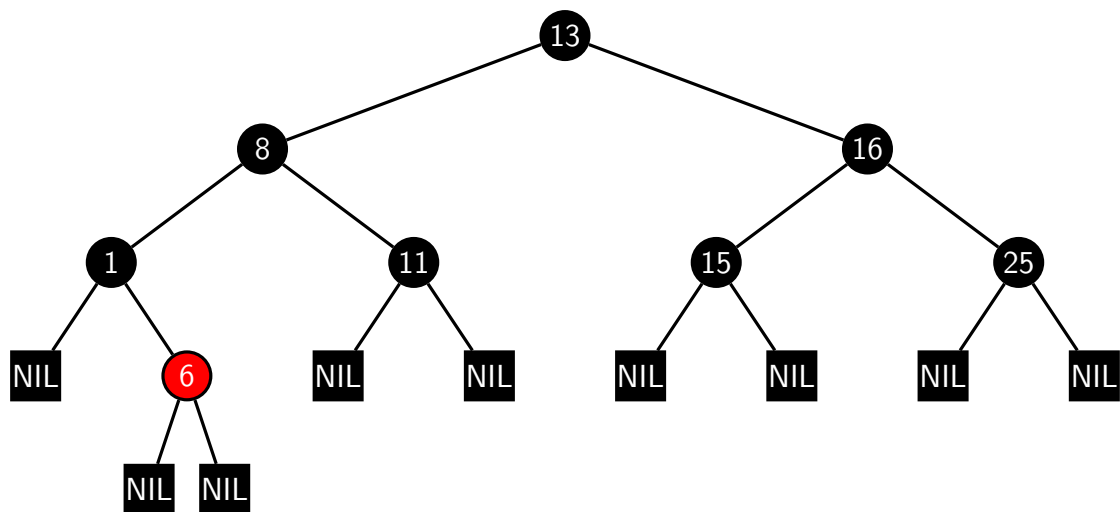


此时，以结点 25 为根的子树符合第三步情况 5 的镜像（结点 NIL 的父结点随意，兄弟结点 15 是黑色左孩子，结点 NIL 的右侄子是红色，左侄子是黑色）。

通过左旋转和变色，将子树转化成第三步情况 6 的镜像（结点 NIL 的父结点随意，兄弟结点 16 是黑色左孩子，结点 NIL 的左侄子是红色）。



通过右旋转和变色，使整棵二叉树重新符合红黑树的规则。



16.2.6 红黑树与 AVL 树的区别

AVL 树是严格平衡的二叉树，要求每个结点的左右子树高度差不超过 1，而红黑树则要宽松一些，要求任何一条路径的长度不超过其它路径长度的 2 倍。

正因为这个差别，AVL 树的查找效率更高，但平衡调整的成本也更高。在需要频繁查找时，选用 AVL 树更合适，在需要频繁插入删除时，选用红黑树更合适。

16.3 B 树

16.3.1 B 树

B-树就是 B 树，中间的是连字符而不是减号，谁要是把 B-树读成“B 减树”，那可就丢人现眼了。

B 树主要应用于文件系统以及部分数据库索引，比如著名的非关系型数据库 MongoDB。数据库索引使用树型结构的原因在于树的查询效率高，而且可以保持有序。既然如此，为什么索引没有采用二叉查找树来实现呢？二叉查找树查询的时间复杂度是 $O(\log n)$ ，性能已经足够高了，难道 B 树可以比它更快？

其实从算法逻辑上来讲，二叉查找树的查找速度和比较次数都是最小的。但是，我们不得不考虑一个现实的问题——磁盘 I/O。数据库索引是存储在磁盘上的，当数据量比较大的时候，索引的大小可能有好几个 G 甚至更多。当利用索引查询的时候，显然不可能把整个索引全部加载到内存，能做的只有逐一加载每一页的磁盘页，这里的磁盘页就对应着索引树的结点。

如果利用二叉查找树作为索引结构，假设需要查询数据 10，一共需要进行 4 次磁盘 I/O。

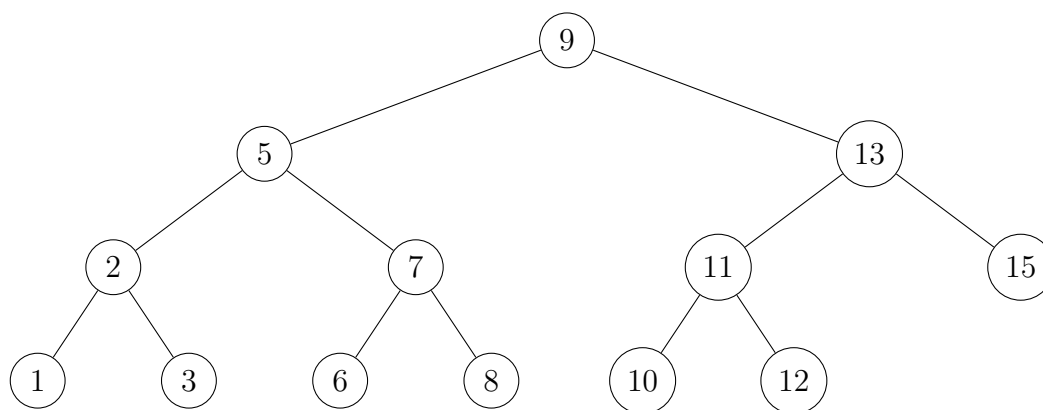


图 16.46: 索引树

最坏情况下，磁盘 I/O 的次数等于索引树的高度。既然如此，为了减少磁盘 I/O 次数，就需要把原本瘦高的树结构变得矮胖，这就是 B 树的特征之一。

B 树是一种多路平衡查找树，它的每一个结点最多包含 k 个孩子， k 被称为 B 树的阶， k 的大小取决于磁盘页的大小。

一个 m 阶的 B 树具有以下特征：

1. 根结点至少有 2 个孩子。
2. 每个中间结点都包含 $k - 1$ 个元素和 k 个孩子 ($\frac{m}{2} \leq k \leq m$)。
3. 每个叶子结点都包含 $k - 1$ 个元素 ($\frac{m}{2} \leq k \leq m$)。
4. 所有叶子结点都位于同一层。
5. 每个结点中的元素从小到大排列，结点中 $k - 1$ 个元素正好是 k 个孩子包含的元素的值域分划。

例如对于一个 3 阶 B 树：

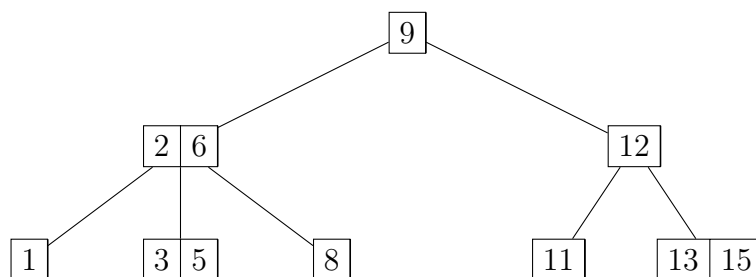


图 16.47: 3 阶 B 树

其中结点 (2, 6) 包含两个元素，并且有三个孩子 1、(3, 5) 和 8，其中 1 小于元素 2，(3, 5) 在元素 2 和 6 之间，8 大于 6。

16.3.2 查询结点

这树长得真奇怪，它真的能实现高效查询吗？

例如查询数据 5，只需进行 3 次磁盘 I/O。B 树在查询中的比较次数其实不比二叉查找树少，尤其当单一结点中的元素数量很多时。可是相比磁盘 I/O 的速度，内存中的比较耗时是几乎可以忽略的。所以只要树的高度足够低，I/O 次数足够少，就可以提升查询效率。相比之下结点内部元素多一些也没有关系，仅仅

是多了几次内存交互，只要不超过磁盘页的大小即可。这些就是 B 树的优势之一。

16.3.3 插入结点

例如在 3 阶 B 树中插入元素 4，由于结点 (3, 5) 已经是两元素结点，无法再增加。父结点 (2, 6) 也是两元素结点，也无法再增加。根结点 9 是单元素结点，可以升级为两元素结点。于是拆分结点 (3, 5) 和结点 (2, 6)，让根结点 9 升级为 (4, 9)，结点 6 独立为根结点的第二个孩子。

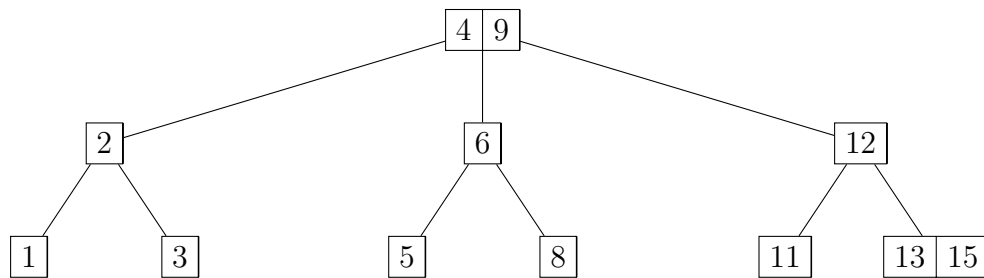


图 16.48: 插入结点 4

为了插入一个元素，让整个 B 树的那么多结点都发生了连锁改变，确实有点麻烦。但也正因为如此，让 B 树能够始终维持多路平衡。因此自平衡是 B 树的另一大优势。

16.3.4 删除结点

例如在 3 阶 B 树中删除元素 11，删除结点 11 后，结点 12 只有一个孩子，不符合 B 树特征。因此找出 12、13 和 15 中的中位数 13，取代结点 12，而结点 12 下移成为第一个孩子。

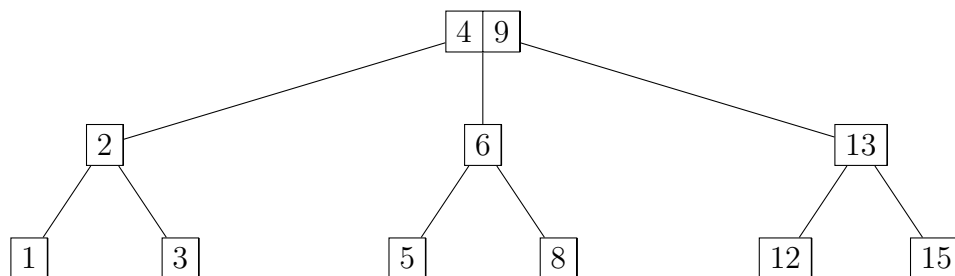


图 16.49: 删除结点 11

16.4 B+ 树

16.4.1 B+ 树

B+ 树是 B 树的一种变体，有着比 B 树更高的查询性能。B+ 树和 B 树有一些共同点，但是 B+ 树也具备一些新的特征。

一个 m 阶的 B+ 树具有以下特征：

1. 有 k 个子树的中间结点包含 k 个元素（B 树中是 $k - 1$ 个元素）。
2. 所有的叶子结点中包含了全部元素的信息，及指向含这些元素记录的指针，且叶子结点本身按照关键字大小连接。
3. 所有中间结点元素都同时存在于子结点，在子结点元素中是最大（或最小）元素。

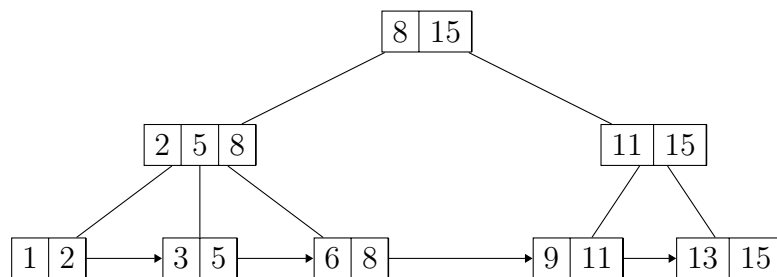


图 16.50: B+ 树

这又是什么怪树？不但结点之间含有重复元素，而且叶子结点之间还用指针连在一起。

这些正是 B+ 树的几个特性。首先，每一个父结点的元素都出现在子结点中，是子结点的最大（或最小）元素。因此根结点的最大元素也就是整个 B+ 树的最大元素，之后无论插入或删除多少元素，始终要保持最大元素在根结点中。至于叶子结点，由于父结点的元素都出现在子结点，因此所有叶子结点包含了全部元素信息。并且每一个叶子结点都带有指向下一个结点的指针，形成了一个有序链表。

16.4.2 单行查询

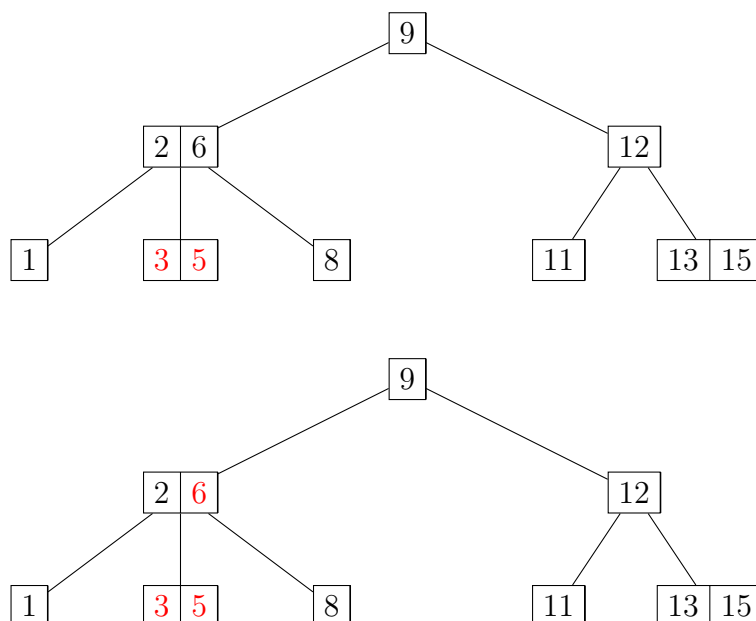
B+ 树的好处主要体现在查询性能上。在单元素查询的时候，B+ 树会自顶向下逐层查找结点，最终找到匹配的叶子结点。例如在 B+ 树中查询元素 3，需要进行 3 次磁盘 I/O。

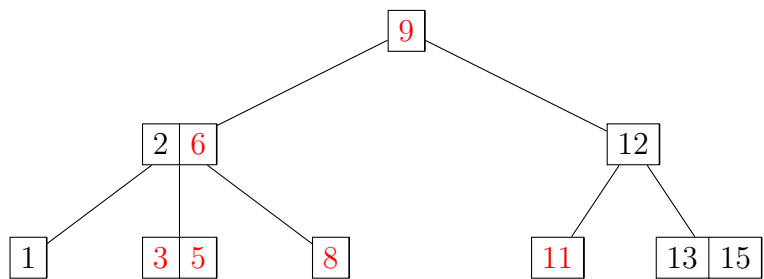
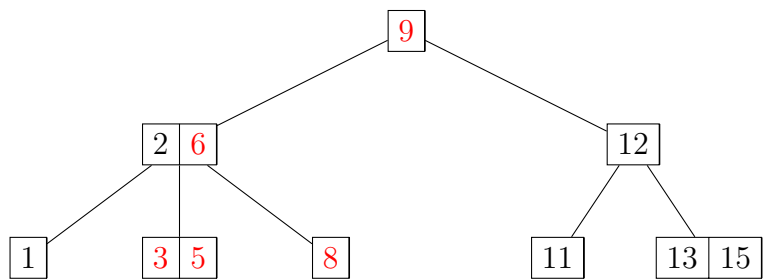
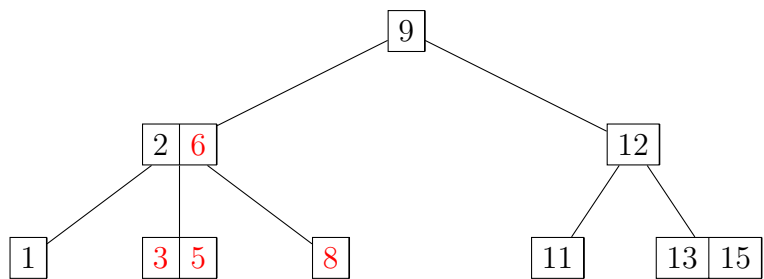
虽然流程看起来和 B 树差不多，但是有两点不同。首先，B+ 树的中间结点并不关联数据，仅仅是索引，所以同样大小的磁盘页可以容纳更多的结点元素。这就意味着，数据量相同的情况下，B+ 树的结构比 B 树更加矮胖，因此查询时 I/O 次数也更少。

其次，B+ 树的查询性能必须最终查找到叶子结点，而 B 树只要找到匹配元素即可。因此，B 树的查询性能并不稳定，最好情况下只查询根结点，最坏情况下是查到叶子结点，而 B+ 树的每一次查询都是稳定的。

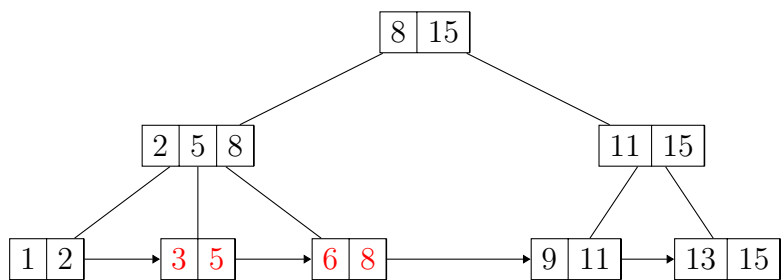
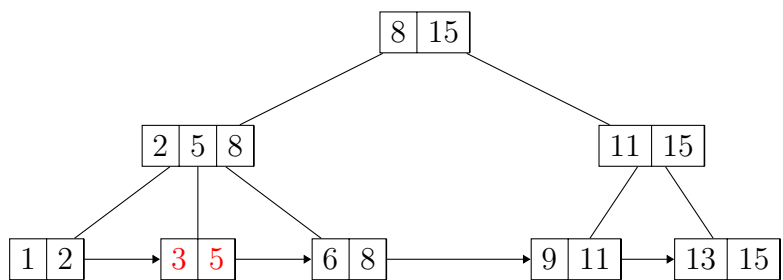
16.4.3 范围查询

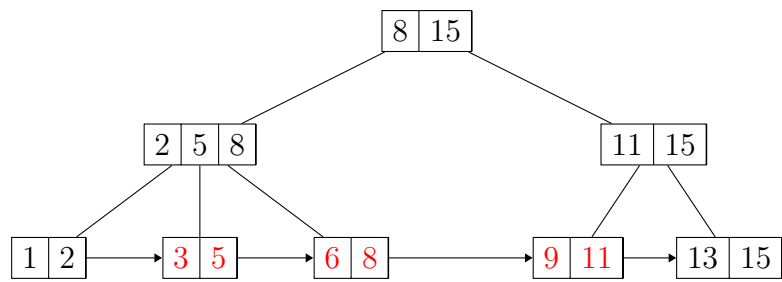
B 树中进行范围查询只能通过繁琐的中序遍历。例如在 B 树中查询范围为 3 到 11 的元素：





反观 B+ 树的范围查询，则要简单的多，只需要在链表上做遍历即可：





16.5 并查集

16.5.1 并查集 (Disjoint Set)

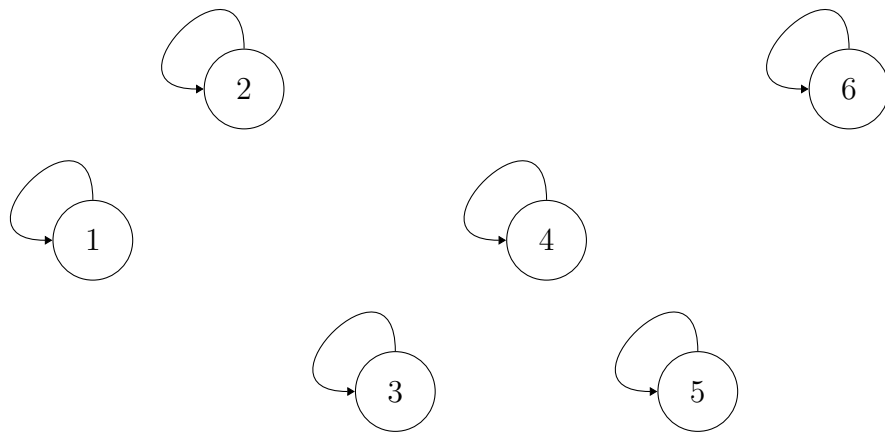
并查集是一种简洁优雅的数据结构，主要用于解决一些元素分组的问题。

它管理一系列不相交的集合，并支持两种操作：

1. 合并 (union)：把两个不相交的集合合并为一个集合。
2. 查询 (find)：查询两个元素是否在同一个集合中。

例如在某个家族中，如果 x 和 y 是亲戚， y 和 z 是亲戚，那么 x 和 z 也是亲戚。如果 x 和 y 是亲戚，那么 x 的亲戚都是 y 的亲戚， y 的亲戚也都是 x 的亲戚。为了判断两个人是否为亲戚，只需判断他们是否属于同一个集合即可。

并查集的思想在于用集合中的一个元素代表集合，类似于把集合看做帮派，代表元素则是帮主，最开始所有元素各自为一个集合（各自的帮主就是自己）。



例如将元素 1 和 3 合并，就是将 1 和 3 比武，假设 1 赢了，3 就认 1 作帮主。

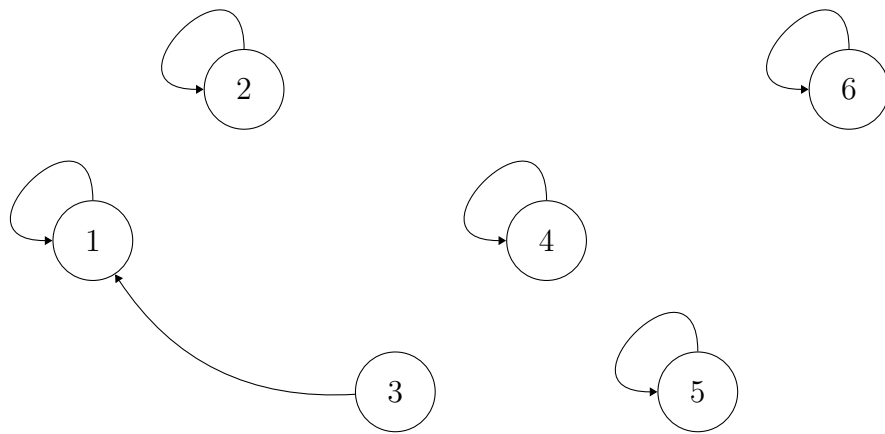


图 16.51: 合并 1 和 3

再合并元素 2 和 3，但 3 表示“别跟我打，让我的帮主来收拾你”。假设还是 1 赢了，2 也认 1 作帮主。

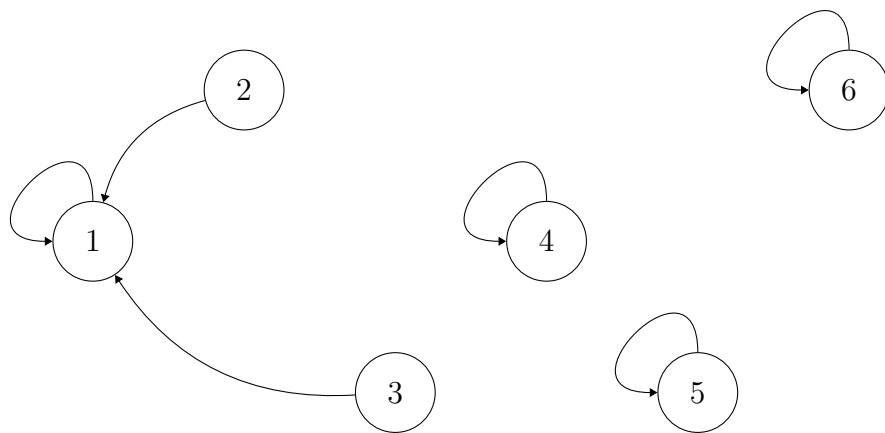


图 16.52: 合并 2 和 3

假设元素 4、5、6 也进行了相关合并：

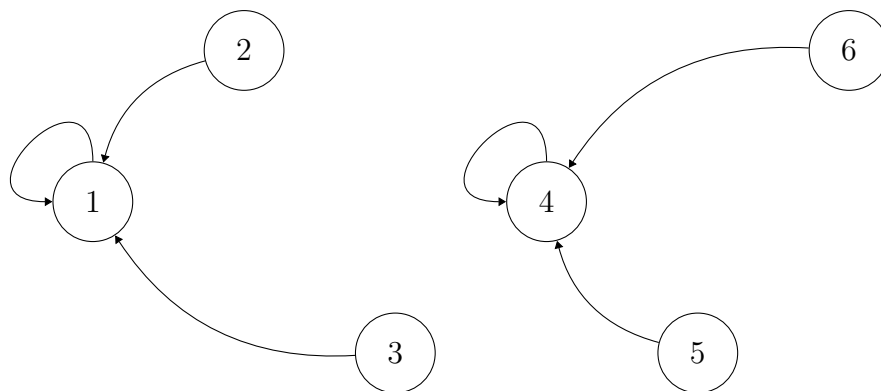


图 16.53: 合并 4、5、6

现在再合并元素 2 和 6，将它们的帮主 1 和 4 比武，假设 1 胜利后，4 认 1 为帮主，当然它的手下也都跟着投降了。

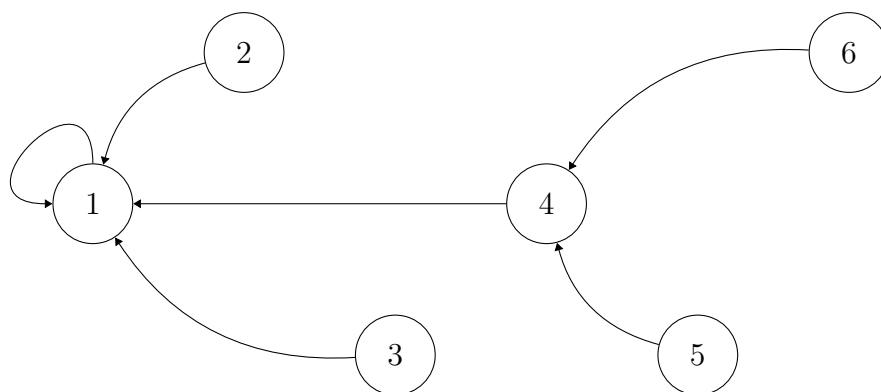


图 16.54: 合并 2 和 6

并查集是一种树型结构，要寻找集合的代表元素，只需要一层一层向上访问父结点，直达树的根结点即可。根结点的父结点就是它自己。

假设有 n 个元素，利用数组 `parent` 存放每个元素的父结点，一开始每个元素的父结点为自己。

初始化

```
1 void init(int n) {
2     parent = (int *)malloc(sizeof(int) * n);
3     for(int i = 0; i < n; i++) {
4         parent[i] = i;
```



```
5     }  
6 }
```

要判断两个元素是否属于同一个集合，只需要看它们的根结点是否相同。利用递归的方法可以实现对代表元素的查询，一层一层访问父结点，直至根结点。

查询

```
1 int find(int val) {  
2     if(parent[val] == val) {  
3         return val;  
4     }  
5     return find(parent[val]);  
6 }
```

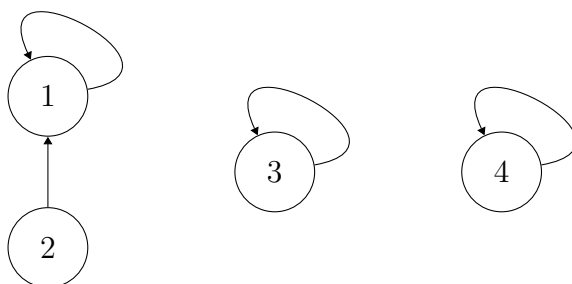
合并操作需要先找到两个集合的根结点，将前者的父结点设置为后者即可（也可将后者的父结点设置为前者）。

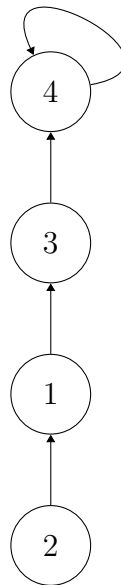
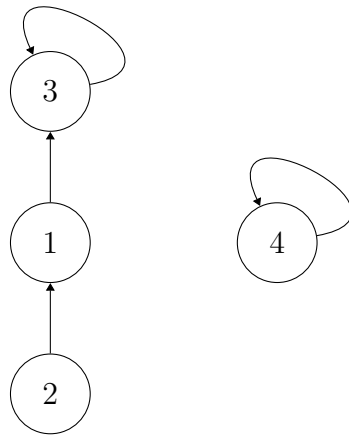
合并

```
1 void merge(int i, int j) {  
2     parent[find(i)] = find(j);  
3 }
```

16.5.2 路径压缩 (Path Compression)

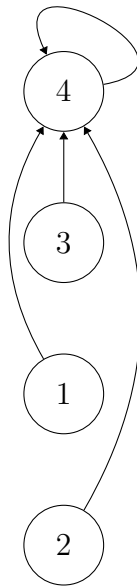
最简单的并查集效率是比较低的，分别进行 `merge(2, 3)` 和 `merge(2, 4)`：





这样可能会形成一条长链，随着链越来越长，想要从底部找到根结点会变得越来
越难。利用路径压缩的方法可以解决这个问题，因为我们只关心一个元素所对应
的根结点，因此每个元素到根结点的路径最好尽可能短。

实现的时候只需在查询过程中，把沿途的每个结点的父结点都设置为根结点即
可。在下一次查询的时候，就可以节省很多时间。

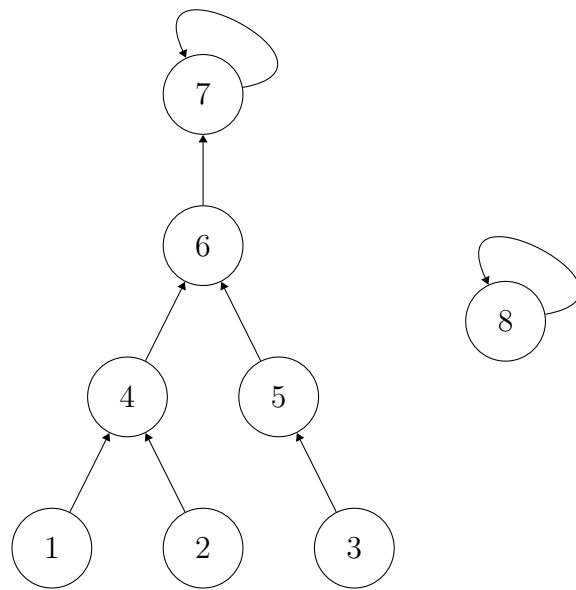


查询（路径压缩）

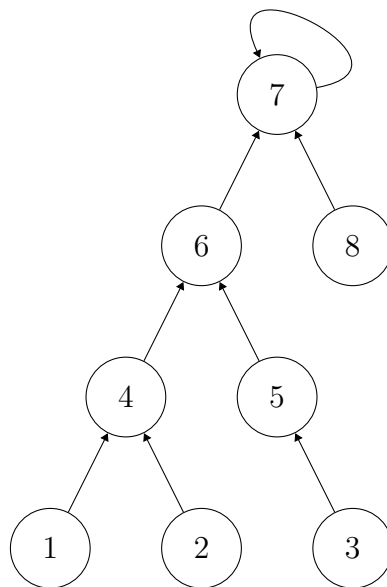
```
1 int find(int val) {  
2     if(parent[val] == val) {  
3         return val;  
4     } else {  
5         parent[val] = find(parent[val]);  
6         return parent[val];  
7     }  
8 }
```

16.5.3 按秩合并

如果需要将一棵比较复杂的树与一个单元素进行合并，例如 `merge(7, 8)` 时，是把元素 7 的父结点设为 8 好，还是把 8 的父结点设为 7 呢？



如果把 7 的父结点设为 8，会使树的深度加深，原来树中的每个元素到根结点的距离都变长了，之后寻找根结点的路径也会变长。虽然有路径压缩，但路径压缩也是要消耗时间的。而把 8 的父结点设为 7，并不会影响到不相关的结点。



因此在合并两个集合时，应该把简单的树往复杂的树上合并。利用数组 rank 记录每个结点对应的树的深度，一开始所有元素的 rank 设为 1。合并时比较两个根结点，把 rank 较小者往较大者上合并。深度相同的情况下，无论如何合并，都会使树的深度增加 1。

按秩合并

```
1 void init(int n) {
2     parent = (int *)malloc(sizeof(int) * n);
3     rank = (int *)malloc(sizeof(int) * n);
4     for(int i = 0; i < n; i++) {
5         parent[i] = i;
6         rank[i] = 1;
7     }
8 }
9
10 void merge(int i, int j) {
11     // 找到对应根结点
12     int x = find(i);
13     int y = find(j);
14     if(rank[x] <= rank[y]) {
15         parent[x] = y;
16     } else {
17         parent[y] = x;
18     }
19     // 如果深度相同且根结点不同，则新的根结点深度+1
20     if(rank[x] == rank[y] && x != y) {
21         rank[y]++;
22     }
23 }
```