

讲排序的时候，都是用整数来举例，但在真正软件开发中，我们要排序的往往不是单纯的整数，而是一组对象，我们需要按照对象的某个 key 来排序。

比如说，

现在要给电商交易系统里的“订单”排序。订单有两个属性，一个是下单时间，另一个是订单金额。如果我们现在有 10 万条订单数据，我们希望按照金额从小到大对订单数据排序。对于金额相同的订单，我们希望按照下单时间从早到晚有序。

最先想到的方法是：我们先按照金额对订单数据进行排序，然后，再遍历排序之后的订单数据，对于每个金额相同的小区再按照下单时间排序。这种排序思路理解起来不难，但是实现起来会很复杂。

还有更好的办法吗？

**原地排序 (Sorted in place)** 就是特指空间复杂度是  $O(1)$  的排序算法。

**稳定性** 如果待排序的序列中存在值相等的元素，经过排序之后，相等元素之间原有的先后顺序不变。

借助稳定排序算法，这个问题可以非常简洁地解决。

解决思路是这样的：我们先按照下单时间给订单排序，注意是按照下单时间，不是金额。排序完成之后，我们用稳定排序算法，按照订单金额重新排序。两遍排序之后，我们得到的订单数据就是按照金额从小到大排序，金额相同的订单按照下单时间从早到晚排序的。为什么呢？

稳定排序算法可以保持金额相同的两个对象，在排序之后的前后顺序不变。第一次排序之后，所有的订单按照下单时间从早到晚有序了。在第二次排序中，我们用的是稳定的排序算法，所以经过第二次排序之后，相同金额的订单仍然保持下单时间从早到晚有序。

## 按下单时间有序

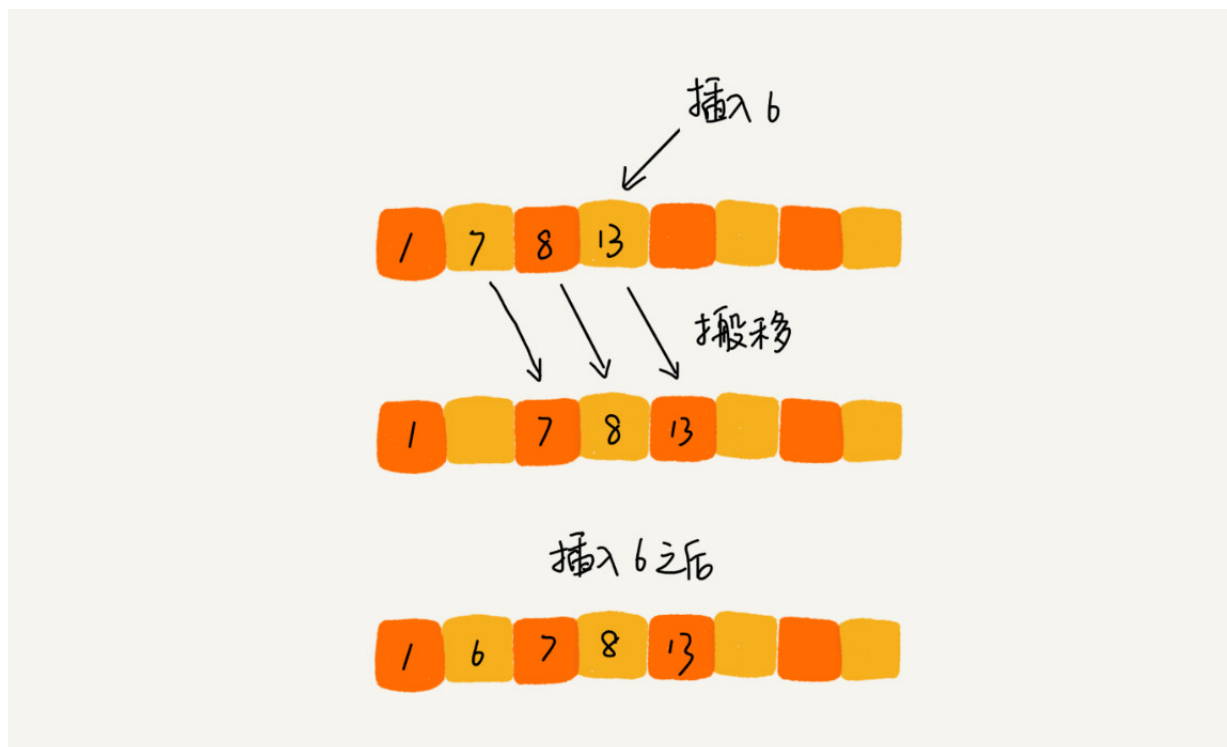
ID	下单时间	金额
1	2018-9-3 15:06:07	50
2	2018-9-3 16:08:10	30
3	2018-9-3 18:01:33	40
4	2018-9-3 20:23:31	30
5	2018-9-3 22:15:13	30
6	2018-9-4 05:07:33	60

## 按金额重新排序

ID	下单时间	金额
1	2018-9-3 16:08:10	30
2	2018-9-3 20:23:31	30
3	2018-9-3 22:15:13	30
4	2018-9-3 18:01:33	40
5	2018-9-3 15:06:07	50
6	2018-9-4 05:07:33	60

## 插入排序

先来看一个问题。一个有序的数组，我们往里面添加一个新的数据后，如何继续保持数据有序呢？很简单，我们只要遍历数组，找到数据应该插入的位置将其插入即可。

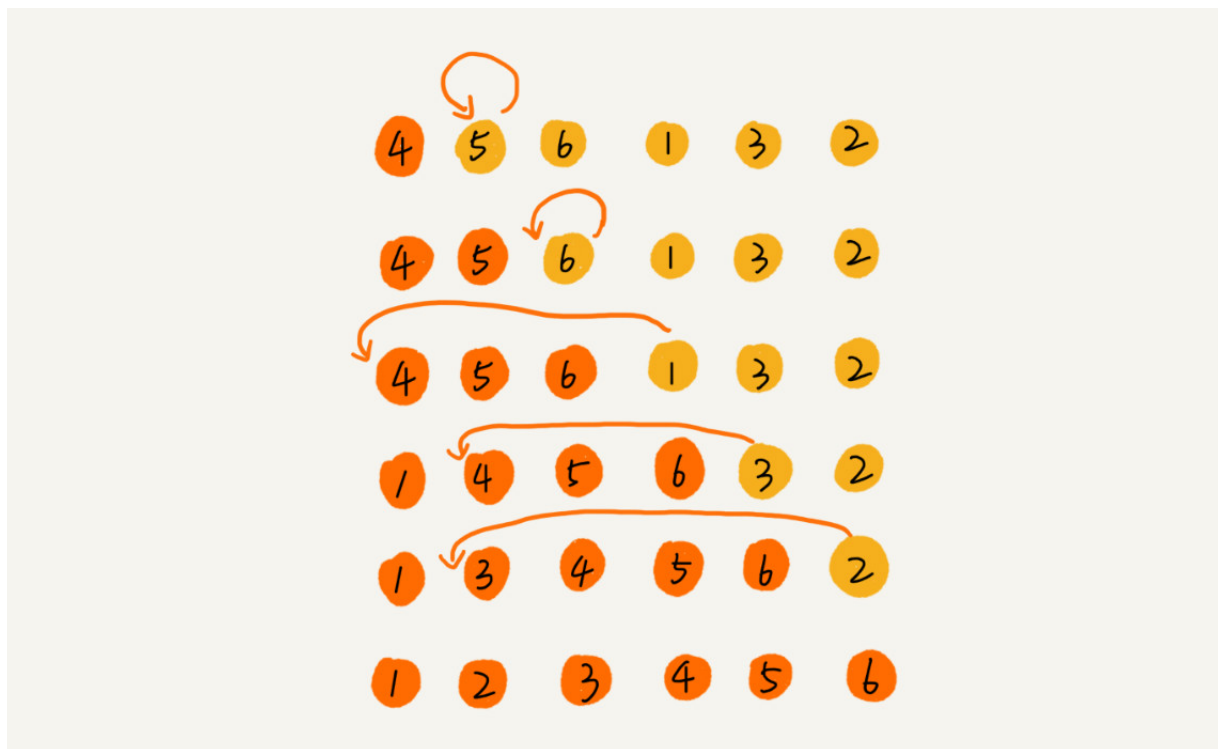


这是一个动态排序的过程，即动态地往有序集合中添加数据，我们可以通过这种方法保持集合中的数据一直有序。而对于一组静态数据，我们也可以借鉴上面讲的插入方法，来进行排序，于是就有了插入排序算法。

那插入排序具体是如何借助上面的思想来实现排序的呢？

首先，我们将数组中的数据分为两个区间，已排序区间和未排序区间。初始已排序区间只有一个元素，就是数组的第一个元素。插入算法的核心思想是取未排序区间中的元素，在已排序区间中找到合适的插入位置将其插入，并保证已排序区间数据一直有序。重复这个过程，直到未排序区间中元素为空，算法结束。

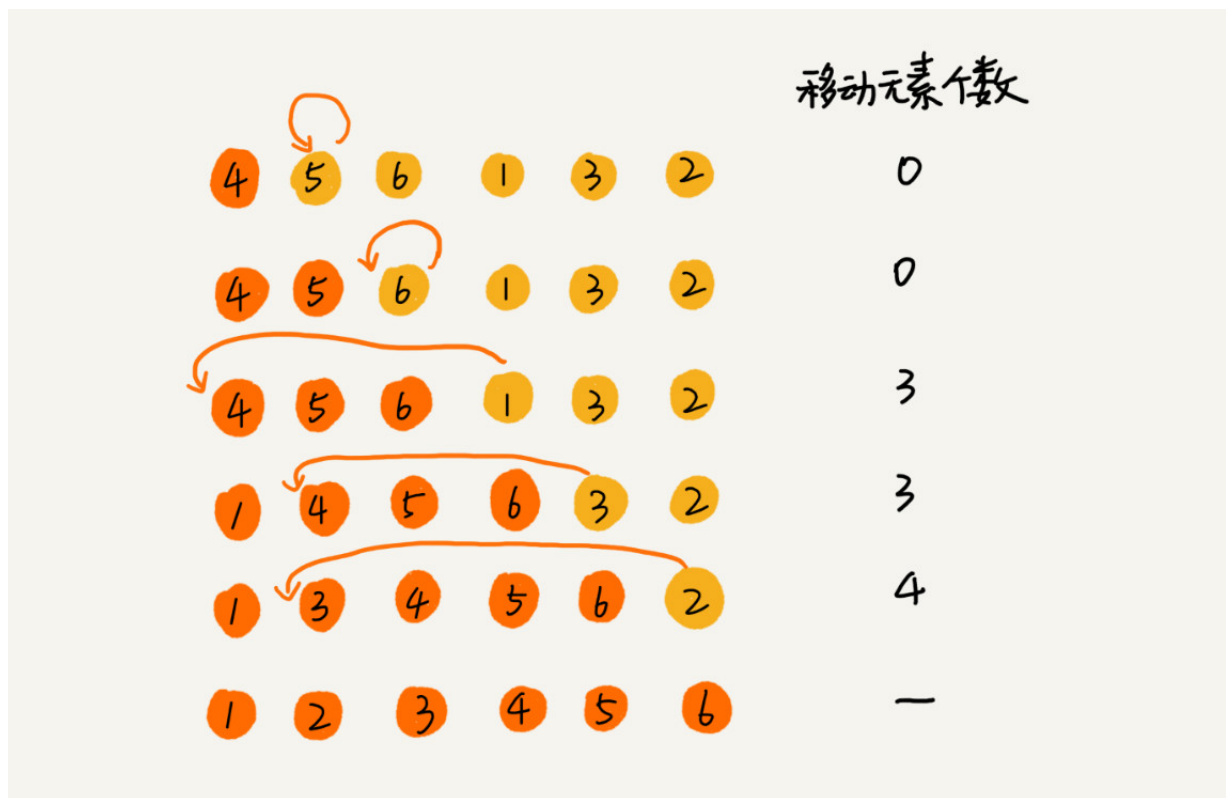
如图所示，要排序的数据是 4, 5, 6, 1, 3, 2，其中左侧为已排序区间，右侧是未排序区间。



插入排序也包含两种操作，一种是**元素的比较**，一种是**元素的移动**。当我们需要将一个数据 a 插入到已排序区间时，需要拿 a 与已排序区间的元素依次比较大小，找到合适的插入位置。找到插入点之后，我们还需要将插入点之后的元素顺序往后移动一位，这样才能腾出位置给元素 a 插入。

对于不同的查找插入点方法（从头到尾、从尾到头），元素的比较次数是有区别的。但对于一个给定的初始序列，移动操作的次数总是固定的，就等于**逆序度**。

为什么说移动次数就等于逆序度呢？我拿刚才的例子画了一个图表，你一看就明白了。满有序度是  $n \cdot (n-1)/2 = 15$ ，初始序列的有序度是 5，所以逆序度是 10。插入排序中，数据移动的个数总和也等于  $10 = 3 + 3 + 4$ 。



代码：

```
// 插入排序，a表示数组，n表示数组大小
public void insertionSort(int[] a, int n) {
    if (n <= 1) return;

    for (int i = 1; i < n; ++i) {
        int value = a[i];
        int j = i - 1;
        // 查找插入的位置
        for (; j >= 0; --j) {
            if (a[j] > value) {
                a[j+1] = a[j]; // 数据移动
            } else {
                break;
            }
        }
        a[j+1] = value; // 插入数据
    }
}
```

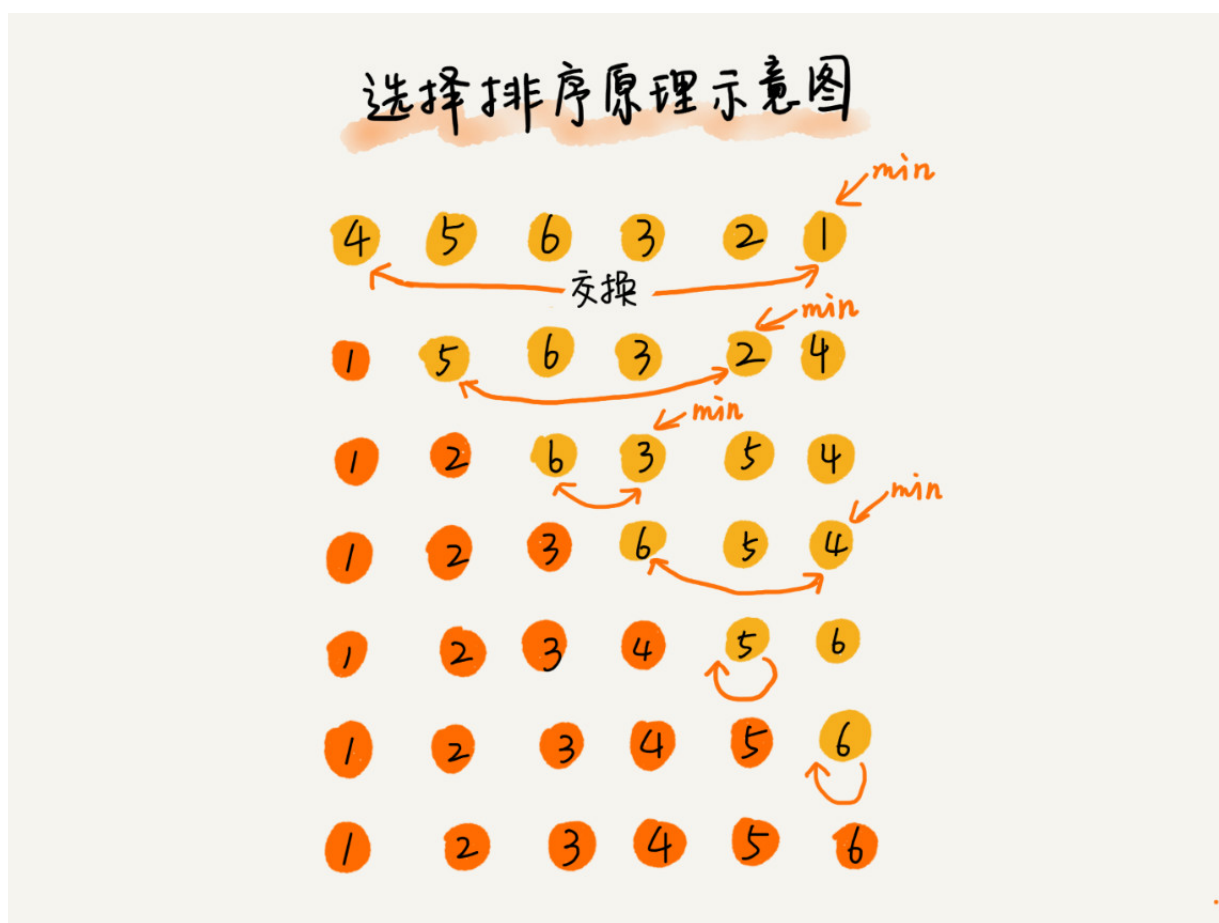
- 插入排序是原地排序算法吗？从实现过程可以很明显地看出，插入排序算法的运行并不需要额外的存储空间，所以空间复杂度是  $O(1)$ ，也就是说，这是一个原地排序算法。
- 插入排序是稳定的排序算法吗？在插入排序中，对于值相同的元素，我们可以选择将后面出现的元

素，插入到前面出现元素的后面，这样就可以保持原有的前后顺序不变，所以插入排序是**稳定的排序算法**。

- 插入排序的时间复杂度是多少？如果要排序的数据已经是有序的，我们并不需要搬移任何数据。如果我们从尾到头在有序数据组里面查找插入位置，每次只需要比较一个数据就能确定插入的位置。所以这种情况下，最好是时间复杂度为  $O(n)$ 。注意，这里是从尾到头遍历已经有序的数据。如果数组是倒序的，每次插入都相当于在数组的第一个位置插入新的数据，所以需要移动大量的数据，所以最坏和平均时间复杂度为  $O(n^2)$ 。

## 选择排序

选择排序算法的实现思路有点类似插入排序，也分已排序区间和未排序区间。但是选择排序每次会从未排序区间中找到最小的元素，将其放到已排序区间的末尾。

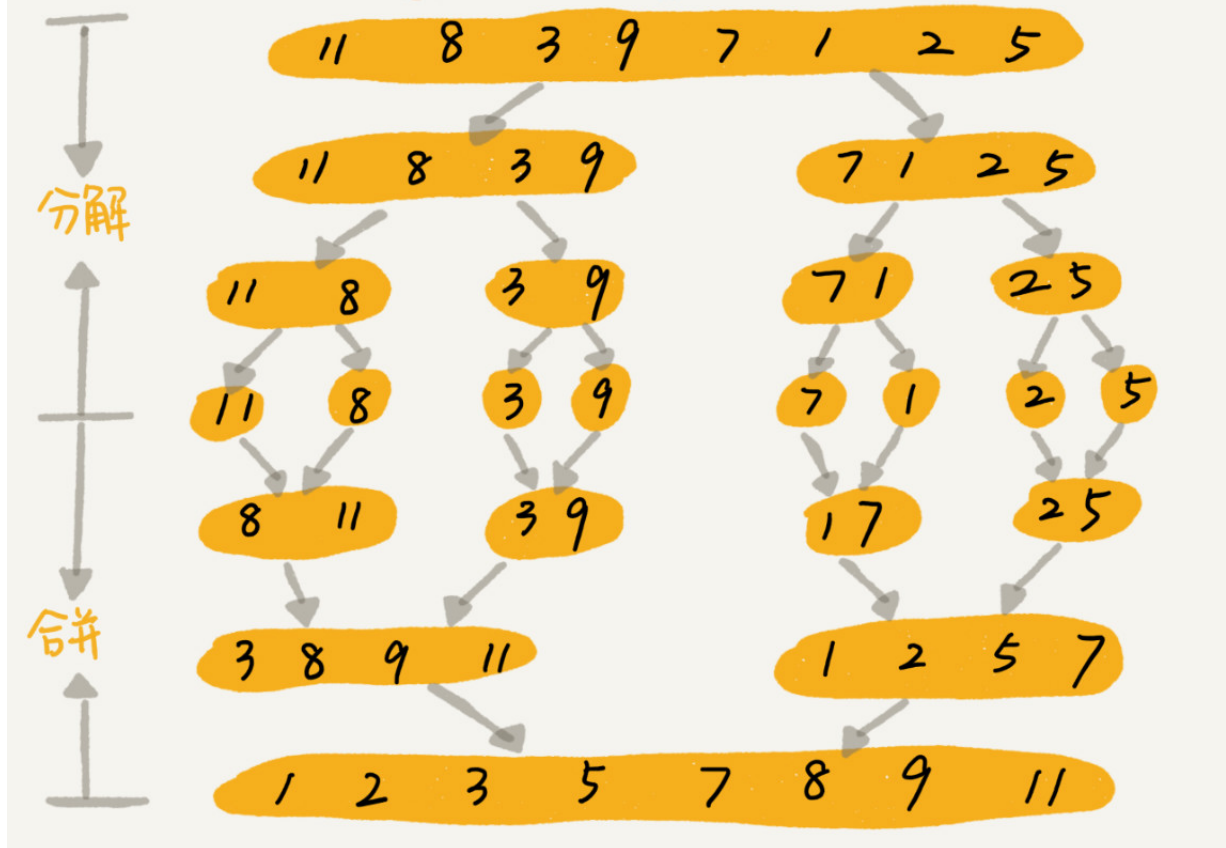


选择排序空间复杂度为  $O(1)$ ，是一种原地排序算法。选择排序的最好情况时间复杂度、最坏情况和平均情况时间复杂度都为  $O(n^2)$ 。选择排序是一种不稳定的排序算法（比如 5，8，5，2，9，第一次交换后，两个 5 的相对位置改变），所以稍逊于插入排序和冒泡排序。

## 归并排序

归并排序的核心思想很简单的。如果要排序一个数组，我们先把数组从中间分成前后两部分，然后对前后两部分分别排序，再将排好序的两部分合并在一起，这样整个数组就都有序了。

## 归并排序分解图



递推公式:

```
merge_sort(p...r) = merge(merge_sort(p...q), merge_sort(q+1...r))
```

终止条件:

`p >= r` 不用再继续分解

// 归并排序算法, A是数组, n表示数组大小

```
merge_sort(A, n) {
    merge_sort_c(A, 0, n-1)
}
```

// 递归调用函数

```
merge_sort_c(A, p, r) {
    // 递归终止条件
    if p >= r then return
```

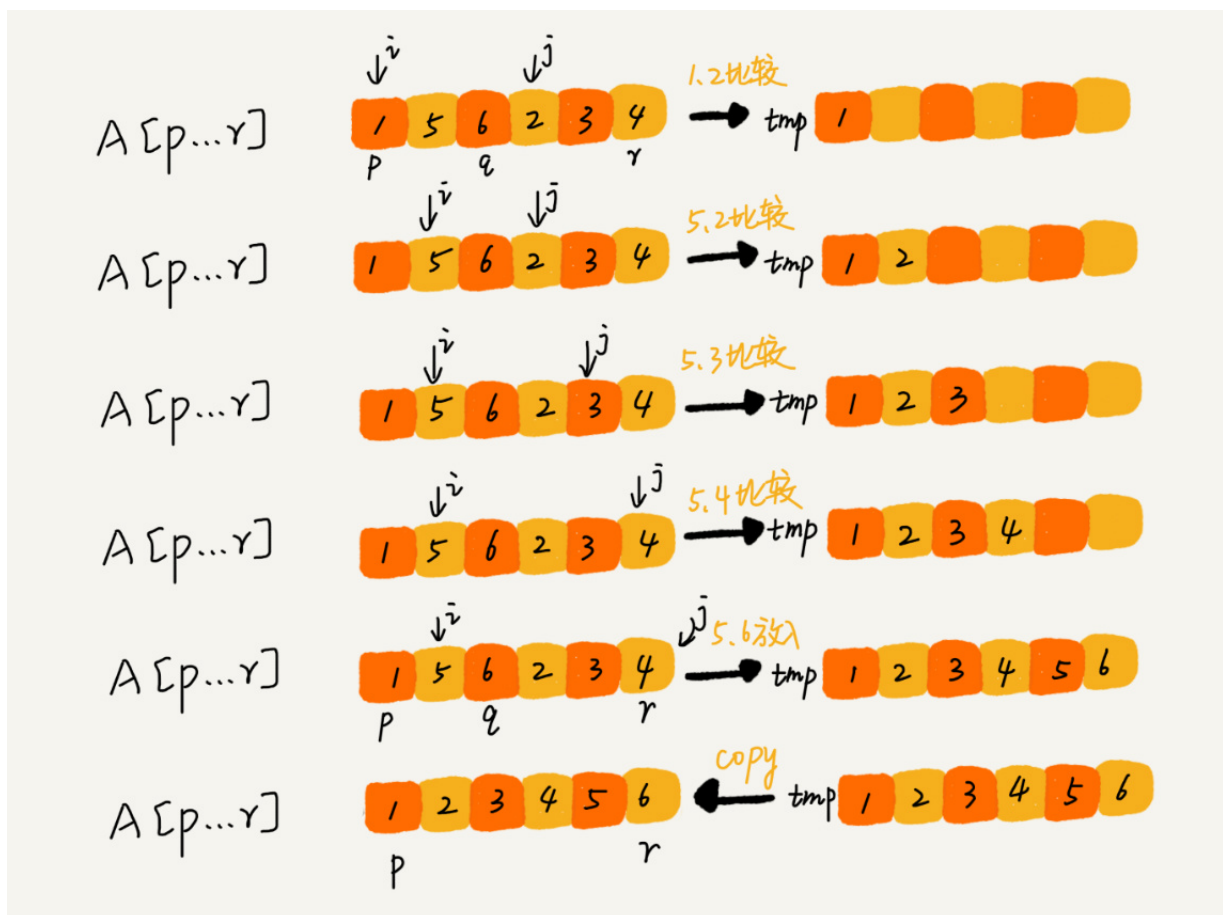
// 取p到r之间的中间位置q

```
q = (p+r) / 2
```



```
// 分治递归
merge_sort_c(A, p, q)
merge_sort_c(A, q+1, r)
// 将A[p...q]和A[q+1...r]合并为A[p...r]
merge(A[p...r], A[p...q], A[q+1...r])
}
```

`merge(A[p...r], A[p...q], A[q+1...r])` 这个函数的作用就是，将已经有序的 `A[p...q]` 和 `A[q+1...r]` 合并成一个有序的数组，并且放入 `A[p...r]`。那这个过程具体该如何做呢？如图所示，我们申请一个临时数组 `tmp`，大小与 `A[p...r]` 相同。我们用两个游标 `i` 和 `j`，分别指向 `A[p...q]` 和 `A[q+1...r]` 的第一个元素。比较这两个元素 `A[i]` 和 `A[j]`，如果 `A[i] ≤ A[j]`，我们就把 `A[i]` 放入到临时数组 `tmp`，并且 `i` 后移一位，否则将 `A[j]` 放入到数组 `tmp`，`j` 后移一位。继续上述比较过程，直到其中一个子数组中的所有数据都放入临时数组中，再把另一个数组中的数据依次加入到临时数组的末尾，这个时候，临时数组中存储的就是两个子数组合并之后的结果了。最后再把临时数组 `tmp` 中的数据拷贝到原数组 `A[p...r]` 中。



伪代码：

```
merge(A[p...r], A[p...q], A[q+1...r]) {
    var i := p, j := q+1, k := 0 // 初始化变量i, j, k
    var tmp := new array[0...r-p] // 申请一个大小跟A[p...r]一样的临时数组
    while i ≤ q AND j ≤ r do {
        if A[i] ≤ A[j] {
```

```

    tmp[k++] = A[i++] // i++等于i:=i+1
  } else {
    tmp[k++] = A[j++]
  }
}

// 判断哪个子数组中有剩余的数据
var start := i, end := q
if j<=r then start := j, end:=r

// 将剩余的数据拷贝到临时数组tmp
while start <= end do {
  tmp[k++] = A[start++]
}

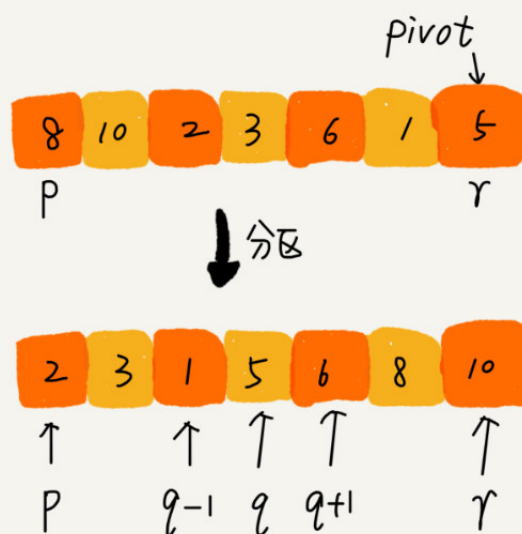
// 将tmp中的数组拷贝回A[p...r]
for i:=0 to r-p do {
  A[p+i] = tmp[i]
}
}

```

## 快速排序

快速排序算法（Quicksort），我们习惯性把它简称为“快排”。快排利用的也是分治思想。乍看起来，它有点像归并排序，但是思路其实完全不一样。待会会讲两者的区别。

快排的思想是这样的：如果要排序数组中下标从  $p$  到  $r$  之间的一组数据，我们选择  $p$  到  $r$  之间的任意一个数据作为 pivot（分区点）。我们遍历  $p$  到  $r$  之间的数据，将小于 pivot 的放到左边，将大于 pivot 的放到右边，将 pivot 放到中间。经过这一步骤之后，数组  $p$  到  $r$  之间的数据就被分成了三个部分，前面  $p$  到  $q-1$  之间都是小于 pivot 的，中间是 pivot，后面的  $q+1$  到  $r$  之间是大于 pivot 的。





根据分治、递归的处理思想，我们可以用递归排序下标从  $p$  到  $q-1$  之间的数据和下标从  $q+1$  到  $r$  之间的数据，直到区间缩小为 1，就说明所有的数据都有序了。如果我们用递推公式来将上面的过程写出来的话，就是这样：

递推公式：

```
quick_sort(p...r) = quick_sort(p...q-1) + quick_sort(q+1... r)
```

终止条件：

```
p >= r
```

```
// 快速排序，a是数组，n表示数组的大小
```

```
quick_sort(A, n) {
```

```
    quick_sort_c(A, 0, n-1)
```

```
}
```

```
// 快速排序递归函数，p,r为下标
```

```
quick_sort_c(A, p, r) {
```

```
    if p >= r then return
```

```
    q = partition(A, p, r) // 获取分区点
```

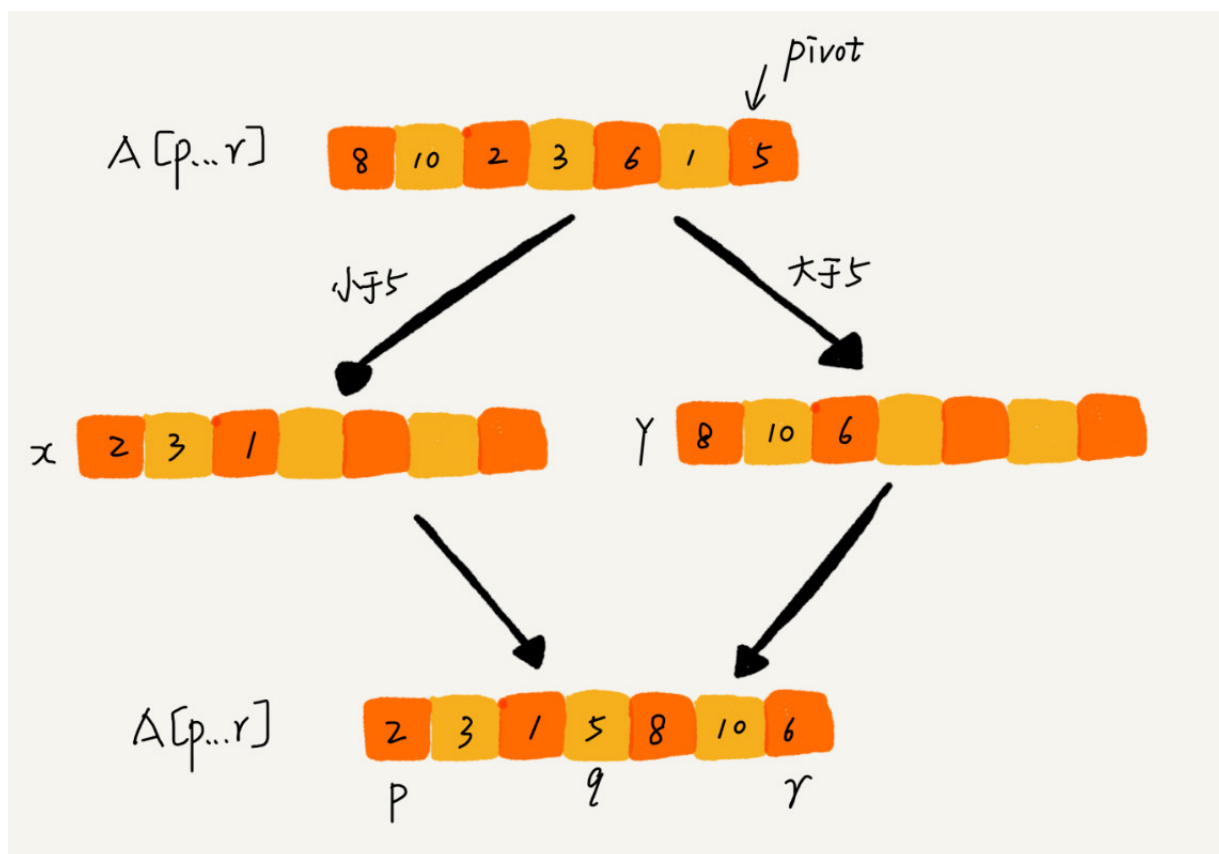
```
    quick_sort_c(A, p, q-1)
```

```
    quick_sort_c(A, q+1, r)
```

```
}
```

归并排序中有一个 `merge()` 合并函数，我们这里有一个 `partition()` 分区函数。`partition()` 分区函数实际上我们前面已经讲过了，就是随机选择一个元素作为 `pivot`（一般情况下，可以选择  $p$  到  $r$  区间的最后一个元素），然后对  $A[p...r]$  分区，函数返回 `pivot` 的下标。

如果我们不考虑空间消耗的话，`partition()` 分区函数可以写得非常简单。我们申请两个临时数组  $X$  和  $Y$ ，遍历  $A[p...r]$ ，将小于 `pivot` 的元素都拷贝到临时数组  $X$ ，将大于 `pivot` 的元素都拷贝到临时数组  $Y$ ，最后再将数组  $X$  和数组  $Y$  中数据顺序拷贝到  $A[p...r]$ 。

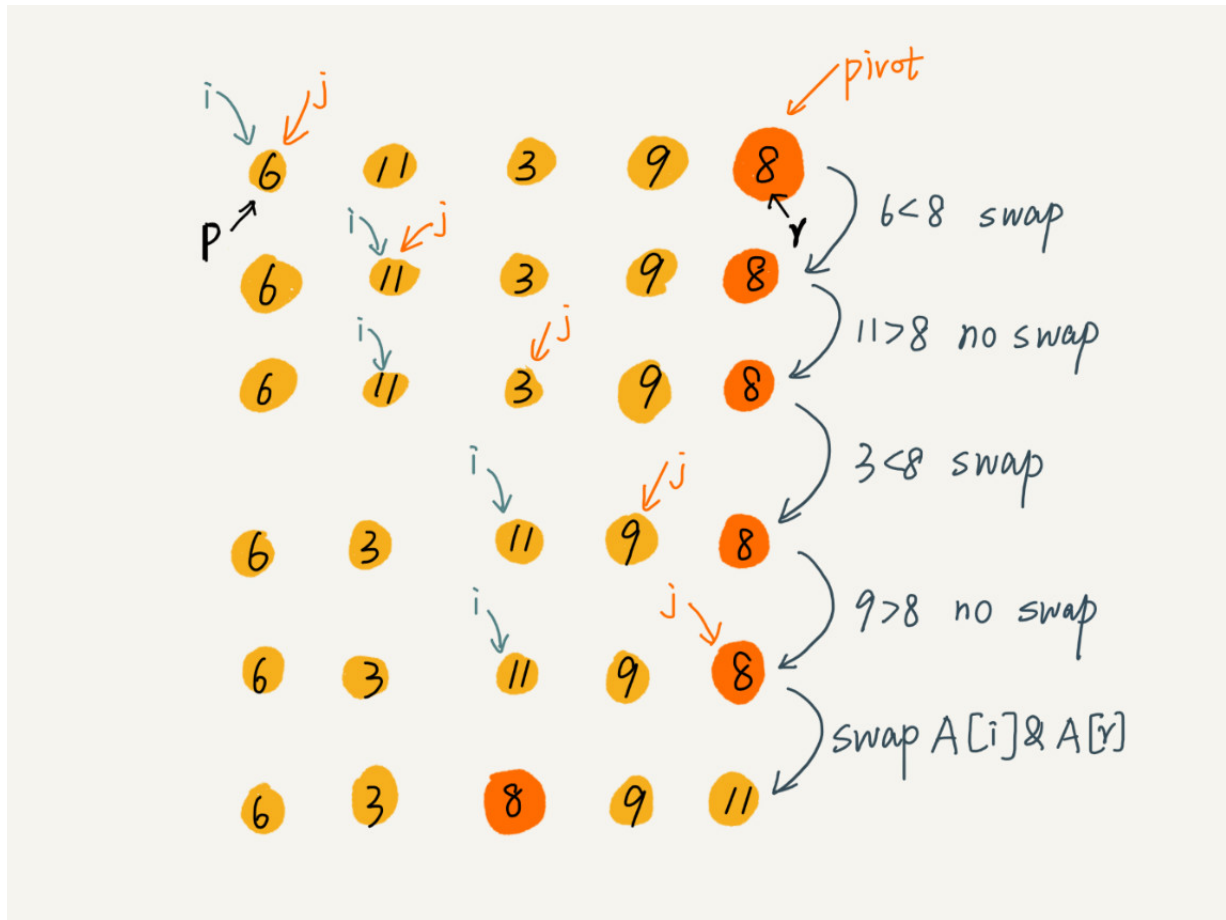


但是，如果按照这种思路实现的话，`partition()` 函数就需要很多额外的内存空间，所以快排就不是原地排序算法了。如果我们希望快排是原地排序算法，那它的空间复杂度得是  $O(1)$ ，那 `partition()` 分区函数就不能占用太多额外的内存空间，我们就需要在  $A[p \dots r]$  的原地完成分区操作。原地分区函数的实现思路非常巧妙，我写成了伪代码，我们一起来看一下。

```
partition(A, p, r) {  
    pivot := A[r]  
    i := p  
    for j := p to r-1 do {  
        if A[j] < pivot {  
            swap A[i] with A[j]  
            i := i+1  
        }  
    }  
    swap A[i] with A[r]  
    return i  
}
```

这里的处理有点类似选择排序。我们通过游标  $i$  把  $A[p \dots r-1]$  分成两部分。 $A[p \dots i-1]$  的元素都是小于 `pivot` 的，我们暂且叫它“已处理区间”， $A[i \dots r-1]$  是“未处理区间”。我们每次都从未处理的区间  $A[i \dots r-1]$  中取一个元素  $A[j]$ ，与 `pivot` 对比，如果小于 `pivot`，则将其加入到已处理区间的尾部，也就是  $A[i]$  的位置。

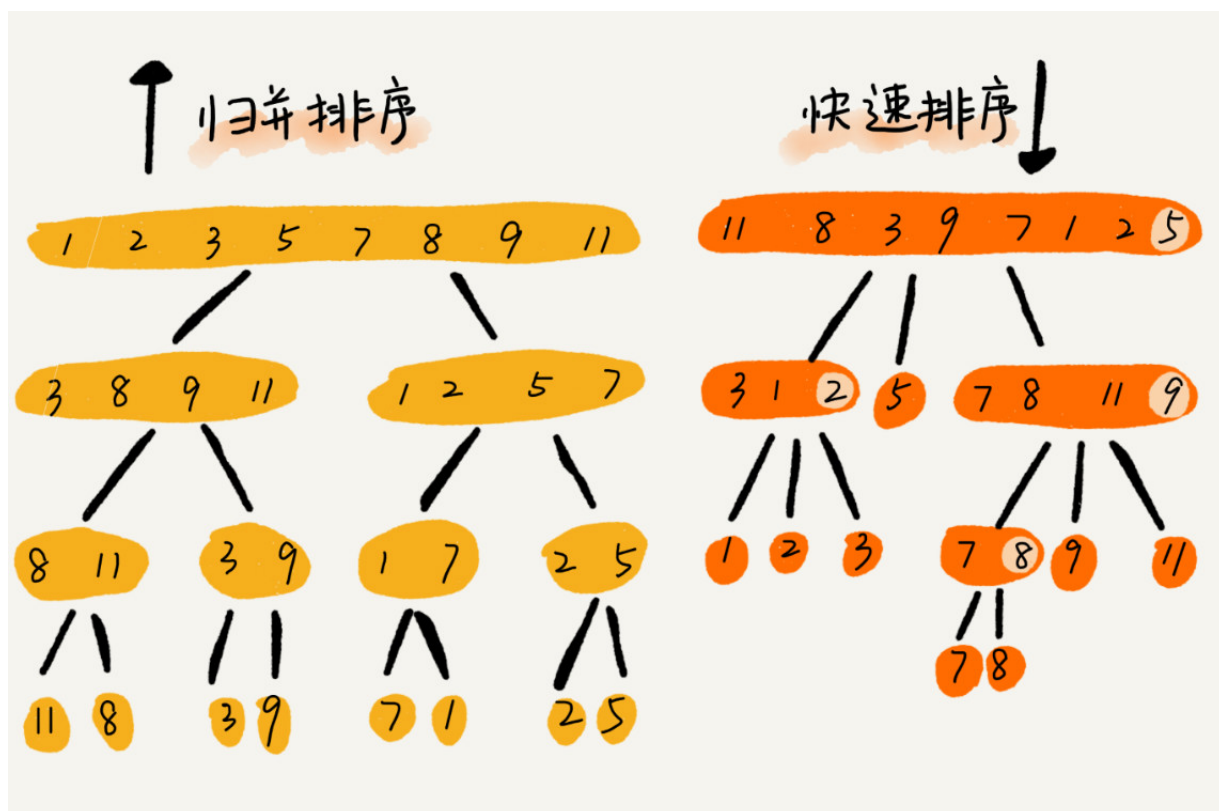
数组的插入操作还记得吗？在数组某个位置插入元素，需要搬移数据，非常耗时。当时我们也讲了一种处理技巧，就是交换，在  $O(1)$  的时间复杂度内完成插入操作。这里我们也借助这个思想，只需要将  $A[i]$  与  $A[j]$  交换，就可以在  $O(1)$  时间复杂度内将  $A[j]$  放到下标为  $i$  的位置。



因为涉及到交换，所以快排不是稳定的排序。

如果数组中有两个相同的元素，比如序列 6, 8, 7, 6, 3, 5, 9, 4，在经过第一次分区操作之后，两个 6 的相对先后顺序就会改变。

快排和归并用的都是分治思想，递推公式和递归代码也非常相似，那它们的区别在哪里呢？



归并排序：处理过程是由下到上的，先处理子问题，然后再合并。

快排：处理过程是由上到下的，先分区，然后再处理子问题。

归并排序：稳定的、时间复杂度为  $O(n\log n)$  的排序算法，但是它是非原地排序算法。

快速排序：不稳定的、时间复杂度为  $O(n\log n)$  的原地排序算法，解决了归并排序占用太多内存的问题。

## 最坏情况

快排的前提是每次分区操作，我们选择的 pivot 都很合适，正好能将大区间对等地一分为二。但实际上这种情况是很难实现的。如果数组中的数据原来已经是有序的了，比如 1, 3, 5, 6, 8。如果我们每次选择最后一个元素作为 pivot，那每次分区得到的两个区间都是不均等的。我们需要进行大约  $n$  次分区操作，才能完成快排的整个过程。每次分区我们平均要扫描大约  $n/2$  个元素，这种情况下，快排的时间复杂度就从  $O(n\log n)$  退化成了  $O(n^2)$ 。

假设每次分区操作都将区间分成大小为 9:1 的两个小区间。我们继续套用递归时间复杂度的递推公式，就会变成这样：

$T(1) = c$ ；  $n=1$  时，只需要常量级的执行时间，所以表示为  $c$ 。

$T(n) = T(n/10) + T(9*n/10) + n$ ；  $n > 1$

推导过程较复杂，最后的时间复杂度也是  $O(n\log n)$ 。

结论： $T(n)$  在大部分情况下的时间复杂度都可以做到  $O(n\log n)$ ，只有在极端情况下，才会退化到  $O(n^2)$ 。

## 思考

1. 根据伪代码写出快速排序的算法。
2. 如何利用快排思想在  $O(n)$  时间内查找第  $K$  大元素? ([http://ybt.ssoier.cn:8088/problem\\_show.php?pid=1167](http://ybt.ssoier.cn:8088/problem_show.php?pid=1167))
3. (选做) 1179 奖学金 [http://ybt.ssoier.cn:8088/problem\\_show.php?pid=1179](http://ybt.ssoier.cn:8088/problem_show.php?pid=1179)