

0.前言

堆（Heap）与栈（Stack）是开发人员必须面对的两个概念，在理解这两个概念时，需要放到具体的场景下，因为不同场景下，堆与栈代表不同的含义。一般情况下，有两层含义：

- (1) 程序 **内存** 布局场景下，堆与栈表示两种内存管理方式；
- (2) **数据结构** 场景下，堆与栈表示两种常用的数据结构。

1.程序内存分区中的堆与栈

1.1 栈简介

栈由操作系统自动分配释放，用于存放函数的参数值、局部变量等，其操作方式类似于数据结构中的栈。参考如下代码：

```
1 | int main() {  
2 |     int b;                                // 栈  
3 |     char s[] = "abc";                    // 栈  
4 |     char *p2;                            // 栈  
5 | }
```

其中函数中定义的局部变量按照先后定义的顺序依次压入栈中，也就是说相邻变量的地址之间不会存在其它变量。栈的内存地址生长方向与堆相反，由高到底，所以后定义的变量地址低于先定义的变量，比如上面代码中变量 s 的地址小于变量 b 的地址，p2 地址小于 s 的地址。栈中存储的数据的生命周期随着函数的执行完成而结束。

1.2 堆简介

堆由开发人员分配和释放，若开发人员不释放，程序结束时由 OS 回收，分配方式类似于链表。参考如下代码：

```
1 |  
2 | int main() {  
3 |     // C 中用 malloc() 函数申请  
4 |     char* p1 = (char *)malloc(10);  
5 |     cout<<(int*)p1<<endl;           // 输出: 00000000003BA0C0  
6 |  
7 |     // 用 free() 函数释放  
8 |     free(p1);  
9 |  
10 |    // C++ 中用 new 运算符申请  
11 |    char* p2 = new char[10];  
12 |    cout << (int*)p2 << endl;       // 输出: 00000000003BA0C0  
13 | }
```

```
14 | // 用 delete 运算符释放
15 | delete[] p2;
    | }
```

其中 p1 所指的 10 字节的内存空间与 p2 所指的 10 字节内存空间都是存在于堆。堆的内存地址生长方向与栈相反，由低到高，但需要注意的是，后申请的内存空间并不一定在先申请的内存空间的后面，即 p2 指向的地址并不一定大于 p1 所指向的内存地址，原因是先申请的内存空间一旦被释放，后申请的内存空间则会利用先前被释放的内存，从而导致先后分配的内存空间在地址上不存在先后关系。堆中存储的数据若未释放，则其生命周期等同于程序的生命周期。

关于堆上内存空间的分配过程，首先应该知道操作系统有一个记录空闲内存地址的链表，当系统收到程序的申请时，会遍历该链表，寻找第一个空间大于所申请空间的堆节点，然后将该节点从空闲节点链表中删除，并将该节点的空间分配给程序。另外，对于大多数系统，会在这块内存空间中的首地址处记录本次分配的大小，这样，代码中的delete语句才能正确地释放本内存空间。由于找到的堆节点的大小不一定正好等于申请的大小，系统会自动地将多余的那部分重新放入空闲链表。

1.3 堆与栈区别

堆与栈实际上是操作系统对进程占用的内存空间的两种管理方式，主要有如下几种区别：

(1) 管理方式不同。栈由操作系统自动分配释放，无需我们手动控制；堆的申请和释放工作由程序员控制，容易产生内存泄漏；

(2) 空间大小不同。每个进程拥有的栈大小要远远小于堆大小。理论上，进程可申请的堆大小为虚拟内存大小，进程栈的大小 64bits 的 Windows 默认 1MB，64bits 的 Linux 默认 10MB；

(3) 生长方向不同。堆的生长方向向上，内存地址由低到高；栈的生长方向向下，内存地址由高到低。

(4) 分配方式不同。堆都是动态分配的，没有静态分配的堆。栈有 2 种分配方式：静态分配和动态分配。静态分配是由操作系统完成的，比如局部变量的分配。动态分配由 `alloca()` 函数分配，但是栈的动态分配和堆是不同的，它的动态分配是由操作系统进行释放，无需我们手工实现。

(5) 分配效率不同。栈由操作系统自动分配，会在硬件层级对栈提供支持：分配专门的寄存器存放栈的地址，压栈出栈都有专门的指令执行，这就决定了栈的效率比较高。堆则是由C/C++提供的库函数或运算符来完成申请与管理，实现机制较为复杂，频繁的内存申请容易产生内存碎片。显然，堆的效率比栈要低得多。

(6) 存放内容不同。栈存放的内容，函数返回地址、相关参数、局部变量和寄存器内容等。当主函数调用另外一个函数的时候，要对当前函数执行断点进行保存，需要使用栈来实现，首先入栈的是主函数下一条语句的地址，即扩展指针寄存器的内容（EIP），然后是当前栈帧的底部地址，即扩展基址指针寄存器内容（EBP），再然后是被调函数的实参等，一般情况下是按照从右向左的顺序入栈，之后是被调函数的局部变量，注意静态变量是存放在数据段或者BSS段，是不入栈的。出栈的

顺序正好相反，最终栈顶指向主函数下一条语句的地址，主程序又从该地址开始执行。堆，一般情况堆顶使用一个字节的空间来存放堆的大小，而堆中具体存放内容是由程序员来填充的。

从以上可以看到，堆和栈相比，由于大量malloc()/free()或new/delete的使用，容易造成大量的内存碎片，并且可能引发用户态和核心态的切换，效率较低。栈相比于堆，在程序中应用较为广泛，最常见的是函数的调用过程由栈来实现，函数返回地址、EBP、实参和局部变量都采用栈的方式存放。虽然栈有众多的好处，但是由于和堆相比不是那么灵活，有时候分配大量的内存空间，主要还是用堆。

无论是堆还是栈，在内存使用时都要防止非法越界，越界导致的非法内存访问可能会摧毁程序的堆、栈数据，轻则导致程序运行处于不确定状态，获取不到预期结果，重则导致程序异常崩溃，这些都是我们编程时与内存打交道时应该注意的问题。

2.数据结构中的堆与栈

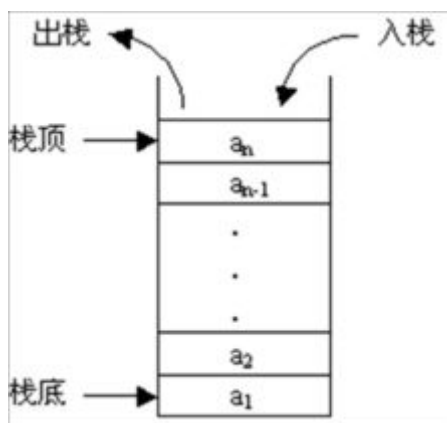
数据结构中，堆与栈是两个常见的数据结构，理解二者的定义、用法与区别，能够利用堆与栈解决很多实际问题。

2.1 栈简介

栈是一种运算受限的线性表，其限制是指只仅允许在表的一端进行插入和删除操作，这一端被称为栈顶（Top），相对地，把另一端称为栈底（Bottom）。把新元素放到栈顶元素的上面，使之成为新的栈顶元素称作进栈、入栈或压栈（Push）；把栈顶元素删除，使其相邻的元素成为新的栈顶元素称作出栈或退栈（Pop）。这种受限的运算使栈拥有“先进后出”的特性（First In Last Out），简称FILO。

栈分顺序栈和链式栈两种。栈是一种线性结构，所以可以使用数组或链表（单向链表、双向链表或循环链表）作为底层数据结构。使用数组实现的栈叫做顺序栈，使用链表实现的栈叫做链式栈，二者的区别是顺序栈中的元素地址连续，链式栈中的元素地址不连续。

栈的结构如下图所示：



栈的基本操作包括初始化、判断栈是否为空、入栈、出栈以及获取栈顶元素等。下面以顺序栈为

用 C++ 给出一个简单的实现。

```

2  #include<stdio.h>
3  #include<malloc.h>
4
5  #define DataType int
6  #define MAXSIZE 1024
7  struct SeqStack {
8      DataType data[MAXSIZE];
9      int top;
10 };
11 // 栈初始化, 成功返回栈对象指针, 失败返回空指针NULL
12 SeqStack* initSeqStack() {
13     SeqStack* s=(SeqStack*)malloc(sizeof(SeqStack));
14     if(!s) {
15         printf("空间不足\n");
16         return NULL;
17     } else {
18         s->top = -1;
19         return s;
20     }
21 }
22
23 // 判断栈是否为空
24 bool isEmptySeqStack(SeqStack* s) {
25     if (s->top == -1)
26         return true;
27     else
28         return false;
29 }
30
31 // 入栈, 返回-1失败, 0成功
32 int pushSeqStack(SeqStack* s, DataType x) {
33     if(s->top == MAXSIZE-1)
34     {
35         return -1; // 栈满不能入栈
36     } else {
37         s->top++;
38         s->data[s->top] = x;
39         return 0;
40     }
41 }
42
43 // 出栈, 返回-1失败, 0成功
44 int popSeqStack(SeqStack* s, DataType* x) {
45     if(isEmptySeqStack(s)) {
46         return -1; // 栈空不能出栈
47     } else {
48         *x = s->data[s->top];
49         s->top--;

```

```

49         return 0;
50     }
51 }
52
53 //取栈顶元素, 返回-1失败, 0成功
54 int topSeqStack(SeqStack* s,DataType* x) {
55     if (isEmptySeqStack(s))
56         return -1;        //栈空
57     else {
58         *x=s->data[s->top];
59         return 0;
60     }
61 }
62
63 //打印栈中元素
64 int printSeqStack(SeqStack* s) {
65     int i;
66     printf("当前栈中的元素:\n");
67     for (i = s->top; i >= 0; i--)
68         printf("%4d",s->data[i]);
69     printf("\n");
70     return 0;
71 }
72
73 //test
74 int main() {
75     SeqStack* seqStack=initSeqStack();
76     if(seqStack) {
77         //将4、5、7分别入栈
78         pushSeqStack(seqStack,4);
79         pushSeqStack(seqStack,5);
80         pushSeqStack(seqStack,7);
81
82         //打印栈内所有元素
83         printSeqStack(seqStack);
84
85         //获取栈顶元素
86         DataType x=0;
87         int ret=topSeqStack(seqStack,&x);
88         if(0==ret) {
89             printf("top element is %d\n",x);
90         }
91
92         //将栈顶元素出栈
93         ret=popSeqStack(seqStack,&x);
94         if(0==ret) {
95             printf("pop top element is %d\n",x);
96         }
97     }
98 }

```

```

96 |         }
97 |         return 0;
98 |     }
99 |

```

运行上面的程序，输出结果：

```

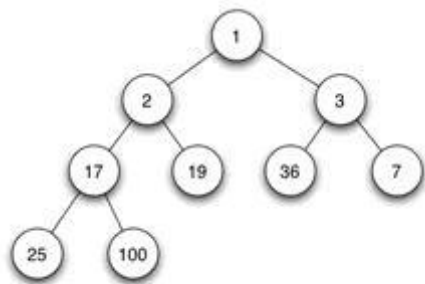
1 | 当前栈中的元素：
2 |     7   5   4
3 | top element is 7
4 | pop top element is 7

```

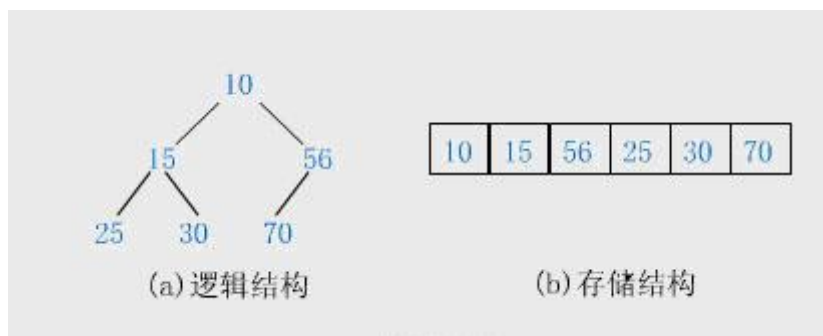
2.2 堆简介

2.2.1 堆的性质

堆是一种常用的树形结构，是一种特殊的完全二叉树，当且仅当满足所有节点的值总是不大于或不小于其父节点的值是完全二叉树被称之为堆。堆的这一特性称之为堆序性。因此，在一个堆中，根节点是最大（或最小）节点。如果根节点最小，称之为小顶堆（或小根堆），如果根节点最大，称之为大顶堆（或大根堆）。堆的左右孩子没有大小的顺序。下面是一个小顶堆示例：



堆的存储一般都用数组来存储堆， i 节点的父节点下标就为 $(i-1)/2$ 。它的左右子节点下标分别为 $2*i+1$ 和 $2*i+2$ 。如第0个节点左右子节点下标分别为1和2。



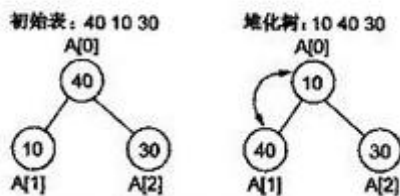
2.2.2 堆的基本操作

- 建立

以最小堆为例，如果以数组存储元素时，一个数组具有对应的树表示形式，但树并不满足堆的条件，需要重新排列元素，可以建立“堆化”的树。

初始表: 40 10 30

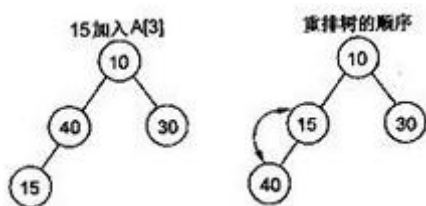
堆化树: 10 40 30



- 插入

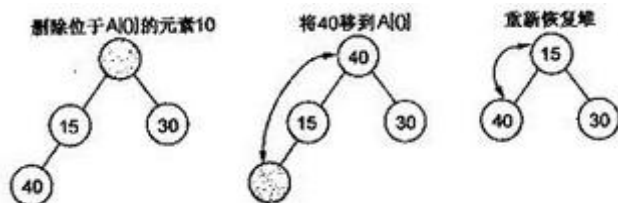
将一个新元素插入到数组末尾，如果新构成的二叉树不满足堆的性质，需要将新元素在其到堆顶的路径上，找到属于自己的位置，即进行上浮操作。

下图演示了插入 15 时，堆的调整。



- 删除

删除一个元素总是发生在堆顶，因为堆顶的元素是最小的（小顶堆中）。数组中最后一个元素用来填补空缺位置，然后对顶部元素进行下沉，如果左右孩子有比自己小的，则选择选择最小的那个进行交换。重复进行下沉操作，以满足堆的条件。



2.2.3 堆操作实现

(1) 插入实现

每次插入都是将新数据放在数组最后。可以发现从这个新数据的父节点到根节点必然为一个有序的数列，现在的任务是将这个新数据插入到这个有序数据中，这就类似于直接插入排序中将一个数据并入到有序区间中，这是节点“上浮”调整。不难写出插入一个新数据时堆的调整代码：

```
1 // 新加入 i 节点, 其父节点为 (i-1)/2
2 // 参数: a: 数组, i: 新插入元素在数组中的下标
3 void minHeapFixUp(int a[], int i) {
4     int j, temp;
5     temp = a[i];
6     j = (i-1)/2; // 父节点
7     while (j >= 0 && i != 0) {
```

```

8         if (a[j] <= temp)//如果父节点不大于新插入的元素, 停止寻找
9             break;
10        a[i]=a[j];           //把较大的子节点往下移动, 替换它的子节点
11        i = j;
12        j = (i-1)/2;
13    }
14    a[i] = temp;
15 }

```

因此, 插入数据到最小堆时:

```

1 // 在最小堆中加入新的数据data
2 // a: 数组, index: 插入的下标,
3 void minHeapAddNumber(int a[], int index, int data) {
4     a[index] = data;
5     minHeapFixUp(a, index);
6 }

```

(2) 删除实现

按照堆删除的说明, 堆中每次都只能删除第0个数据。为了便于重建堆, 实际的操作是将数组最后一个数据与根节点交换, 然后再从根节点开始进行一次从上向下的调整。

调整时先在左右儿子节点中找最小的, 如果父节点不大于这个最小的子节点说明不需要调整了, 反之将最小的子节点换到父节点的位置。此时父节点实际上并不需要换到最小子节点的位置, 因为这不是父节点的最终位置。但逻辑上父节点替换了最小的子节点, 然后再考虑父节点对后面的节点的影响。堆元素的删除导致的堆调整, 其整个过程就是将根节点进行“下沉”处理。下面给出代码:

```

1
2 // minHeapFixDown 小顶堆结点下沉操作。
3 // a 为数组, len 为结点总数; 从 index 结点开始调整, index 从 0 开始计算 index 其子:
4 void minHeapFixDown(int a[],int len,int index) {
5     // index 为叶子节点不用调整。
6     if(index>(len/2-1)) return;
7     int tmp=a[index];
8     lastIndex=index;
9
10    // 当下沉到叶子节点时, 就不用调整了。
11    while(index<=len/2-1) {
12        // 如果左子节点小于待调整节点
13        if(a[2*index+1]<tmp) {
14            lastIndex = 2*index+1;
15        }
16        // 如果存在右子节点且小于左子节点和待调整节点
17        if(2*index+2<len && a[2*index+2]<a[2*index+1]&& a[2*index+

```



```

18         lastIndex=2*index+2;
19     }
20     //如果左右子节点有一个小于待调整节点，选择最小子节点进行上浮
21     if(lastIndex!=index) {
22         a[index]=a[lastIndex];
23         index=lastIndex;
24     } else break;           // 否则待调整节点不用下沉调整
25 }
26 // 将待调整节点放到最后的位置。
27 a[lastIndex]=tmp;
}

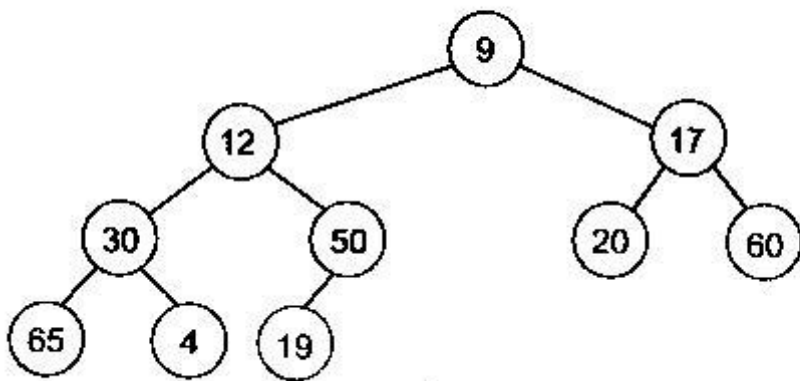
```

根据堆删除的下沉思想，可以有不同版本的代码实现，以上是和孙凛同学一起讨论出的一个版本，在这里感谢他的参与，读者可另行给出。个人体会，这里建议大家根据对堆调整过程的理解，写出自己的代码，切勿看示例代码去理解算法，而是理解算法思想写出代码，否则很快就会忘记。

(3) 建堆

有了堆的插入和删除后，再考虑下如何对一个数据进行堆化操作。要一个一个的从数组中取出数据来建立堆吧，不用！先看一个数组，如下图：

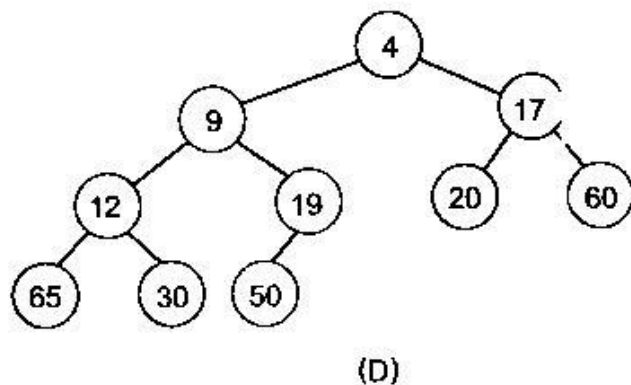
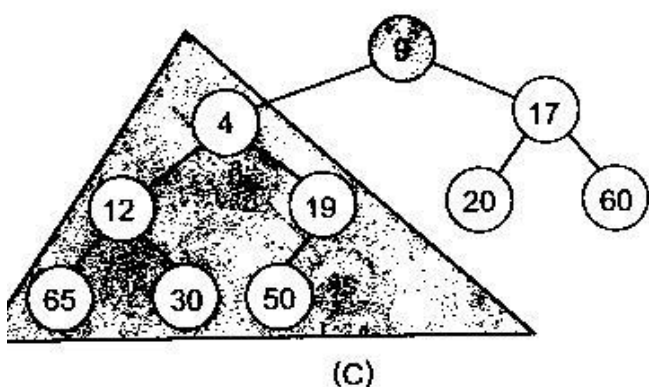
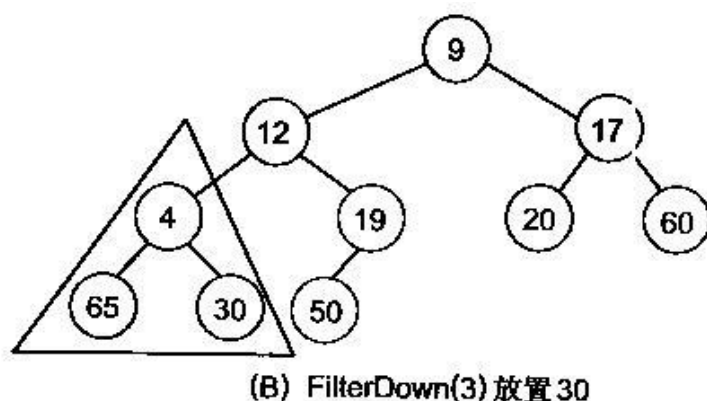
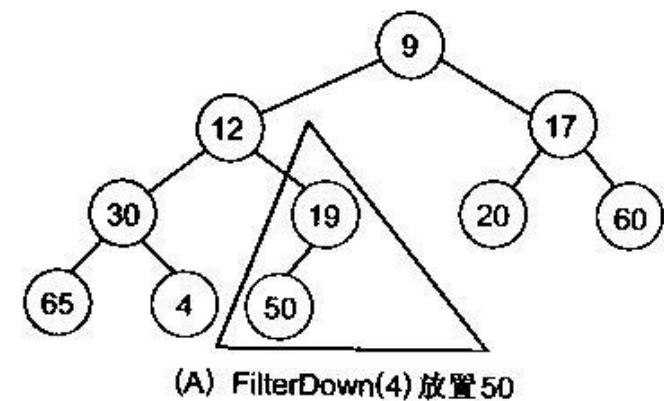
```
int A[0] = {9,12,17,30,50,20,60,65,4,49};
```



初始表

很明显，对叶子结点来说，可以认为它已经是一个合法的堆了，即 20, 60, 65, 4, 49 都分别是一个合法的堆。只要从 $A[4]=50$ 开始向下调整就可以了。然后再取 $A[3]=30$, $A[2]=17$, $A[1]=12$, $A[0]=9$ 分别作一次向下调整操作就可以了。

下图展示了这些步骤：



写出堆化数组的代码：

```
1 // makeMinHeap 建立最小堆。
2 // a: 数组, n: 数组长度。
3 void makeMinHeap(int a[], int n) {
4     for (int i = n/2-1; i >= 0; i--) {
5         minHeapFixDown(a, i, n);
6     }
7 }
```

2.2.4 堆的具体应用——堆排序

堆排序 (Heapsort) 是堆的一个经典应用，有了上面对堆的了解，不难实现堆排序。由于堆也是用数组来存储的，故对数组进行堆化后，第一次将 $A[0]$ 与 $A[n-1]$ 交换，再对 $A[0 \dots n-2]$ 重新恢复堆。第二次将 $A[0]$ 与 $A[n-2]$ 交换，再对 $A[0 \dots n-3]$ 重新恢复堆，重复这样的操作直到 $A[0]$ 与 $A[1]$ 交换。由于每次都是将最小的数据并入到后面的有序区间，故操作完成后整个数组就有序了。有点类似于直接选择排序。

因此，完成堆排序并没有用到前面说明的插入操作，只用到了建堆和节点向下调整的操作，堆排序的操作如下：

```
1 // array: 待排序数组, len: 数组长度
2 void heapSort(int array[], int len) {
3     // 建堆
4     makeMinHeap(array, len);
5 }
```

```

6
7      // 最后一个叶子节点和根节点交换，并进行堆调整，交换次数为len-1次
8      for(int i=len-1;i>0;--i) {
9          //最后一个叶子节点交换
10         array[i]=array[i]+array[0];
11         array[0]=array[i]-array[0];
12         array[i]=array[i]-array[0];
13
14         // 堆调整
15         minHeapFixDown(array, 0, len-i-1);
16     }
}

```

(1) 稳定性。堆排序是不稳定排序。

(2) 堆排序性能分析。由于每次重新恢复堆的时间复杂度为 $O(\log N)$ ，共 $N-1$ 次堆调整操作，再加上前面建立堆时 $N/2$ 次向下调整，每次调整时间复杂度也为 $O(\log N)$ 。两次操作时间复杂度相加还是 $O(N\log N)$ ，故堆排序的时间复杂度为 $O(N\log N)$ 。

最坏情况：如果待排序数组是有序的，仍然需要 $O(N\log N)$ 复杂度的比较操作，只是少了移动的操作；

最好情况：如果待排序数组是逆序的，不仅需要 $O(N\log N)$ 复杂度的比较操作，而且需要 $O(N\log N)$ 复杂度的交换操作，总的时间复杂度还是 $O(N\log N)$ 。

因此，堆排序和快速排序在效率上是差不多的，但是堆排序一般优于快速排序的重要一点是数据的初始分布情况对堆排序的效率没有大的影响。

参考文献

- [1] 浅谈堆和栈的区别
- [2] 栈内存和堆内存的区别
- [3] 浅谈内存分配方式以及堆和栈的区别（很清楚）
- [4] C++函数调用过程深入分析
- [5] 十种排序算法

杂注

我的博客即将搬运同步至腾讯云+社区，邀请大家一同入驻：

https://cloud.tencent.com/developer/support-plan?invite_code=2z2o0f9ecoo44